

1 **A SMOOTHING ACTIVE SET METHOD FOR LINEARLY**
2 **CONSTRAINED NON-LIPSCHITZ NONCONVEX OPTIMIZATION ***

3 CHAO ZHANG[†] AND XIAOJUN CHEN[‡]

4 **Abstract.** We propose a novel smoothing active set method for linearly constrained non-
5 Lipschitz nonconvex problems. At each step of the proposed method, we approximate the objective
6 function by a smooth function with a fixed smoothing parameter and employ a new active set method
7 for minimizing the smooth function over the original feasible set, until a special updating rule for the
8 smoothing parameter meets. The updating rule is always satisfied within finite number of iterations
9 since the new active set method for smooth problems proposed in this paper forces at least one sub-
10 sequence of projected gradients to zero. Any accumulation point of the smoothing active set method
11 is a stationary point associated with the smoothing function used in the method, which is necessary
12 for local optimality of the original problem. And any accumulation point for the $\ell_2 - \ell_p$ ($0 < p < 1$)
13 sparse optimization model is a limiting stationary point, which is a local minimizer under a certain
14 second-order condition. Numerical experiments demonstrate the efficiency and effectiveness of our
15 smoothing active set method for hyperspectral unmixing on 3D image cube of large size.

16 **Key words.** Non-Lipschitz, nonconvex, linearly constrained, smoothing active set method,
17 stationary point

18 **AMS subject classifications.** 65K10 90C26 90C46

19 **1. Introduction.** Active set methods have been successfully used for linearly
20 constrained smooth optimization problems of large size; see [8, 13, 17, 18, 25, 42]
21 and references therein. Hager and Zhang developed a novel active set algorithm for
22 the bound constrained smooth optimization problems in [17], and ten years later they
23 extended the method to solve linearly constrained smooth optimization problems [18].
24 The active set method in [18] switches between phase one that employs the gradient
25 projection algorithm for the original problem and phase two that uses an algorithm
26 with certain requirements for solving linearly constrained optimization problems on a
27 face of the original feasible set. Hager and Zhang [18] showed that any accumulation
28 point of the sequence generated by their method is a stationary point, and only phase
29 two is performed after a finite number of iterations under certain conditions.

30 For linearly constrained nonsmooth convex optimization problems, Panier pro-
31 posed an active set method [29], in which the search direction is computed by a
32 bundle principle. And the convergence result is obtained under a certain nondegener-
33 acy assumption. Wen et al. developed an active set algorithm for the unconstrained
34 ℓ_1 minimization with good numerical performance and convergence results [36, 37].
35 For bound-constrained nonsmooth nonconvex optimization, Keskar and Wächter pro-
36 posed a limited-memory quasi-Newton algorithm which uses an active set selection
37 strategy to define the subspace in which search directions are computed [21]. Numerical
38 experiments were conducted to show the efficacy of the algorithm, but theoretical
39 convergence guarantees are elusive even for the unconstrained case. To the best of
40 our knowledge, there is no active set method that tackles linearly constrained non-
41 Lipschitz nonconvex optimization problems with solid convergence results.

*Submitted to the editors April 1, 2019.

Funding: C. Zhang's work was supported in part by NSFC grants 11571033, 11431002. X. Chen's work was supported in part by Hong Kong Research Council Grant PolyU15300/17P.

[†]Department of Applied Mathematics, Beijing Jiaotong University, Beijing 100044, China (zc.njtu@163.com).

[‡]Department of Applied Mathematics, The Hong Kong Polytechnic University, Hung Hum, Kowloon, Hong Kong (xiaojun.chen@polyu.edu.hk).

42 One effective way to overcome the nonsmoothness in optimization is the type
 43 of smoothing methods, which uses the structure of the problem to define smoothing
 44 functions and the algorithms for solving smooth problems. Nesterov proposed a s-
 45 moothing scheme [27] for minimizing a nonsmooth convex function over a convex set.
 46 Zhang and Chen proposed a smoothing projected gradient method [41] for minimiz-
 47 ing a Lipschitz continuous function over a convex set. Bian and Chen developed a
 48 smoothing quadratic regularization method [4] for a class of linearly constrained non-
 49 Lipschitz optimization problems arising from image restoration. Xu et al. proposed
 50 a smoothing sequential quadratic programming method [38] for solving degenerate
 51 nonsmooth and nonconvex constrained optimization problems with applications to
 52 bilevel programs. Liu et al. proposed a smoothing sequential quadratic programming
 53 framework [26] for a class of composite ℓ_p ($0 < p < 1$) minimization over polyhedron.

54 Inspired by the active set method [18] and the smoothing technique, we develop
 55 a novel smoothing active set method with solid convergence results for the following
 56 minimization problem

$$57 \quad (1.1) \quad \min f(x) \quad \text{s.t.} \quad x \in \Omega,$$

58 where $f : R^n \rightarrow R$ is continuous but not necessarily Lipschitz continuous and

$$59 \quad (1.2) \quad \Omega = \{x \in R^n : c_i^T x = d_i, i \in \mathcal{M}_E; c_i^T x \leq d_i, i \in \mathcal{M}_I\}.$$

60 Here $\mathcal{M}_E = \{1, 2, \dots, m_e\}$, $\mathcal{M}_I = \{m_e + 1, m_e + 2, \dots, m\}$, $\mathcal{M} = \mathcal{M}_E \cup \mathcal{M}_I$, and
 61 $c_i \in R^n$, $d_i \in R$ for $i = 1, 2, \dots, m$.

62 Problem (1.1) involving a sparsity penalized term in the objective function has
 63 recently intrigued a lot of interests. It serves as a basic model for a variety of im-
 64 portant applications, including the compressed sensing [1], the edge-preserving image
 65 restoration [4, 28], the sparse nonnegative matrix factorization for data classification
 66 [40], and the sparse portfolio selection [9, 15]. For example, the widely used $\ell_2 - \ell_p$
 67 ($0 < p < 1$) sparse optimization model

$$68 \quad (1.3) \quad \min \|Ax - b\|^2 + \tau \|x\|_p^p \quad \text{s.t.} \quad x \geq 0,$$

69 where $\|\cdot\|$ refers to the Euclidean norm, $\|x\|_p^p = \sum_{i=1}^n |x_i|^p$, and $A \in R^{l \times n}$, $b \in R^l$,
 70 and $\tau > 0$ are given. The non-Lipschitz nonconvex term $\|x\|_p^p$ in the objective function
 71 and the nonnegative constraints benefit to recover some prior knowledge such as the
 72 sparsity of the signal, or the range of pixels. It is worth mentioning that in typical
 73 compressive sensing or image restoration, the dimension of optimization problems is
 74 large.

75 In order to develop the smoothing active set method, we first assume f is smooth
 76 in (1.1) in section 2 and develop an efficient new active set method for the linearly
 77 constrained smooth problems, which can be considered as a modification of the active
 78 set algorithm [18]. The new active set method combines the projected gradient (PG)
 79 method [8] and a linearly constrained optimizer (LCO) that satisfies mild require-
 80 ments. We show in Theorem 2.2 that the new active set method forces at least one
 81 subsequence of projected gradients to zero. This property is essential in developing
 82 the smoothing active set method with global convergence in section 3. It is guaran-
 83 teed that any accumulation point of the sequence generated by the new active set
 84 method is a stationary point. Moreover, if the sequence generated by the new active
 85 set method converges to a stationary point x^* , then the sequence can identify the set
 86 of strongly active constraints and hence is trapped by the face exposed by $-\nabla f(x^*)$

87 after a finite number of iterations. The convergence and identification properties are
 88 not guaranteed by the active set method in [18] for the smooth problems. Based on
 89 the identification properties, we also prove the local convergence result that if the
 90 sequence converges to x^* and the strong second-order sufficient optimality condition
 91 holds, then only the LCO is executed after a finite number of iterations.

92 Combining the new active set method for linearly constrained smooth minimiza-
 93 tion problem with delicate smoothing strategies, we then develop in section 3 a novel
 94 smoothing active set method that solves the linearly constrained non-Lipschitz min-
 95 imization problem (1.1). The new active set method for smooth problems is used to
 96 solve the smoothing problems. We give the concept of a stationary point associated
 97 with the smoothing function and show that it is necessary for optimality of the original
 98 problem. We show that any accumulation point generated by the smoothing active
 99 set method is a stationary point of the original problem. Moreover, it is a limiting
 100 stationary point of problem (1.3). If in addition a second-order condition holds, it is
 101 also a strict local minimizer of (1.3).

102 We conduct numerical experiments on real applications of large scale in hyper-
 103 spectral unmixing in section 4. The numerical results manifest that the smoothing
 104 active set method performs favorably in comparison to [several state-of-the-art meth-](#)
 105 [ods](#) in hyperspectral unmixing.

106 Throughout the paper, we use the following notation. $\langle x, y \rangle = x^T y$ presents the
 107 inner product of two vectors x and y of the same dimension. $R_+^n = \{x \in R^n : x \geq 0\}$
 108 and $R_{++}^n = \{x \in R^n : x > 0\}$. $|\mathcal{S}|$ corresponds to the cardinality of a finite set \mathcal{S} . If
 109 \mathcal{S} is a subset of $\{1, 2, \dots, n\}$, then for any vector $u \in R^n$ and $M \in R^{n \times n}$, $u_{\mathcal{S}}$ is the
 110 subvector of u whose entries lie in u indexed by \mathcal{S} , and $M_{\mathcal{S}\mathcal{S}}$ denotes the submatrix
 111 of M whose rows and columns lie in \mathcal{S} . $\mathcal{N}(M)$ is the null space of M . Let \mathbb{N} be the
 112 set of all natural numbers and $\mathcal{N}_{\infty}^{\#}$ be the infinite subsets of \mathbb{N} . We use the notation
 113 \xrightarrow{N} for the convergence indexed by $N \in \mathcal{N}_{\infty}^{\#}$. The normal cone to a closed convex set
 114 Ω at x is denoted by $N_{\Omega}(x)$, and $P_{\Omega}[x] = \operatorname{argmin}\{\|z - x\| : z \in \Omega\}$ is the orthogonal
 115 projection from x into Ω . The ball with center x^* and radius δ is denoted by $B_{\delta}(x^*)$.
 116 For any $x \in R^n$, the active and free index sets are defined by

$$117 \quad \mathcal{A}(x) := \mathcal{M}_E \cup \{i \in \mathcal{M}_I : c_i^T x = d_i\}, \quad \mathcal{F}(x) := \{i \in \mathcal{M}_I : c_i^T x < d_i\}.$$

118 **2. A new active set method for linearly constrained smooth minimiza-**
 119 **tion.** In this section, we consider the following linearly constrained smooth problem

$$120 \quad (2.1) \quad \min f(x) \quad \text{s.t.} \quad x \in \Omega,$$

121 where f is continuously differentiable and Ω is defined in (1.2).

122 Recall that the projected gradient $\nabla_{\Omega} f(x)$ is defined by

$$123 \quad \nabla_{\Omega} f(x) \equiv P_{T(x)}[-\nabla f(x)] = \operatorname{argmin}\{\|v + \nabla f(x)\| : v \in T(x)\},$$

124 where $T(x)$ is the tangent cone to Ω at x . Calamai and Moré (Lemma 3.1 of [8])
 125 showed that $x^* \in \Omega$ is a stationary point of (2.1) if and only if $\nabla_{\Omega} f(x^*) = 0$. It is
 126 worth mentioning that $\|\nabla_{\Omega} f(x)\|$ can be bounded away from zero in a neighborhood of
 127 a stationary point x^* , since $\|\nabla_{\Omega} f(\cdot)\|$ is not continuous, but only lower semicontinuous
 128 on Ω according to Lemma 3.3 of [8]. That is, for any $\{x^k\} \subset \Omega$ converging to x ,

$$129 \quad \|\nabla_{\Omega} f(x)\| \leq \liminf_{k \rightarrow \infty} \|\nabla_{\Omega} f(x^k)\|.$$

130 A stationary point x^* of (2.1) is often characterized as

$$131 \quad d^1(x^*) := P_{\Omega}[x^* - \nabla f(x^*)] - x^* = 0.$$

We find that convergence of most existing active set methods for (2.1) is to show $\liminf_{k \rightarrow \infty} \|d^1(x^k)\| = 0$, such as the active set method in [18]. However, since the norm of projected gradient is not continuous, $\liminf_{k \rightarrow \infty} \|d^1(x^k)\| = 0$ does not imply $\liminf_{k \rightarrow \infty} \|\nabla_{\Omega} f(x^k)\| = 0$. See Example 1 in section 2. The new active set method proposed in this section aims to have

$$\liminf_{k \rightarrow \infty} \|\nabla_{\Omega} f(x^k)\| = 0,$$

132 which is essential for showing the convergence result of the smoothing active set
133 method for solving nonsmooth problem (1.1) proposed in section 3.

134 **2.1. Structure of the new active set method.** Now we introduce the neces-
135 sary notation used in the new active set method. Let us denote $g(x) = \nabla f(x)$. Given
136 an index set \mathcal{S} satisfying $\mathcal{M}_E \subseteq \mathcal{S} \subseteq \mathcal{M}$, we define $g^{\mathcal{S}}(x) \in R^n$ by

$$137 \quad (2.2) \quad g^{\mathcal{S}}(x) = P_{N(C_{\mathcal{S}}^T)}[g(x)] = \arg \min \{\|y - g(x)\| : y \in R^n \text{ and } C_{\mathcal{S}}^T y = 0\},$$

138 where $C_{\mathcal{S}} \in R^{n \times |\mathcal{S}|}$ is the matrix whose columns are $c_i, i \in \mathcal{S}$. In particular, we denote
139 $g^{\mathcal{A}}(x)$ for $\mathcal{S} = \mathcal{A}(x)$ and if $\mathcal{A}(x) = \emptyset$, then $g^{\mathcal{A}}(x) = g(x)$. Thus $g^{\mathcal{A}}(x)$ is the unique
140 optimal solution of the strongly convex problem

$$141 \quad (2.3) \quad \min \quad \frac{1}{2} \|y - g(x)\|^2 \quad \text{s.t.} \quad c_i^T y = 0, \quad i \in \mathcal{A}(x).$$

142 From the first-order optimality conditions, it is easy to find that for $x \in \Omega$, $g^{\mathcal{A}}(x) = 0$
143 if and only if x is a stationary point of f on its associated face

$$144 \quad (2.4) \quad \check{\Omega}(x) := \{y \in \Omega : c_i^T y = d_i \text{ for all } i \in \mathcal{A}(x)\}.$$

145 Let x^* be a stationary point of (2.1) and $\Lambda(x^*)$ be the set of Lagrange multipliers
146 associated with the constraints. That is, $x^* \in \Omega$ and for any $\lambda^* \in \Lambda(x^*)$, (x^*, λ^*)
147 satisfies

$$148 \quad (2.5) \quad \begin{aligned} & g(x^*) + \sum_{i \in \mathcal{M}} \lambda_i^* c_i = 0, \\ & \lambda_i^* \geq 0 \text{ if } i \in \mathcal{M}_I \cap \mathcal{A}(x^*), \quad \lambda_i^* = 0 \text{ if } i \in \mathcal{F}(x^*), \\ & \lambda_i^* (c_i^T x^* - d_i) = 0 \text{ for all } i \in \mathcal{M}_I. \end{aligned}$$

149 Consider

$$150 \quad (2.6) \quad y(x, \alpha) = P_{\Omega}[x - \alpha g(x)] = \operatorname{argmin} \{\|x - \alpha g(x) - y\|^2 : y \in \Omega\},$$

151 where $\alpha > 0$ is a given number. Thus there exists $\lambda \in R^m$ such that $(y(x, \alpha), \lambda)$
152 satisfies

$$153 \quad (2.7) \quad \begin{aligned} & y(x, \alpha) - (x - \alpha g(x)) + \sum_{i \in \mathcal{M}} \lambda_i c_i = 0, \\ & \lambda_i \geq 0 \text{ if } i \in \mathcal{M}_I \cap \mathcal{A}(y(x, \alpha)), \quad \lambda_i = 0 \text{ if } i \in \mathcal{F}(y(x, \alpha)), \\ & \lambda_i (c_i^T y(x, \alpha) - d_i) = 0 \text{ for all } i \in \mathcal{M}_I. \end{aligned}$$

154 Let $\Lambda(x, \alpha)$ be the set of Lagrange multipliers satisfying (2.7) at the solution $y =$
155 $y(x, \alpha)$ of (2.6). It is easy to see that

$$156 \quad (2.8) \quad y(x^*, \alpha) = x^* \quad \text{and} \quad \Lambda(x^*, \alpha) = \alpha \Lambda(x^*).$$

157 In the new active set method, it employs either the iteration of the PG method
 158 or the iteration of the LCO by given rules. Let x^k be the current iterate and the LCO
 159 be chosen to get the new iterate. Then the LCO solves the problem

$$160 \quad (2.9) \quad \min f(y) \quad \text{s.t.} \quad y \in \check{\Omega}(x^k),$$

161 which operates on the faces of Ω . Compared to the original problem (2.1), there are
 162 usually much more equality constraints in (2.9) which may lead the efficiency of the
 163 LCO. This is obviously true when the feasible set is defined by the bound constraints
 164 or the simplex constraint (which are sometimes called ‘‘hard constraints’’ and it is
 165 better to satisfy them strictly rather than penalize them into the objective function).
 166 The PG step comes from the classic ‘‘piecewise PG method’’ proposed in [8], and an
 167 arbitrary LCO can be chosen as long as it satisfies certain requirements listed below.

- 168 • PG method

169 Given $\rho, \beta \in (0, 1)$. For $k = 1, 2, \dots$,
 170 set $d^k = -g(x^k)$ and let $x^{k+1} = P_{\Omega}[x^k + \alpha_k d^k]$ where α_k is determined by
 171 the Armijo line search, i.e., $\alpha_k = \max\{\rho^0, \rho^1, \dots\}$ is chosen such that

$$172 \quad (2.10) \quad f(x^{k+1}) \leq f(x^k) + \beta \langle g(x^k), x^{k+1} - x^k \rangle.$$

- 173 • LCO Requirements

174 For $k = 1, 2, \dots$,
 175 F1: $x^k \in \Omega$ and $f(x^{k+1}) \leq f(x^k)$ for each k .
 176 F2: $\mathcal{A}(x^k) \subseteq \mathcal{A}(x^{k+1})$ for each k .
 177 F3: If $\exists \bar{k} > 0$ such that $\mathcal{A}(x^j) \equiv \bar{\mathcal{A}}$ for all $j \geq \bar{k}$, then $\liminf_{j \rightarrow \infty} \|g^{\bar{\mathcal{A}}}(x^j)\| = 0$.

178 F1 and F2 of the LCO Requirements are satisfied, as long as the LCO adopts
 179 a monotone line search, and whenever a new constraint becomes active, it changes
 180 the corresponding inequality constraint to the equality constraint in (2.9). Later we
 181 always assume the two strategies are incorporated in the LCO. F3 requires that if the
 182 active set becomes stable as $\mathcal{A}(x^j) \equiv \bar{\mathcal{A}}$, then at least one accumulation point x^* of
 183 the sequence $\{x^k\}$ generated by the LCO is a stationary point of problem (2.9) with
 184 $\check{\Omega}(x^k) = \{y \in \Omega : c_i^T y = d_i \text{ for all } i \in \bar{\mathcal{A}}\}$. Note that in this case x^* is a stationary
 185 point if and only if $g^{\bar{\mathcal{A}}}(x^*) = 0$. And since $g^{\bar{\mathcal{A}}}(x) = P_{\mathcal{N}(C_{\bar{\mathcal{A}}}^T)}[g(x)]$, we know that
 186 $g^{\bar{\mathcal{A}}}(\cdot) : R^n \rightarrow R^n$ is a continuous function. Thus $g^{\bar{\mathcal{A}}}(x^*) = 0$ indicates

$$187 \quad \liminf_{j \rightarrow \infty} \|g^{\bar{\mathcal{A}}}(x^j)\| = \liminf_{j \rightarrow \infty} \|g^{\bar{\mathcal{A}}}(x^j)\| = 0.$$

188 Therefore the LCO Requirements can be easily fulfilled by many algorithms based
 189 on gradient or Newton type iterations that employ a monotone line search and add
 190 constraints to the active set whenever a new constraint becomes active, e.g., the pro-
 191 jected gradient method [8], the method of Zoutendijk (section 10.1 of [2]), the Frank-
 192 Wolfe algorithm [16], the first-order interior-point method [33], and the affine-scaling
 193 interior-point method [19]. When $\Omega = R_+^n$, we can employ the LCO using essentially
 194 unconstrained optimization methods such as the conjugate gradient method as in [17].

195 Now we are ready to outline the new active set method for problem (2.1).

196 2.2. Convergence analysis.

Assumption 2.1. For any $\Gamma \in R$, the level set

$$\mathcal{L}_{\Gamma} = \{x \in \Omega : f(x) \leq \Gamma\}$$

197 is bounded.

Algorithm 2.1 A new active set method

-
- 1: **Parameters:** $\epsilon \in [0, \infty)$, θ and $\eta \in (0, 1)$. $x^1 = P_\Omega[x^0]$, $k = 1$.
 - 2: **Phase one:**
 - 3: **while** $\|\nabla_\Omega f(x^k)\| > \epsilon$, **do**
 - 4: Execute the PG step to obtain x^{k+1} from x^k . Let $k \leftarrow k + 1$.
 - 5: If $\|g^A(x^k)\| \leq \theta \|\nabla_\Omega f(x^k)\|$, then $\theta \leftarrow \eta\theta$.
 - 6: If $\|g^A(x^k)\| > \theta \|\nabla_\Omega f(x^k)\|$, then go to phase two.
 - 7: **end while**
 - 8: **Phase two:**
 - 9: **while** $\|\nabla_\Omega f(x^k)\| > \epsilon$, **do**
 - 10: Execute the LCO step to obtain x^{k+1} from x^k . Let $k \leftarrow k + 1$.
 - 11: If $\|g^A(x^k)\| \leq \theta \|\nabla_\Omega f(x^k)\|$, then go to phase one and $\theta \leftarrow \eta\theta$.
 - 12: **end while**
-

198 In the remainder of this paper, we assume that the LCO satisfies the LCO Re-
 199 quirements F1-F3, and Assumption 2.1 holds. We now show the global convergence
 200 of Algorithm 2.1 for problem (2.1).

201 **THEOREM 2.2.** *Let $\{x^k\}$ be the sequence generated by Algorithm 2.1 with $\epsilon = 0$.*
 202 *Then there exists at least one accumulation point of $\{x^k\}$,*

$$203 \quad (2.11) \quad \liminf_{k \rightarrow \infty} \|\nabla_\Omega f(x^k)\| = 0,$$

204 and any accumulation point of $\{x^k\}$ is a stationary point of (2.1).

205 *Proof.* By Assumption 2.1, there exists at least one accumulation point x^* of
 206 $\{x^k\}$. Let $\{x^k\}_{k \in K}$ be an infinite subsequence of $\{x^k\}$ such that $\lim_{k \rightarrow \infty, k \in K} x^k = x^*$.

207 If only phase one is performed for k sufficiently large, then by Assumption 2.1
 208 and Theorem 2.4 of [8],

$$209 \quad \lim_{k \rightarrow \infty, k \in K} \frac{x^{k+1} - x^k}{\alpha_k} = 0.$$

210 Hence for $k \rightarrow \infty$, $k \in K$,

$$211 \quad \|x^{k+1} - x^*\| \leq \|x^{k+1} - x^k\| + \|x^k - x^*\| \rightarrow 0,$$

which indicates $\lim_{k \rightarrow \infty, k \in K} x^{k+1} = x^*$. According to Theorem 3.4 of [8],

$$\lim_{k \rightarrow \infty, k \in K} \|\nabla_\Omega f(x^{k+1})\| = 0.$$

212 By the lower semicontinuity of $\|\nabla_\Omega f(\cdot)\|$ shown in Lemma 3.3 of [8],

$$213 \quad \|\nabla_\Omega f(x^*)\| \leq \lim_{k \rightarrow \infty, k \in K} \|\nabla_\Omega f(x^{k+1})\| = 0,$$

214 which guarantees that x^* is a stationary point of (2.1).

215 If only phase two is performed for k sufficiently large, then there exists $\hat{\theta} > 0$ such
 216 that $\theta \equiv \hat{\theta}$ for k sufficiently large, because θ is never reduced in phase two. Hence for
 217 k sufficiently large,

$$218 \quad (2.12) \quad \|g^A(x^k)\| \geq \hat{\theta} \|\nabla_\Omega f(x^k)\|.$$

219 Note that the LCO works on the faces of Ω and no index in the active set can be
 220 freed from x^k to x^{k+1} using the LCO. By F2 of the LCO Requirements, the active set
 221 becomes stable for k large enough and hence $\liminf_{k \rightarrow \infty} \|g^A(x^k)\| = 0$ according to
 222 F3. From (2.12) we then have (2.11) holds. By the lower semicontinuity of $\|\nabla_{\Omega} f(\cdot)\|$,
 223 x^* is a stationary point of (2.1).

224 The remaining case is that there are an infinite number of branches from phase
 225 two to phase one for $\{x^k\}_{k \in K}$. Then phase one is performed an infinite number
 226 of times at $k_1 < k_2 < \dots < \dots$, where $\{k_i\} \subseteq K$. By Theorem 3.4 of [8],
 227 $\lim_{k_i \rightarrow \infty} \|\nabla_{\Omega} f(x^{k_i+1})\| = 0$. Again we find x^* is a stationary point by using $\{x^{k_i+1}\} \rightarrow$
 228 x^* and the lower semicontinuity of $\|\nabla_{\Omega} f(\cdot)\|$. The proof is completed. \square

229 Identification properties of an algorithm for linearly constrained problems are
 230 significant from both a theoretical and a practical point of view [14, 25]. For a
 231 stationary point x^* , the set of strongly active constraints is defined by

$$232 \quad \mathcal{A}_+(x^*) = \mathcal{M}_E \cup \{i \in \mathcal{M}_I : c_i^T x^* = d_i, \text{ and } \exists \lambda^* \in \Lambda(x^*) \text{ such that } \lambda_i^* > 0\}.$$

In convex analysis, the face of a convex set Ω exposed by the vector $w \in R^n$ is

$$E[w] \equiv \operatorname{argmax}\{w^T x : x \in \Omega\}.$$

233 A computation based on the definition of a face shows that for the polyhedral set Ω
 234 given in (1.2),

$$235 \quad (2.13) \quad E[-\nabla f(x^*)] = \{x \in \Omega : c_i^T x = d_i \text{ if } \lambda_i^* > 0 \text{ for } i \in \mathcal{M}_I\},$$

236 where $\lambda^* \in \Lambda(x^*)$. Note that this expression is valid for any choice of Lagrange
 237 multipliers $\lambda^* \in \Lambda(x^*)$.

238 We say that the linear independence constraint qualification (LICQ) holds at a
 239 point $x \in \Omega$, if the gradients $c_i, i \in \mathcal{A}(x)$ are linearly independent.

THEOREM 2.3. *Let $\{x^k\}$ be a sequence generated by Algorithm 2.1 with $\epsilon = 0$
 which converges to x^* . Suppose that the LICQ holds at x^* , and for some $\delta > 0$, g
 is Lipschitz continuous in $B_{\delta}(x^*)$ with a Lipschitz constant ϱ . Then there exists an
 integer $\hat{k}_0 > 0$ such that*

$$\mathcal{A}_+(x^*) \subseteq \mathcal{A}(x^k) \quad \text{and} \quad x^k \in E[-\nabla f(x^*)], \quad \text{for } k \geq \hat{k}_0.$$

240 *Proof.* Since $\{x^k\}$ converges to x^* , there exists $k_0 > 0$ such that $x^k \in B_{\delta}(x^*)$ for
 241 any $k \geq k_0$. Using the definition of $y(x, \alpha)$ in (2.6) and the Lipschitz continuity of g
 242 with the Lipschitz constant ϱ in $B_{\delta}(x^*)$, we have for any $\alpha > 0$ and $k \geq k_0$,

$$\begin{aligned} 243 \quad \|y(x^k, \alpha) - x^*\| &= \|y(x^k, \alpha) - y(x^*, \alpha)\| \\ 244 \quad &= \|P_{\Omega}[x^k - \alpha g(x^k)] - P_{\Omega}[x^* - \alpha g(x^*)]\| \\ 245 \quad &\leq \|x^k - x^* + \alpha(g(x^*) - g(x^k))\| \\ 246 \quad &\leq (1 + \alpha\varrho)\|x^k - x^*\|. \end{aligned}$$

247 Since $\{x^k\}$ converges to x^* , there is an integer $\bar{k} > 0$ such that $\mathcal{F}(x^*) \subseteq \mathcal{F}(y(x^k, \alpha))$ for
 248 $k \geq \bar{k}$. We know that $\Lambda(x^*)$ is a singleton, since the gradients of the active constraints
 249 at x^* are linearly independent. Thus $\Lambda(x^*, \alpha) = \alpha\Lambda(x^*)$ is also a singleton for any
 250 given $\alpha > 0$. Moreover, $\Lambda(x^k, \alpha)$ is a singleton for $k \geq \bar{k}$, because $\mathcal{A}(y(x^k, \alpha)) \subseteq$
 251 $\mathcal{A}(x^*)$ for $k \geq \bar{k}$ and the gradients of the active constraints at $y(x^k, \alpha)$ are linearly
 252 independent.

253 Consider the linear system

$$254 \quad (2.14) \quad q + \sum_{i \in \mathcal{M}} \lambda_i c_i = 0, \quad \lambda_i \geq 0 \text{ for } i \in \mathcal{M}_I, \quad \lambda_i = 0 \text{ for } i \in \mathcal{F}(x^*).$$

Let

$$p_1 = y(x^k, \alpha) - x^k + \alpha g(x^k), \quad \text{and} \quad p_2 = y(x^*, \alpha) - x^* + \alpha g(x^*).$$

According to (2.7), $\lambda^k \in \Lambda(x^k, \alpha)$ is feasible in the linear system (2.14) with $q = p_1$. And by (2.5) and (2.8), it is easy to see that for $\lambda^* \in \Lambda(x^*)$, $\alpha \lambda^* \in \Lambda(x^*, \alpha)$ is also feasible in the same system (2.14) but with $q = p_2$. Hence by Hoffman's result (see, e.g., Theorem 7.12 of [32]) and the fact that $\Lambda(x^*, \alpha)$ is a singleton, there exists a positive constant ν , independent of p_1 and p_2 and depending only on c_i , $i \in \mathcal{M}$, such that

$$\|\lambda^k - \alpha \lambda^*\| \leq \nu \|p_1 - p_2\| \leq 2\nu(1 + \alpha \varrho) \|x^k - x^*\|.$$

255 For any $i_0 \in \mathcal{M}_I \cap \mathcal{A}_+(x^*)$, the Lagrange multiplier $\lambda^* \in \Lambda(x^*)$ satisfies $\lambda_{i_0}^* > 0$.
 256 Thus there exists an integer $\tilde{k}_{i_0} > 0$ such that $\lambda_{i_0}^k > 0$ for all $k \geq \tilde{k}_{i_0}$. Now we
 257 consider (2.6) and its first-order optimality conditions given in (2.7). We find that
 258 $c_{i_0}^T y(x^k, \alpha) = d_{i_0}$ by complementarity and hence $i_0 \in \mathcal{A}(y(x^k, \alpha))$. Let

$$259 \quad \tilde{k} = \max\{\tilde{k}_i, i \in \mathcal{M}_I \cap \mathcal{A}_+(x^*)\} \quad \text{and} \quad \hat{k} = \max\{\bar{k}, \tilde{k}\}.$$

Clearly for any $i \in \mathcal{A}_+(x^*)$ and any given $\alpha > 0$,

$$i \in \mathcal{A}(y(x^k, \alpha)) \quad \text{for all } k \geq \hat{k}.$$

260 We need to consider two possible cases.

261 *Case 1:* There exists an integer $\hat{k}_1 \geq \hat{k}$ such that $x^{\hat{k}_1+1}$ is obtained from the PG
 262 step in Algorithm 2.1. Then for any $k \geq \hat{k}_1$ such that x^{k+1} is obtained from x^k by
 263 the PG step in Algorithm 2.1, we know by (2.6)

$$264 \quad x^{k+1} = P_\Omega[x^k - \alpha_k g(x^k)] = y(x^k, \alpha_k),$$

265 and consequently $i \in \mathcal{A}(x^{k+1})$ for any $i \in \mathcal{A}_+(x^*)$. Since no active constraint can be
 266 freed by the LCO step in phase two, we get

$$267 \quad i \in \mathcal{A}(x^k) \quad \text{for any } k \geq \hat{k}_1 + 1.$$

268 *Case 2:* x^{k+1} is obtained from the LCO step in phase two for any $k \geq \hat{k}$. By F2
 269 of the LCO Requirements, we find $\mathcal{A}(x^k) \subseteq \mathcal{A}(x^{k+1})$ for all $k \geq \hat{k}$. Then the active
 270 constraints become unchanged after a finite number of steps. Thus there exists an
 271 integer $\hat{k}_2 > \hat{k}$ such that

$$272 \quad \mathcal{A}(x^k) \equiv \tilde{\mathcal{A}} \subseteq \mathcal{A}(x^*) \text{ for all } k \geq \hat{k}_2.$$

273 By the definition of $g^{\tilde{\mathcal{A}}}(x^k)$, and the first-order optimality conditions at the global
 274 optimizer $g^{\tilde{\mathcal{A}}}(x^k)$, there exists a unique vector $\pi^k \in R^m$ such that

$$275 \quad (2.15) \quad \begin{aligned} & g^{\tilde{\mathcal{A}}}(x^k) - g(x^k) - \sum_{i \in \tilde{\mathcal{A}}} \pi_i^k c_i = 0, \\ & c_i^T g^{\tilde{\mathcal{A}}}(x^k) = 0, \quad i \in \tilde{\mathcal{A}}, \quad \pi_i^k = 0 \text{ if } i \notin \tilde{\mathcal{A}}. \end{aligned}$$

276 Here the vector π^k is unique because the column vectors c_i , $i \in \tilde{\mathcal{A}} \subseteq \mathcal{A}(x^*)$ are
 277 linearly independent. Similarly, by the strong convexity of problem (2.3) with x

278 being replaced by x^* , and the linear independence of $\{c_i, i \in \mathcal{A}(x^*)\}$, there exist a
 279 unique vector $g^{\mathcal{A}}(x^*) \in R^n$ and a unique vector $\lambda \in R^m$ such that

$$280 \quad (2.16) \quad \begin{aligned} &g^{\mathcal{A}}(x^*) - g(x^*) - \sum_{i \in \mathcal{A}(x^*)} \lambda_i c_i = 0, \\ &c_i^T g^{\mathcal{A}}(x^*) = 0, \quad i \in \mathcal{A}(x^*), \quad \lambda_i = 0 \text{ if } i \notin \mathcal{A}(x^*). \end{aligned}$$

281 And there exists a unique vector $\lambda^* \in R^m$ such that

$$282 \quad (2.17) \quad g(x^*) = - \sum_{i \in \mathcal{A}(x^*)} \lambda_i^* c_i, \quad \lambda_i^* = 0 \text{ if } i \notin \mathcal{A}(x^*),$$

283 since x^* is a stationary point of (2.1) and the gradients of the active constraints at
 284 x^* are linearly independent.

285 We get $g^{\mathcal{A}}(x^*) = 0$ and $\lambda = \lambda^*$, by comparing (2.16), (2.17) and using the
 286 uniqueness of $g^{\mathcal{A}}(x^*)$ and λ in (2.16). Moreover, $\liminf_{k \rightarrow \infty} g^{\tilde{\mathcal{A}}}(x^k) = 0$ according to
 287 F3 of the LCO Requirements. Let $\{k_j\} \subseteq \{k\}$ be an infinite subsequence such that
 288 $\lim_{k_j \rightarrow \infty} g^{\tilde{\mathcal{A}}}(x^{k_j}) = 0$. Taking limit to the first linear system in (2.15), we have

$$289 \quad (2.18) \quad 0 = \lim_{k_j \rightarrow \infty} g^{\tilde{\mathcal{A}}}(x^{k_j}) = g(x^*) + \sum_{i \in \tilde{\mathcal{A}}} \lim_{k_j \rightarrow \infty} \pi_i^{k_j} c_i.$$

290 Comparing (2.17) and (2.18), and noting the uniqueness of λ^* in (2.17), we find

$$291 \quad \lim_{k_j \rightarrow \infty} \pi_i^{k_j} = \lambda_i^* > 0 \quad \text{for any } i \in \mathcal{A}_+(x^*) \setminus \mathcal{M}_E.$$

292 Since $\pi_i^k = 0$ if $i \notin \tilde{\mathcal{A}}$ for k sufficiently large according to (2.15), we know

$$293 \quad \lim_{k_j \rightarrow \infty} \pi_i^{k_j} = 0 \quad \text{for any } i \in \mathcal{A}_+(x^*) \setminus \tilde{\mathcal{A}}.$$

294 This indicates $\mathcal{A}_+(x^*) \setminus \tilde{\mathcal{A}} = \emptyset$. Hence for any $i \in \mathcal{A}_+(x^*)$, we get $i \in \tilde{\mathcal{A}} \equiv \mathcal{A}(x^k)$ for
 295 $k \geq \hat{k}_2$.

Thus in any case, there exists an index \hat{k}_0 ($\hat{k}_0 = \hat{k}_1 + 1$ if Case 1 occurs, and
 $\hat{k}_0 = \hat{k}_2$ if Case 2 happens otherwise) such that

$$\mathcal{A}_+(x^*) \subseteq \mathcal{A}(x^k) \quad \text{for } k \geq \hat{k}_0.$$

This, combined with (2.13), implies

$$x^k \in E[-\nabla f(x^*)] \quad \text{for } k \geq \hat{k}_0.$$

296 We complete the proof. □

297 Based on the identification properties analyzed above, we will show the local
 298 convergence result that only iterations in phase two occur for k sufficiently large, if
 299 we further assume that the strong second-order sufficient optimality condition holds at
 300 x^* . A stationary point x^* of (2.1) satisfies the strong second-order sufficient optimality
 301 condition if there exists $\sigma > 0$ such that

$$302 \quad (2.19) \quad v^T \nabla^2 f(x^*) v \geq \sigma \|v\|^2$$

303 for all $v \in R^n$ such that $c_i^T v = 0$ for all $i \in \mathcal{A}_+(x^*)$.

304 LEMMA 2.4. Let $\{x^k\}$ be a sequence generated by Algorithm 2.1 with $\epsilon = 0$ which
 305 converges to x^* . Suppose that the LICQ holds at x^* , and for some $\delta > 0$, g is Lipschitz
 306 continuous in $B_\delta(x^*)$ with a Lipschitz constant ϱ . Then

$$307 \quad (2.20) \quad \|\nabla_\Omega f(x^k)\| \leq \varrho \|x^k - x^*\| \quad \text{for } k \text{ sufficiently large.}$$

308 *Proof.* From the nonexpansive property of the projection operator,

$$309 \quad (2.21) \quad \begin{aligned} \|\nabla_\Omega f(x^k)\| &= \|P_{T(x^k)}[-g(x^k)] - P_{T(x^k)}[-g(x^*)] + P_{T(x^k)}[-g(x^*)]\| \\ &\leq \|g(x^k) - g(x^*)\| + \|P_{T(x^k)}[-g(x^*)]\|. \end{aligned}$$

310 Similarly,

$$311 \quad (2.22) \quad \begin{aligned} &\|P_{T(x^k)}[-g(x^*)]\| \\ &= \|P_{T(x^k)}[-g(x^*)] - P_{T(x^k)}[-g(x^k)] + P_{T(x^k)}[-g(x^k)]\| \\ &\leq \|g(x^k) - g(x^*)\| + \|\nabla_\Omega f(x^k)\|. \end{aligned}$$

312 From (2.21) and (2.22),

$$313 \quad \|\nabla_\Omega f(x^k)\| - \|g(x^k) - g(x^*)\| \leq \|P_{T(x^k)}[-g(x^*)]\| \leq \|\nabla_\Omega f(x^k)\| + \|g(x^k) - g(x^*)\|.$$

314 Theorem 2.3 guarantees that there is an integer \hat{k}_0 such that $x^k \in E[-\nabla f(x^*)]$ for all
 315 $k \geq \hat{k}_0$. Thus according to Theorem 3.1 of [25], $\lim_{k \rightarrow \infty} \|\nabla_\Omega f(x^k)\| = 0$. This, com-
 316 bined with (2.22) and the facts that $\{x^k\} \rightarrow x^*$ and g is locally Lipschitz continuous
 317 at x^* , yields

$$318 \quad (2.23) \quad \lim_{k \rightarrow \infty} \|P_{T(x^k)}[-g(x^*)]\| = 0.$$

319 By direct computation,

$$320 \quad (2.24) \quad T(x^k) = \{v : c_i^T v = 0, i \in \mathcal{M}_E; \quad c_i^T v \leq 0, i \in \mathcal{M}_I \cap \mathcal{A}(x^k)\}.$$

321 When x^k is sufficiently near x^* , we know $\mathcal{F}(x^*) \subseteq \mathcal{F}(x^k)$. Then by Theorem 2.3, we
 322 find

$$323 \quad (2.25) \quad \mathcal{A}_+(x^*) \subseteq \mathcal{A}(x^k) \subseteq \mathcal{A}(x^*).$$

From the inclusions in (2.25) and the fact that $\mathcal{A}(x^*)$ has finite number of subsets,
 there are only a finite number of index sets $\mathcal{A}_1, \dots, \mathcal{A}_\nu$ for $\mathcal{A}(x^k)$, $k = 1, 2, \dots$. From
 the expression of $T(x^k)$ in (2.24), let us define

$$T_j = \{v : c_i^T v = 0, i \in \mathcal{M}_E; \quad c_i^T v \leq 0, i \in \mathcal{M}_I \cap \mathcal{A}_j\} \quad \text{for } j = 1, 2, \dots, \nu.$$

324 Without loss of generality, we assume

$$325 \quad \{T_1, T_2, \dots, T_t\} \subseteq \{T_1, T_2, \dots, T_\nu\}$$

326 is composed by all the elements in $\{T_1, T_2, \dots, T_\nu\}$ such that each T_j , $j = 1, 2, \dots, t$,
 327 contains an infinite number of $T(x^k)$, $k = 1, 2, \dots$. Hence we get $P_{T_j}[-g(x^*)] = 0$ for
 328 $j = 1, 2, \dots, t$, according to (2.23). Consequently, for all k sufficiently large, we have

$$329 \quad P_{T(x^k)}[-g(x^*)] \in \{P_{T_1}[-g(x^*)], P_{T_2}[-g(x^*)], \dots, P_{T_t}[-g(x^*)]\},$$

330 which indicates

$$331 \quad (2.26) \quad P_{T(x^k)}[-g(x^*)] = 0 \quad \text{for all } k \text{ sufficiently large.}$$

332 Substituting (2.26) into (2.21) and using the Lipschitz continuity of g with the Lips-
 333 chitz constant ϱ in $B_\delta(x^*)$, we get our desired result (2.20). \square

334 LEMMA 2.5. Let $\{x^k\}$ be a sequence generated by Algorithm 2.1 with $\epsilon = 0$ which
 335 converges to x^* . If f is twice continuously differentiable in a neighborhood of x^* ,
 336 the LICQ holds at x^* , and the strong second-order sufficient optimality condition in
 337 (2.19) holds at x^* , then there exists $\theta^* > 0$ such that

$$338 \quad (2.27) \quad \|g^{\mathcal{A}}(x^k)\| \geq \theta^* \|\nabla_{\Omega} f(x^k)\| \quad \text{for all } k \text{ sufficiently large.}$$

339 *Proof.* By Theorem 2.3, $\mathcal{A}_+(x^*) \subseteq \mathcal{A}(x^k)$ for $k \geq k_0$. Thus $x^k - x^*$ satisfies
 340 $c_i^T(x^k - x^*) = 0$ for all $i \in \mathcal{A}_+(x^*)$ and $k \geq k_0$. By the strong second-order suffi-
 341 cient optimality condition, we find x^* is a strict local minimizer of (2.1), and for k
 342 sufficiently large,

$$343 \quad (2.28) \quad (x^k - x^*)^T (g(x^k) - g(x^*)) \geq 0.5\sigma \|x^k - x^*\|^2.$$

344 Using the first-order necessary optimality conditions for a local minimizer of (2.1), we
 345 know that there exists a multiplier $\lambda^* \in R^m$ such that

$$346 \quad (2.29) \quad g(x^*) + \sum_{i \in \mathcal{M}} \lambda_i^* c_i = 0, \quad (d_i - c_i^T x^*) \lambda_i^* = 0, \quad i \in \mathcal{M}; \quad \lambda_i^* \geq 0, \quad i \in \mathcal{M}_I.$$

347 We have for k sufficiently large, $\mathcal{A}_+(x^*) \subseteq \mathcal{A}(x^k)$ and $d_i - c_i^T x^* = 0 = d_i - c_i^T x^k$ when
 348 $i \in \mathcal{A}_+(x^*)$, and $\lambda_i^* = 0$ when $i \notin \mathcal{A}_+(x^*)$. Hence

$$349 \quad \lambda_i^* c_i^T (x^k - x^*) = \lambda_i^* [(d_i - c_i^T x^*) - (d_i - c_i^T x^k)] = 0 \quad \text{for all } i \in \mathcal{M}.$$

350 This, combined with (2.29), yields

$$351 \quad (2.30) \quad (x^k - x^*)^T g(x^*) = (x^k - x^*)^T [g(x^*) + \sum_{i \in \mathcal{M}} \lambda_i^* c_i] = 0.$$

352 Denote here $\mathcal{S} = \mathcal{A}(x^k)$ for simplicity. The first-order optimality conditions for
 353 the minimizer $g^{\mathcal{S}}(x^k)$ in (2.2) implies the existence of $\lambda_{\mathcal{S}} \in R^{|\mathcal{S}|}$ such that

$$354 \quad (2.31) \quad g^{\mathcal{S}}(x^k) - g(x^k) + C_{\mathcal{S}} \lambda_{\mathcal{S}} = 0.$$

355 Because $\mathcal{A}(x^k) \subseteq \mathcal{A}(x^*)$ for $k \geq k_0$, we have $c_i^T(x^k - x^*) = 0$ for all $i \in \mathcal{S}$. Hence

$$356 \quad (2.32) \quad [C_{\mathcal{S}}^T (x^k - x^*)]^T \lambda_{\mathcal{S}} = 0, \quad \text{for all } k \text{ sufficiently large.}$$

357 By (2.31) and (2.32), we find

$$358 \quad (2.33) \quad (x^k - x^*)^T g(x^k) = (x^k - x^*)^T [g^{\mathcal{S}}(x^k) + C_{\mathcal{S}} \lambda_{\mathcal{S}}] = (x^k - x^*)^T g^{\mathcal{S}}(x^k).$$

359 Using the Cauchy-Schwarz inequality, (2.33), (2.28) and (2.30) sequentially, we get

$$\begin{aligned} 360 \quad \|x^k - x^*\| \|g^{\mathcal{S}}(x^k)\| &\geq (x^k - x^*)^T g^{\mathcal{S}}(x^k) \\ 361 \quad &= (x^k - x^*)^T g(x^k) \\ 362 \quad &= (x^k - x^*)^T [g(x^k) - g(x^*) + g(x^*)] \\ 363 \quad &\geq 0.5\sigma \|x^k - x^*\|^2. \end{aligned}$$

364 Reminding that $\mathcal{S} = \mathcal{A}(x^k)$, we have

$$365 \quad (2.34) \quad \|g^{\mathcal{A}}(x^k)\| \geq 0.5\sigma \|x^k - x^*\| \quad \text{for } k \text{ sufficiently large.}$$

366 This, together with Lemma 2.4, deduces (2.27) with $\theta^* = 0.5 \frac{\sigma}{\rho}$. \square

367 We are ready to show that the new active set method given in Algorithm 2.1 will
368 only perform the LCO within a finite number of iterations.

369 **THEOREM 2.6.** *Let $\{x^k\}$ be a sequence generated by Algorithm 2.1 with $\epsilon = 0$
370 which converges to x^* . If the assumptions in Lemma 2.5 hold, then within a finite
371 number of iterations, only phase two is executed.*

372 *Proof.* First we claim that phase two must occur within a finite number of iter-
373 erations. If on the contrary only phase one is occurred, then θ is decreased in each
374 iteration, and will be decreased to $\theta < \theta^*$ after a finite number of iterations. Then
375 according to Lemma 2.5, $\|g^A(x^k)\| > \theta\|\nabla_{\Omega}f(x^k)\|$ will occur. Once this holds, phase
376 one branches to phase two. This is a contradiction.

377 Once phase two is invoked, then phase two cannot branch to phase one infinite
378 times. Otherwise, θ will be reduced to $\theta < \theta^*$ and again $\|g^A(x^k)\| > \theta\|\nabla_{\Omega}f(x^k)\|$ will
379 occur, and after that phase two cannot branch to phase one. \square

380 Now we make clear the novelty of our new active set method in Algorithm 2.1,
381 compared to the active set method proposed by Hager and Zhang [18]. Algorithm
382 2.1 adopts the so-called piecewise PG method with $x^{k+1} = P_{\Omega}[x^k - \alpha_k g(x^k)]$ so that
383 the search direction within one iteration is along the projection arc [8]. While the
384 active set method by Hager and Zhang [18] chooses the so-called gradient projection
385 algorithm (GPA) in which the single projection is used to define the feasible search
386 direction $d^k = P_{\Omega}[x^k - \bar{\alpha}g(x^k)] - x^k$ where $\bar{\alpha} > 0$ is a fixed parameter, and the
387 next iterate point $x^{k+1} = x^k + s_k d^k$ is obtained by backtracking toward the starting
388 point along d^k . As pointed out by Bertsekas in subsection 2.3 of [3] that the iterates
389 obtained by the piecewise PG method used in this paper are more likely to be at the
390 boundary than the GPA used in [18]. Moreover, the finite identification property of
391 the new active set method is shown in Theorem 2.3. On contrast, after Lemma 6.1 of
392 [18], the authors stated that there is a fundamental difference between the GPA and
393 the PG method, and consequently they can not show the finite identification property
394 of the active set method in [18].

395 The main motivation of such modification lies in that Algorithm 2.1 guarantees
396 $\liminf_{k \rightarrow \infty} \|\nabla_{\Omega}f(x^k)\| = 0$, which is novel and essential in providing the convergence
397 result of the new smoothing active set method given in the next section. This conver-
398 gence result is stronger than that of the active set method in [18] which guarantees
399 $\liminf_{k \rightarrow \infty} \|d^1(x^k)\| = 0$, since by Lemma 2.2 of [8],

$$400 \quad (2.35) \quad \|\nabla_{\Omega}f(x^k)\| = \lim_{\alpha \downarrow 0} \frac{\|P_{\Omega}[x^k - \alpha \nabla f(x^k)] - P_{\Omega}[x^k]\|}{\alpha} \geq \|d^1(x^k)\|.$$

401 But $\liminf_{k \rightarrow \infty} \|d^1(x^k)\| = 0$ does not imply $\liminf_{k \rightarrow \infty} \|\nabla_{\Omega}f(x^k)\| = 0$ because the
402 norm of projected gradient is not continuous and can be large near the solution. This
403 can be explained by the following simple example.

404 **Example 1** Let us consider the linearly constrained strongly convex quadratic
405 programming

$$406 \quad \min \quad 0.01(10x_1 + x_2)^2 + 10(x_1 + 10.1x_2 + 1)^2 + x_3^2$$

$$407 \quad \text{s.t.} \quad x_2 \geq 1, \quad x_3 \geq 0.$$

We know $\mathcal{M} = \mathcal{M}_I = \{1, 2\}$ for this problem. It is easy to calculate that $x^* = (x_1^*, x_2^*, x_3^*)^T = (-10.1, 1, 0)^T$ is the unique global minimizer. The Lagrangian multipliers corresponding to the constraint $x_2 \geq 1$ and $x_3 \geq 0$ at x^* are $\lambda_1^* = 200$ and

$\lambda_2^* = 0$, respectively. Hence $\mathcal{A}_+(x^*) = \{1\}$, $\mathcal{A}(x^*) = \mathcal{M} = \{1, 2\}$, and x^* is a degenerate stationary point. The tangent cone to the feasible region at x^* , and the gradient at x^* are

$$T(x^*) = \{(d_1, d_2, d_3)^T \in R^3 : d_2 \geq 0, d_3 \geq 0\}, \quad \nabla f(x_1^*, x_2^*, x_3^*) = (0, 200, 0)^T.$$

408 Let $x^k = (x_1^k, x_2^k, x_3^k)^T = (-10.1 + (0.5)^{k/2}, 1 + (0.5)^k, (0.5)^k)^T \rightarrow x^*$ as $k \rightarrow +\infty$.
 409 By direct computation, the tangent cone to the feasible region at x^k is $T(x^k) = R^3$.
 410 Since f is twice continuously differentiable near x^* , we know that

$$411 \quad \nabla f(x^k) \rightarrow \nabla f(x^*) = (0, 200, 0)^T \quad \text{as } k \rightarrow \infty,$$

412 and consequently

$$413 \quad \|\nabla_{\Omega} f(x^k)\| = \|P_{T(x^k)}[-\nabla f(x^k)]\| = \|\nabla f(x^k)\| \rightarrow 200 \quad \text{as } k \rightarrow \infty.$$

414 Hence $\lim_{k \rightarrow \infty} \|\nabla_{\Omega} f(x^k)\| = 200 > 0$, although $\lim_{k \rightarrow \infty} \|d^1(x^k)\| = 0$.

415 *Remark 2.7.* Suppose $\{x^k\} \rightarrow x^*$, ∇f is locally Lipschitz continuous at x^* , and
 416 the active constraints are identified after finite iterations. Then there exists $k_0 > 0$
 417 such that $T(x^k) \equiv T(x^*)$ for all $k \geq k_0$, and consequently $\lim_{k \rightarrow \infty} \|\nabla_{\Omega} f(x^k)\| = 0$
 418 and $\lim_{k \rightarrow \infty} \|d^1(x^k)\| = 0$ are equivalent. However, the active set method in [18] may
 419 not identify the active constraints, but only owns the property in Lemma 6.2 of [18]
 420 that the violation of the constraints $c_i^T x - d_i = 0$ for $i \in \mathcal{A}_+(x^*)$ by iterate x^k is on
 421 the order of the error in x^k squared under certain conditions. Using Example 1, we
 422 find

$$423 \quad \bar{x}^k = \operatorname{argmin}_y \{\|x^k - y\| : y_2 = 1\} = (x_1^k, 1, x_3^k)^T,$$

424 and

$$425 \quad \lim_{k \rightarrow \infty} \frac{\|x^k - \bar{x}^k\|}{\|x^k - x^*\|^2} = \lim_{k \rightarrow \infty} \frac{|x_2^k - 1|}{\|x^k - x^*\|^2} = \lim_{k \rightarrow \infty} \frac{(0.5)^k}{(0.5)^k + (0.5)^{2k} + (0.5)^{2k}} = 1.$$

426 This indicates that although under certain conditions any sequence generated by the
 427 active set method [18] satisfies the property in Lemma 6.2 of [18], this property does
 428 not guarantee $\liminf_{k \rightarrow \infty} \|\nabla_{\Omega} f(x^k)\| = 0$ that we need in designing the smoothing
 429 active set method with convergence result.

430 **3. Smoothing active set method.** In this section, we develop a smoothing
 431 active set method for solving (1.1) with solid convergence result. Here the objective
 432 function f is continuous, but not necessarily Lipschitz continuous.

433 To characterize the stationary points of (1.1), we review first the concepts of rly
 434 several subdifferentials that are often used in nonsmooth analysis [6, 31] and references
 435 therein. Let $f : R^n \rightarrow R$ be a proper lower semi-continuous function and $x \in R^n$ be a
 436 point where $f(x)$ is finite. The Fréchet subdifferential, the limiting (or Mordukhovich)
 437 subdifferential, the horizontal (or singular Mordukhovich) subdifferential, and the
 438 Clarke subdifferential (Definition 1 of [6]) are defined respectively as

$$439 \quad \hat{\partial} f(x) := \{v : f(y) \geq f(x) + v^T(y - x) + o(\|y - x\|), \forall y\},$$

$$440 \quad \partial f(x) := \left\{ v : \exists x^k \xrightarrow{f} x, v^k \rightarrow v \text{ with } v^k \in \hat{\partial} f(x^k), \forall k \right\},$$

$$441 \quad \partial^{\infty} f(x) := \left\{ v : \exists x^k \xrightarrow{f} x, t_k v^k \rightarrow v, t_k \downarrow 0 \text{ with } v^k \in \hat{\partial} f(x^k), \forall k \right\},$$

$$442 \quad \partial^{\circ} f(x) := \operatorname{co}\{\partial f(x) + \partial^{\infty} f(x)\},$$

443 where $x^k \xrightarrow{f} x$ means that $x^k \rightarrow x$ and $f(x^k) \rightarrow f(x)$, and “cō” means the closure
 444 of convex hull. We say that x^* is a Clarke stationary point of (1.1), if there is
 445 $V \in \partial^\circ f(x^*)$ such that

$$446 \quad (3.1) \quad \langle V, x^* - z \rangle \leq 0 \quad \text{for all } z \in \Omega.$$

447 If there exists $V \in \partial f(x^*)$ such that (3.1) holds, then x^* is a limiting stationary point
 448 of (1.1). Under the basic qualification (BQ)

$$449 \quad (3.2) \quad -\partial^\infty f(x^*) \cap N_\Omega(x^*) = \{0\},$$

450 if x^* is a local minimizer, then x^* is a limiting stationary point (Rockafellar and Wets,
 451 Theorem 8.15 of [31]). It is easy to see that BQ in (3.2) holds if f is locally Lipschitz
 452 continuous at x^* , or x^* is an interior point of Ω . However, BQ often fails if f is
 453 non-Lipschitz at a boundary point x^* as pointed out in [9].

454 We use the following definition for smoothing function.

455 **DEFINITION 3.1.** *Let $f : R^n \rightarrow R$ be a continuous function. We call $\tilde{f} : R^n \times$
 456 $R_+ \rightarrow R$ a smoothing function of f , if $\tilde{f}(\cdot, \mu)$ is continuously differentiable in R^n for
 457 any $\mu \in R_{++}$, and for any $x \in R^n$,*

$$458 \quad (3.3) \quad \lim_{z \rightarrow x, \mu \downarrow 0} \tilde{f}(z, \mu) = f(x),$$

459 and there exists a constant $\kappa > 0$ and a function $\omega : R_{++} \rightarrow R_{++}$ such that

$$460 \quad (3.4) \quad |\tilde{f}(x, \mu) - f(x)| \leq \kappa \omega(\mu) \quad \text{with} \quad \lim_{\mu \downarrow 0} \omega(\mu) = 0.$$

461 For each fixed $\mu > 0$, the smooth subproblem is then defined by

$$462 \quad (3.5) \quad \min \tilde{f}(x, \mu) \quad \text{s.t.} \quad x \in \Omega,$$

463 and the projected gradient $\nabla_\Omega \tilde{f}(x, \mu)$ is defined by

$$464 \quad \nabla_\Omega \tilde{f}(x, \mu) \equiv P_{T(x)}[-\nabla_x \tilde{f}(x, \mu)] = \operatorname{argmin}\{\|v + \nabla_x \tilde{f}(x, \mu)\| : v \in T(x)\},$$

465 where $T(x)$ is the tangent cone to Ω at x . Now we present our smoothing active set
 466 method, Algorithm 3.1.

Algorithm 3.1 Smoothing active set method

- 1: Let $\hat{\gamma}$ be a positive constant, ζ be a constant in $(0, 1)$, and $n_1 > 0$ be a positive integer. Choose $x^0 \in \Omega$ and $\mu_0 > 0$.
 For $k \geq 0$:
 - 2: Let $y^{0,k} = x^k$, $j := 0$.
 - 3: **while** $\|\nabla_\Omega \tilde{f}(y^{j,k}, \mu_k)\| > \hat{\gamma} \mu_k$ or $j < n_1$, **do**
 - 4: Execute one iterate of the active set method in Algorithm 2.1 for (3.5) with
 $\mu = \mu_k$ from the initial point $y^{j,k}$ and get the new point $y^{j+1,k}$.
 Set $j := j + 1$.
 - 5: **end while**
 - 6: Set $x^{k+1} = y^{j,k}$.
 - 7: Choose $\mu_{k+1} \leq \zeta \mu_k$.
-

467 *Remark 3.2.* It is worth mentioning that Algorithm 3.1 can be extended to a
 468 general framework of smoothing method, since the new active set method in Algorithm
 469 2.1 that used in Algorithm 3.1 can be substituted by any other type of algorithm
 470 for minimizing smooth function (SA for short) on a closed convex set, as long as the
 471 algorithm satisfies the SA Requirement defined below. And then the same convergence
 472 result developed in this section can be obtained without difficulty.

473 **SA Requirement** For any fixed $\mu > 0$, let $\{x^k\}$ be generated by the SA that
 474 solves (3.5). Then

$$475 \quad \liminf_{k \rightarrow \infty} \nabla_{\Omega} \tilde{f}(x^k, \mu) = 0.$$

476
 477 When $\Omega = R^n$, then (3.5) reduces to unconstrained smooth optimization and
 478 hence $\nabla_{\Omega} \tilde{f}(x, \mu) = -\nabla f(x, \mu)$. Many unconstrained algorithms (UAs) for (3.5)
 479 meet the SA Requirement, e.g., the steepest descent method, the accelerated gra-
 480 dient method proposed by Nesterov, the conjugate gradient method, the trust region
 481 method, and the quasi-Newton method. When Ω is a general closed convex set, the
 482 projected gradient method satisfies the SA Requirement. When Ω is constructed by
 483 linear constraints defined in (1.2), the new active set method developed in section 2
 484 meets the SA Requirement as we desired. Although the proposed active set method
 485 is in spirit very similar to Hager and Zhang's approach [18], the satisfaction of the SA
 486 Requirement makes it necessary and novelty for building up the convergence of the s-
 487 moothing active set method that tackles linearly constrained non-Lipschitz nonconvex
 488 optimization problems.

489 Since we use a smoothing function in Algorithm 3.1, the convergence result is
 490 natural to connect with the smoothing function employed.

491 **DEFINITION 3.3.** We say that x^* is a stationary point of (1.1) associated with a
 492 smoothing function \tilde{f} , if

$$493 \quad (3.6) \quad \liminf_{x \rightarrow x^*, x \in \Omega, \mu \downarrow 0} \langle \nabla_x \tilde{f}(x, \mu), x - z \rangle \leq 0 \quad \text{for all } z \in \Omega.$$

494 For any fixed $x \in \Omega$, denote

$$495 \quad (3.7) \quad G_{\tilde{f}}(x) := \{V : \exists N \in \mathcal{N}_{\infty}^{\sharp}, x^{\nu} \xrightarrow{N} x, \mu_{\nu} \downarrow 0 \quad \text{with} \quad \nabla_x \tilde{f}(x^{\nu}, \mu_{\nu}) \xrightarrow{N} V\}.$$

By Corollary 8.47 (b) in [31], we have

$$\partial f(x) \subseteq G_{\tilde{f}}(x).$$

When f is Lipschitz continuous, it is shown in [7, 10, 31] that many smoothing func-
 tions satisfy the gradient consistency property

$$\partial^{\circ} f(x^*) = G_{\tilde{f}}(x^*).$$

496 Then the stationary point of (1.1) associated with \tilde{f} coincides to the Clarke stationary
 497 point, i.e., there exists $V \in \partial^{\circ} f(x^*)$ such that (3.1) holds. When f is continuously dif-
 498 ferentiable at x^* , then $\partial^{\circ} f(x^*) = \{\nabla f(x^*)\}$ and x^* coincides to the classic stationary
 499 point for smooth minimization problems.

500 Now we show that x^* being a stationary point of (1.1) associated with a smoothing
 501 function \tilde{f} is a necessary optimality condition for x^* being a local minimizer, without
 502 the requirement for BQ.

503 PROPOSITION 3.4. *For any given smoothing function \tilde{f} defined in Definition 3.1,*
 504 *if x^* is a local minimizer of (1.1), then x^* is a stationary point of (1.1) associated*
 505 *with \tilde{f} .*

506 *Proof.* Since x^* is a local minimizer of (1.1), there exists a constant $\delta > 0$ such
 507 that

$$508 \quad f(x^*) \leq f(x) \quad \text{for any } x \in B_\delta(x^*) \cap \Omega.$$

509 This, combined with (3.4) in Definition 3.1 for the smoothing function, yields that for
 510 all $x \in B_\delta(x^*) \cap \Omega$,

$$511 \quad (3.8) \quad \tilde{f}(x^*, \mu) \leq f(x^*) + \kappa\omega(\mu) \leq f(x) + \kappa\omega(\mu) \leq \tilde{f}(x, \mu) + 2\kappa\omega(\mu).$$

512 For any $z \in \Omega$, let $x_\mu = x^* + \sqrt{\omega(\mu)}(z - x^*)$. Since Ω is a convex set and $\lim_{\mu \downarrow 0} \omega(\mu) =$
 513 0 , we get $x_\mu \in B_\delta(x^*) \cap \Omega$ for all μ sufficiently small and $x_\mu \rightarrow x^*$ as $\mu \downarrow 0$. By Taylor's
 514 theorem,

$$515 \quad \begin{aligned} \tilde{f}(x^*, \mu) &= \tilde{f}(x_\mu, \mu) + \nabla_x \tilde{f}(x_\mu, \mu)^T (x^* - x_\mu) + o(\|x^* - x_\mu\|) \\ 516 \quad (3.9) \quad &= \tilde{f}(x_\mu, \mu) + \sqrt{\omega(\mu)} \nabla_x \tilde{f}(x_\mu, \mu)^T (x^* - z) + o(\sqrt{\omega(\mu)}). \end{aligned}$$

517 Substituting (3.9) into the left side of (3.8) and replacing x by x_μ into the right side
 518 of (3.8), we get

$$519 \quad \sqrt{\omega(\mu)} \nabla_x \tilde{f}(x_\mu, \mu)^T (x^* - z) + o(\sqrt{\omega(\mu)}) \leq 2\kappa\omega(\mu).$$

520 Dividing both sides of the above inequality by $\sqrt{\omega(\mu)}$, and taking the limit as $\mu \downarrow 0$,
 521 we find

$$522 \quad (3.10) \quad \limsup_{\mu \downarrow 0} \langle \nabla_x \tilde{f}(x_\mu, \mu), x^* - z \rangle \leq 0.$$

523 Note that

$$524 \quad \langle \nabla_x \tilde{f}(x_\mu, \mu), x_\mu - z \rangle = (1 - \sqrt{\omega(\mu)}) \langle \nabla_x \tilde{f}(x_\mu, \mu), x^* - z \rangle.$$

525 This, together with (3.10), yields that

$$526 \quad \liminf_{\mu \downarrow 0} \langle \nabla_x \tilde{f}(x_\mu, \mu), x_\mu - z \rangle = \liminf_{\mu \downarrow 0} (1 - \sqrt{\omega(\mu)}) \langle \nabla_x \tilde{f}(x_\mu, \mu), x^* - z \rangle \leq 0,$$

527 which indicates

$$528 \quad (3.11) \quad \liminf_{x \rightarrow x^*, x \in \Omega, \mu \downarrow 0} \langle \nabla_x \tilde{f}(x, \mu), x - z \rangle \leq 0 \quad \text{for all } z \in \Omega.$$

529 Hence (3.6) holds and x^* is a stationary point of (1.1) with respect to \tilde{f} . \square

530 Now we are ready to give the global convergence result of Algorithm 3.1.

531 THEOREM 3.5. *Assume Assumption 2.1 holds. Then any accumulation point x^**
 532 *of $\{x^k\}$ generated by Algorithm 3.1 is a stationary point of (1.1) associated with the*
 533 *smoothing function f .*

534 *Proof.* By (3.4) of Definition 3.1, for each fixed $\mu > 0$,

$$535 \quad f(x) - \kappa\omega(\mu) \leq \tilde{f}(x, \mu) \leq f(x) + \kappa\omega(\mu).$$

Then for each fixed $\mu > 0$,

$$\mathcal{L}_{\mu, \Gamma} = \{x \in \Omega : \tilde{f}(x, \mu) \leq \Gamma\}$$

536 is bounded for any Γ , because $\tilde{f}(x, \mu) \leq \Gamma$ implies $f(x) \leq \Gamma + \kappa\omega(\mu)$ and $\mathcal{L}_{\Gamma + \kappa\omega(\mu)}$ is
537 bounded by Assumption 2.1.

538 By (2.11) of Theorem 2.2, we know Algorithm 3.1 is well-defined and

$$539 \quad (3.12) \quad \|\nabla_{\Omega} \tilde{f}(x^{k+1}, \mu_k)\| \leq \hat{\gamma}\mu_k, \quad \text{and} \quad \lim_{k \rightarrow \infty} \mu_k = 0.$$

540 According to Calamai and Moré [8],

$$541 \quad (3.13) \quad \min\{\langle \nabla_x \tilde{f}(x^{k+1}, \mu_k), v \rangle : v \in T(x^{k+1}), \|v\| \leq 1\} = -\|\nabla_{\Omega} \tilde{f}(x^{k+1}, \mu_k)\|.$$

542 For any $z \in \Omega$, it is easy to see that

$$543 \quad v = \frac{z - x^{k+1}}{\|z - x^{k+1}\|} \in T(x^{k+1}) \quad \text{and} \quad \|v\| = 1,$$

544 and hence by (3.13)

$$545 \quad \langle \nabla_x \tilde{f}(x^{k+1}, \mu_k), x^{k+1} - z \rangle \leq \|\nabla_{\Omega} \tilde{f}(x^{k+1}, \mu_k)\| \|z - x^{k+1}\|.$$

546 This, combined with (3.12), yields

$$547 \quad (3.14) \quad \langle \nabla_x \tilde{f}(x^{k+1}, \mu_k), x^{k+1} - z \rangle \leq \hat{\gamma}\mu_k \|z - x^{k+1}\| \quad \text{for any } z \in \Omega.$$

548 Since x^* is an accumulation point of $\{x^k\}$, there exists an infinite sequence $\hat{K} \in$
549 $\mathcal{N}_{\infty}^{\#}$ such that $\lim_{k \rightarrow \infty, k \in \hat{K}} x^k = x^*$. Let us denote $K = \{k - 1 : k \in \hat{K}\}$ and then
550 $\lim_{k \rightarrow \infty, k \in K} x^{k+1} = x^*$. We get from (3.14) that

$$551 \quad (3.15) \quad \liminf_{k \rightarrow \infty, k \in K} \langle \nabla_x \tilde{f}(x^{k+1}, \mu_k), x^{k+1} - z \rangle \leq 0 \quad \text{for any } z \in \Omega.$$

552 Therefore x^* is a stationary point of (1.1) associated with \tilde{f} . □

553 The objective function f in this paper is a general non-Lipschitz nonconvex function
554 which is broader than that considered in [4, 5, 11, 26]. In [5], the optimality and
555 complexity for the convexly-constrained minimization problem are considered with
556 the objective function in the following form

$$557 \quad f(x) := \Theta(x) + c(h(x)), \quad \text{with} \quad h(x) := (h_1(D_1^T x), h_2(D_2^T x), \dots, h_m(D_m^T x))^T.$$

558 Here $\Theta : R^n \rightarrow R$ and $c : R^m \rightarrow R$ are continuously differentiable, $D_i \in R^{n \times r}$, and
559 $h_i : R^r \rightarrow R$, $i = 1, \dots, m$ are continuous, but not necessarily Lipschitz continuous.
560 This type of functions include all the objective functions considered in [4, 11, 26]. A
561 generalized stationary point based on the generalized directional derivative is proposed
562 in Definition 2 of [5], which is shown to be a necessary optimality condition, and
563 satisfies the necessary optimality conditions given or used in [4, 11, 26]. Note that

564 any $v \in T(x^{k+1})$ and $\|v\| \leq 1$, there exists $z \in \Omega$ such that $v = z - x^{k+1} \in T(x^{k+1})$.
 565 By (3.14) of Theorem 3.5 and $\|z - x^{k+1}\| \leq 1$,

$$566 \quad \langle \nabla_x \tilde{f}(x^{k+1}, \mu_k), v \rangle = \langle \nabla_x \tilde{f}(x^{k+1}, \mu_k), z - x^{k+1} \rangle \geq -\hat{\gamma} \mu_k \|z - x^{k+1}\| \geq -\hat{\gamma} \mu_k,$$

567 which implies that (44) in Corollary 2 of [5] holds, and consequently any accumulation
 568 point of $\{x^k\}$ generated by the smoothing active set method is also a generalized
 569 stationary point of (1.1) defined in [5] for the same type of functions in [5] and Ω
 570 defined in (1.2).

571 *Remark 3.6.* In Algorithm 3.1, we require for each fixed μ_k , the iterations of the
 572 inner loop is no less than n_1 . This strategy has no effect for convergence analysis, but
 573 aims to enhance the computational performance of finding a better stationary point
 574 with respect to \tilde{f} .

575 **3.1. $\ell_2 - \ell_p$ sparse optimization model.** Problem (1.3) is a special case of
 576 problem (1.1), for which we show that Algorithm 3.1 has stronger convergence results
 577 than that in Theorem 3.5.

578 For $|t|$, we construct its smoothing function as follows,

$$579 \quad (3.16) \quad s_\mu(t) = \begin{cases} |t| & \text{if } |t| \geq \mu, \\ \frac{t^2}{2\mu} + \frac{\mu}{2} & \text{if } |t| < \mu. \end{cases}$$

580 By simple computation, for any $p \in (0, 1)$ and any $t \in \mathbb{R}$, we have $|s_\mu(t)^p - |t|^p| \leq 2\mu^p$.
 581 We then easily find that

$$582 \quad \tilde{f}(x, \mu) = \|Ax - b\|^2 + \tau \sum_{i=1}^n (s_\mu(x_i))^p$$

583 is a smoothing function of the objective function f in (1.3), and for any $x \in \mathbb{R}^n$,

$$584 \quad (3.17) \quad |\tilde{f}(x, \mu) - f(x)| \leq \kappa \mu^p, \quad \text{with } \kappa = 2\tau n.$$

585 The gradient of $\tilde{f}(x, \mu)$ is

$$586 \quad (3.18) \quad \nabla_x \tilde{f}(x, \mu) = 2A^T(Ax - b) + \tau p \sum_{i=1}^n (s_\mu(x_i))^{p-1} s'_\mu(x_i).$$

587 **THEOREM 3.7.** *There exists at least one accumulation point x^* of $\{x^k\}$ generated*
 588 *by Algorithm 3.1 with the smoothing function \tilde{f} . Suppose $\lim_{k \rightarrow \infty, k \in K} x^{k+1} = x^*$. Then*

589 *$\{\lim_{k \rightarrow \infty, k \in K} \nabla_x \tilde{f}(x^{k+1}, \mu_k)\}$ is nonempty and bounded, and x^* is a limiting stationary*
 590 *point of (1.3).*

591 *Proof.* Assumption 2.1 holds for f in (1.3), since the objective function in (1.3)
 592 satisfies that $f(x) \rightarrow +\infty$ if $\|x\| \rightarrow +\infty$. Moreover, we know from (3.17) that

$$593 \quad \tilde{f}(x^{j+1}, \mu_j) - f(x^{j+1}) \geq -\kappa \mu_j^p \quad \text{and} \quad \tilde{f}(x^j, \mu_j) - f(x^j) \leq \kappa \mu_j^p.$$

594 Therefore for any natural number k ,

$$\begin{aligned} 595 \quad f(x^{k+1}) &\leq \tilde{f}(x^{k+1}, \mu_k) + \kappa \mu_k^p \leq \tilde{f}(x^k, \mu_k) + \kappa \mu_k^p \leq f(x^k) + 2\kappa \mu_k^p \\ 596 &\leq \dots \\ 597 &\leq f(x^0) + 2\kappa[\mu_0^p + (\zeta \mu_0)^p + (\zeta^2 \mu_0)^p + \dots + (\zeta^k \mu_0)^p] \\ 598 &\leq f(x^0) + 2\kappa \mu_0^p \frac{1}{1 - \zeta^p}. \end{aligned}$$

599 Hence $\{x^k\}$ is bounded and there exists at least one accumulation point x^* of $\{x^k\}$
 600 generated by Algorithm 3.1.

601 For any index i_0 such that $x_{i_0}^* > 0$, by direct computation,

$$602 \quad \lim_{k \rightarrow \infty, k \in K} (\nabla_x \tilde{f}(x^{k+1}, \mu_k))_{i_0} = (2A^T(Ax^* - b))_{i_0} + \tau p (x_{i_0}^*)^{p-1}.$$

603 For i_0 such that $x_{i_0}^* = 0$, let $K_2 = \{k \in K : x_{i_0}^{k+1} > 0\}$. If K_2 is an infinite
 604 subsequence, then we define $z^{k+1,1}$ and $z^{k+1,2}$ in R_+^n for each $k \in K_2$, where

$$605 \quad z_i^{k+1,1} = \begin{cases} x_i^{k+1} & \text{if } i \neq i_0, \\ 0 & \text{if } i = i_0, \end{cases} \quad \text{and} \quad z_i^{k+1,2} = \begin{cases} x_i^{k+1} & \text{if } i \neq i_0, \\ 2x_i^{k+1} & \text{if } i = i_0. \end{cases}$$

606 Replacing $z^{k+1,1}$ and $z^{k+1,2}$ in (3.14) of Theorem 3.5 respectively, we get eventually

$$607 \quad -\hat{\gamma}\mu_k \leq (\nabla_x \tilde{f}(x^{k+1}, \mu_k))_{i_0} \leq \hat{\gamma}\mu_k \quad \text{for any } k \in K_2,$$

608 and consequently

$$609 \quad (3.19) \quad \lim_{k \rightarrow \infty, k \in K_2} (\nabla_x \tilde{f}(x^{k+1}, \mu_k))_{i_0} = 0.$$

610 Otherwise, there exists an integer $\bar{k} > 0$ such that $x_{i_0}^{k+1} = 0$ for all $k \geq \bar{k}$, $k \in K$. In
 611 this case

$$\begin{aligned} 612 \quad (\nabla_x \tilde{f}(x^{k+1}, \mu_k))_{i_0} &= (2A^T(Ax^{k+1} - b))_{i_0} + \tau p (s_{\mu_k}(x_{i_0}^{k+1}))^{p-1} s'_{\mu_k}(x_{i_0}^{k+1}) \\ 613 &= (2A^T(Ax^{k+1} - b))_{i_0} + \tau p \left(\frac{\mu_k}{2}\right)^{p-1} \frac{x_{i_0}^{k+1}}{\mu_k} \\ 614 &= (2A^T(Ax^{k+1} - b))_{i_0} \quad \text{for all } k \geq \bar{k}, k \in K. \end{aligned}$$

615 Consequently

$$616 \quad (3.20) \quad \lim_{k \rightarrow \infty, k \in K} (\nabla_x \tilde{f}(x^{k+1}, \mu_k))_{i_0} = (2A^T(Ax^* - b))_{i_0}.$$

617 Combining (3.19) and (3.20), we can easily find that any accumulation point $V \in R^n$
 618 of $\{\nabla_x \tilde{f}(x^{k+1}, \mu_k)\}_K$ is of the special form

$$619 \quad (3.21) \quad V_i = \begin{cases} (2A^T(Ax^* - b))_i + \tau p (x_i^*)^{p-1} & \text{if } x_i^* > 0 \\ (2A^T(Ax^* - b))_i \text{ or } 0, & \text{if } x_i^* = 0 \end{cases}$$

620 that is bounded.

621 Furthermore, we know $V \in \partial f(x)$ by the definition of the limiting subdifferential,
 622 which indicates that x^* is also a limiting stationary point of (1.3). \square

623 **THEOREM 3.8.** *Let x^* be an accumulation point of a sequence $\{x^k\}$ generated by*
 624 *Algorithm 3.1 for solving (1.3). If $\mathcal{F}(x^*) = \emptyset$, then $x^* = 0$ is a local minimizer of*
 625 *(1.3). If $\mathcal{F}(x^*) \neq \emptyset$ and*

$$626 \quad (3.22) \quad 2(A^T A)_{\mathcal{F}(x^*)\mathcal{F}(x^*)} + \tau p(p-1) \text{diag}((x_{\mathcal{F}(x^*)}^*)^{p-2}) \quad \text{is positive definite,}$$

627 *then x^* is a strict local minimizer of (1.3).*

Proof. By Theorem 3.7, and (3.15) in the proof of Theorem 3.5, there exists an accumulation point V of $\{\lim_{k \rightarrow \infty, k \in K} \nabla_x \tilde{f}(x^{k+1}, \mu_k)\}$ in the form of (3.21) such that

$$\langle V, x^* - z \rangle \leq 0 \quad \text{for all } z \geq 0.$$

628 This indicates $V_i = 0$ for all $i \in \mathcal{F}(x^*)$.

629 Let us define $\varsigma_i := \frac{2}{\tau} \left(\max\{-(A^T(Ax^* - b))_i, 0\} + 1 \right)$ for all $i \in \mathcal{A}(x^*)$, and

$$630 \quad (3.23) \quad \bar{f}(x) := \|Ax - b\|^2 + \tau \sum_{i \in \mathcal{F}(x^*)} |x_i|^p + \tau \sum_{i \in \mathcal{A}(x^*)} \varsigma_i x_i.$$

631 Now we consider the minimization problem

$$632 \quad (3.24) \quad \min \bar{f}(x) \quad \text{s.t.} \quad x \geq 0,$$

633 whose objective function is twice continuously differentiable around $x^* \in R_+^n$. By
634 direct computation, $\bar{f}(x^*) = f(x^*)$ and the gradient $\nabla \bar{f}(x^*)$ has the form

$$635 \quad (\nabla \bar{f}(x^*))_i = \begin{cases} (2A^T(Ax^* - b))_i + \tau p(x_i^*)^{p-1} & \text{if } i \in \mathcal{F}(x^*), \\ (2A^T(Ax^* - b))_i + \tau \varsigma_i & \text{if } i \in \mathcal{A}(x^*). \end{cases}$$

636 Clearly, $(\nabla \bar{f}(x^*))_i = V_i = 0$ for all $i \in \mathcal{F}(x^*)$ and $(\nabla \bar{f}(x^*))_i \geq 2$ for all $i \in \mathcal{A}(x^*)$.

637 Therefore, x^* is a stationary point of (3.24) since

$$638 \quad (3.25) \quad x^* \geq 0, \quad \nabla \bar{f}(x^*) \geq 0, \quad x^{*T} \nabla \bar{f}(x^*) = 0.$$

Note that for any $p \in (0, 1)$,

$$\lim_{t \downarrow 0, t \neq 0} \frac{t^p}{t} = \lim_{t \downarrow 0, t \neq 0} t^{p-1} = +\infty.$$

639 Thus there exists $\delta_1 > 0$ such that for any $x \in B_{\delta_1}(x^*) \cap R_+^n$

$$640 \quad \varsigma_i x_i \leq x_i^p \quad \text{for all } i \in \mathcal{A}(x^*).$$

641 Consequently for any $x \in B_{\delta_1}(x^*) \cap R_+^n$,

$$642 \quad (3.26) \quad \bar{f}(x) - f(x) = \tau \sum_{i \in \mathcal{A}(x^*)} (\varsigma_i x_i - x_i^p) \leq 0.$$

If $\mathcal{F}(x^*) = \emptyset$, then $x^* = 0$ and $\bar{f}(x)$ in (3.23) is a convex function. Any stationary point of (3.24) is a global minimizer of (3.24). Hence

$$\bar{f}(x^*) \leq \bar{f}(x) \quad \text{for any } x \in R_+^n.$$

643 This, combined with (3.26), yields

$$644 \quad f(x^*) = \bar{f}(x^*) \leq \bar{f}(x) \leq f(x) \quad \text{for any } x \in B_{\delta_1}(x^*) \cap R_+^n.$$

645 Hence x^* is a local minimizer of (1.3).

646 Now we consider $\mathcal{F}(x^*) \neq \emptyset$. Noting (3.25), we know that (x^*, λ^*) satisfies
647 the KKT conditions if and only if $\lambda^* = \nabla \bar{f}(x^*)$. Since for any $i \in \mathcal{A}(x^*)$, $\lambda_i^* =$
648 $(\nabla \bar{f}(x^*))_i \geq 2$, it follows that the critical cone

$$649 \quad \mathcal{C}(x^*, \lambda^*) = \{d \in R^n : d_i = 0 \text{ for } i \in \mathcal{A}(x^*), \text{ and } d_i \geq 0 \text{ for } i \in \mathcal{F}(x^*)\}.$$

650 It is easy to see that (3.22) is equivalent to

$$651 \quad d^T \nabla^2 \bar{f}(x^*) d > 0 \quad \text{for any } d \in \mathcal{C}(x^*, \lambda^*), \quad d \neq 0,$$

652 which are the second-order sufficient conditions for x^* being a strict local minimizer
653 of (3.24). Then there exists $\delta > 0$ such that

$$654 \quad (3.27) \quad f(x^*) = \bar{f}(x^*) < \bar{f}(x) \quad \text{for any } x \in B_\delta(x^*) \cap R_+^n.$$

655 This, combined with (3.26), yields

$$656 \quad f(x^*) < f(x) \quad \text{for any } x \in B_{\check{\delta}}(x^*) \cap R_+^n,$$

657 where $\check{\delta} = \min\{\delta, \delta_1\}$. Hence x^* is a strict local minimizer of (1.3). \square

658 **4. Numerical experiments.** Hyperspectral image is a 3D image cube at hun-
659 dreds of contiguous and narrow spectral channels often used in earth observation and
660 remote sensing. Due to the low spatial resolution of hyperspectral cameras, pixels
661 are often a mixture of several spectra of materials in a scene. This, together with
662 the 3D image cube, makes the hyperspectral image hard to display and understand.
663 Hyperspectral unmixing is the process of estimating a common set of spectral bases
664 (called endmembers) and their corresponding composite percentages (called abun-
665 dance) at each pixel so that people can better visualize, analyze and understand the
666 hyperspectral image.

667 In this section, we apply Algorithm 3.1 with Algorithm 2.1 to the constrained
668 sparse nonnegative matrix factorization (NMF) used in hyperspectral unmixing. The
669 mathematical model is as follows.

$$670 \quad (4.1) \quad \min_{W, H} \quad \frac{1}{2} \|V - WH\|_F^2 + \tau \|H\|_p^p$$

$$671 \quad (4.2) \quad \text{s.t. } W \geq 0, \quad H \geq 0,$$

$$672 \quad (4.3) \quad 1_K^T H = 1_N^T,$$

673 where $V = [v_1, v_2, \dots, v_N] \in R_+^{L \times N}$ is the given hyperspectral image data with L
674 channels and N pixels, $W = [w_1, w_2, \dots, w_K] \in R_+^{L \times K}$ is the endmember matrix
675 including K endmember vectors with $K \ll \min\{L, N\}$, and $H = [h_1, h_2, \dots, h_N] \in$
676 $R_+^{K \times N}$ is the corresponding abundance matrix. Here 1_K and 1_N are the column
677 vectors of all ones of dimension K and N , respectively.

678 In the objective function in (4.1), the parameter $\tau > 0$ balances the data fidelity
679 term $\frac{1}{2} \|V - WH\|_F^2$ and the sparse regularization term $\|H\|_p^p$, $p \in (0, 1)$ that forces
680 the sparsity of the abundance matrix. The sparse regularization term is effective
681 for spectral unmixing since only a few endmembers can contribute to representing
682 an observed pixel. To be physically meaningful, the nonnegative constraints in (4.2)
683 are necessary. Moreover, the abundance sum-to-one constraints (ASC) in (4.3) are
684 required since each column of H is the abundance vector whose components are the
685 proportions of each endmember contributing to the mixed pixel. Let H_{ij} denote the
686 (i, j) -entry of the matrix H . The existence of ASC makes the usually used sparsity-
687 induced regularization term $\|H\|_1 = \sum_{i,j} |H_{ij}|$ meaningless since in this case $\|H\|_1$
688 equals a constant N .

689 To solve the constrained sparse NMF model, the two block coordinate descent
690 method is adopted. That is, W and H are considered to be two separate block
691 variables, and the scheme alternatively solves the two subproblems of matrix-based

692 optimization problems. The difficulty of solving problem (4.1)-(4.3) for block H lies
 693 in two aspects: the non-Lipschitz regularization term of the objective function in (4.1)
 694 and the numerous N constraints defined by ASC in (4.3).

695 In [30], Qian et al. considered the special case $p = \frac{1}{2}$ and called the model $L_{1/2}$ -
 696 NMF. To deal with the ASC, Qian et al. adopted the strategy akin to that in [20] by
 697 augmenting the data matrix V and the endmember matrix W to V_a and W_a as

$$698 \quad (4.4) \quad V_a = \begin{pmatrix} V \\ \delta \mathbf{1}_N^T \end{pmatrix} \quad \text{and} \quad W_a = \begin{pmatrix} W \\ \delta \mathbf{1}_K^T \end{pmatrix},$$

699 where $\delta > 0$ controls the impact of the additivity constraint over the endmember
 700 abundances. This strategy, in fact, leads to solve the penalized counterpart

$$701 \quad (4.5) \quad \min_{W \geq 0, H \geq 0} \frac{1}{2} \|V - WH\|_F^2 + \tau \|H\|_p^p + \frac{1}{2} \delta^2 \|\mathbf{1}_K^T H - \mathbf{1}_N^T\|_F^2.$$

702 The multiplicative update (MU) method [23] for classic NMF is extended to solve
 703 the $L_{1/2}$ -NMF, by alternatively updating W and H as

$$704 \quad (4.6) \quad W \leftarrow W * (VH^T) ./ (WHH^T),$$

$$705 \quad (4.7) \quad H \leftarrow H * (W_a^T V_a) ./ (W_a^T W_a H + \frac{\tau}{2} T_\xi(H)^{-\frac{1}{2}}),$$

706 where $(T_\xi(H)^{-\frac{1}{2}})_{ij} = H_{ij}^{-\frac{1}{2}}$ if $H_{ij} > \xi$ and $(T_\xi(H)^{-\frac{1}{2}})_{ij} = 0$ otherwise for a predefined
 707 threshold $\xi > 0$ to avoid computationally instability. Here “ $*$ ” and “ $./$ ” denote the
 708 elementwise matrix multiplication and division, respectively.

709 Here we use the two block proximal alternating optimization (PAO) framework
 710 to solve (4.5). Let W_a^k be the augmented matrix in (4.4) where the block W in W_a is
 711 replaced by W^k .

Algorithm 4.1 PAO Framework

- 1: Initialize $W^1 \geq 0$, $H^1 \geq 0$, and parameters $\tau_1 > 0$ and $\tau_2 > 0$.
- 2: Repeat until a stopping criterion is satisfied
 - 2.1 Find W^{k+1} and H^{k+1} such that

$$(4.8) \quad W^{k+1} = \arg \min_{W \geq 0} \left\{ \frac{1}{2} \|V - WH^k\|_F^2 + \frac{1}{2} \tau_1 \|W - W^k\|_F^2 \right\},$$

$$(4.9) \quad H^{k+1} = \arg \min_{H \geq 0} \left\{ \frac{1}{2} \|V_a - W_a^k H\|_F^2 + \tau \|H\|_p^p + \frac{1}{2} \tau_2 \|H - H^k\|_F^2 \right\}.$$

2.2 Set $k := k + 1$.

712 We combine Algorithm 2.1 and Algorithm 3.1 proposed in this paper to solve the
 713 two subproblems (4.8) and (4.9) in Algorithm 4.1.

- 714 • To solve the W -subproblem in (4.8), we use ASCG, i.e., Algorithm 2.1 with
 715 the LCO employing the conjugate gradient (CG) method [12].
- 716 • To solve the H -subproblem in (4.9) that involves the non-Lipschitz term,
 717 we use SASCG, i.e., Algorithm 3.1 with ASCG that solves the smoothing
 718 H -subproblem of (4.9). The smoothing function of $\|H\|_p^p$ is constructed by
 719 using (3.16).

720 We denote the method as PAO-ASCG-SASCG for short.

721 We also use the two block proximal alternating optimization (PAO) framework
 722 to solve (4.1)-(4.3) directly without penalization to the equality constraints, by sub-
 723 stituting (4.9) in Algorithm 4.1 by

$$724 \quad (4.10) \quad H^{k+1} = \arg \min_{H \geq 0, \mathbf{1}_K^T H = \mathbf{1}_N^T} \{F_{W^k, H^k}(H)\},$$

725 where

$$726 \quad (4.11) \quad F_{W^k, H^k}(H) := \frac{1}{2} \|V - W^k H\|_F^2 + \tau \|H\|_p^p + \frac{1}{2} \tau_2 \|H - H^k\|_F^2.$$

727
 728 We then combine Algorithm 2.1 and Algorithm 3.1 proposed in this paper to solve
 729 (4.8) and (4.10) in the PAO framework.

- 730 • To solve the W -subproblem in (4.8), we use the projected gradient method.
- 731 • To solve the H -subproblem in (4.10), we use SASPG, i.e., Algorithm 3.1,
 732 together with Algorithm 2.1 in which the LCO being the projected gradient
 733 method. The smoothing function of $\|H\|_p^p$ is also constructed by using (3.16).

734 We denote the method as PAO-PG-SASPG-O for short. Here ‘-O’ indicates that the
 735 original $L_{1/2}$ -NMF problem (4.1)-(4.3) is solved.

736 It is worth mentioning that the constraints in (4.10) are N independent simplex
 737 $h_j \geq 0, \sum_{i=1}^K H_{ij} = 1, j = 1, 2, \dots, N$. Let

$$738 \quad \mathcal{A}(H^k) := \{(i, j) : H_{ij}^k = 0\},$$

$$739 \quad \check{\Omega}(H^k) := \{H \in \Omega : H_{ij} = 0 \text{ if } (i, j) \in \mathcal{A}(H^k)\}.$$

740
 741 The efficiency of Algorithm 2.1 depends on the fast computation of matrices
 742 $P_\Omega[H], P_{\check{\Omega}(H^k)}[H], \nabla_\Omega F_{W^k, H^k}(H)$, and $g^A(H)$. Here $P_{\check{\Omega}(H^k)}[H]$ is used for the pro-
 743 jected gradient method that works on the faces $\check{\Omega}(H^k)$ of Ω . All the four types of
 744 matrices are essentially composed by projections of a vector on a certain polyhedron.
 745 The projections of a vector on a polyhedron can be obtained efficiently, e.g., [18]. Here
 746 we compute them in matrix form directly, since N is in general no less than 10000.
 747 We adopt the Matlab code **SimplexProj** in [34] for obtaining $P_\Omega[H]$. And by using
 748 the grouping idea of inactive indices as in [22], we use **SimplexProj** for computing
 749 $P_{\check{\Omega}(H^k)}[H]$ on each group with the same inactive constraints. Moreover, the projected
 750 gradient $\nabla_\Omega F_{W^k, H^k}(H)$, and $g^A(H)$ can be computed efficiently in matrix form using
 751 the KKT conditions.

752 We use two real-world data in the experiment.

753 **Jasper Ridge**, is a popular hyperspectral data. There are 512×614 pixels in
 754 it. In this image, each pixel is recorded at 224 channels ranging from 0.38 to $2.5\mu m$,
 755 and the spectral resolution is up to $9.46nm$. Because this hyperspectral image is
 756 too complex to get the groundtruth, we consider a subimage of 100×100 as in [43],
 757 the first pixel of which is the (105, 269)-th pixel in the original image. After the
 758 channels 1–3, 108–112, 154–166 and 220–224 are removed (due to dense water vapor
 759 and atmospheric effects), we remain 198 channels (this is a common preprocess for
 760 hyperspectral unmixing analysis). There are 4 endmembers in groundtruth: #1 Tree,
 761 #2 Soil, #3 Water, #4 Road.

762 **Urban**, is one of the most widely used hyperspectral data in the hyperspectral
 763 unmixing study. There are 307×307 pixels, each of which corresponds to a 2×2

764 m² area. In this image, there are 210 wavelengths ranging from 400 nm to 2500 nm,
 765 resulting in a spectral resolution of 10 nm. After the channels 1–4, 76, 87, 101–111,
 766 136–153 and 198–210 are removed, we remain 162 channels. There are 4 endmembers
 767 in ground truth: #1 Asphalt, #2 Grass, #3 Tree, #4 Roof.

768 We choose $p = \frac{1}{2}$ and consider the $L_{1/2}$ -NMF problem. We compare our methods
 769 (PAO-ASCG-SASCG and PAO-PG-SASPG-O) with the other three methods. The
 770 information of all the methods are summarized as follows.

- 771 1) Our method: PAO-ASCG-SASCG that solves the penalized counterpart of
 772 $L_{1/2}$ -NMF problem in (4.5).
- 773 2) Our method: PAO-PG-SASPG-O that solves the original $L_{1/2}$ -NMF problem
 774 in (4.1)-(4.3).
- 775 3) PAO-PG-SPG-O: this method solves the original $L_{1/2}$ -NMF problem in (4.1)-
 776 (4.3). It employs the PAO framework in Algorithm 4.1 with (4.9) substituted
 777 by (4.10). The W -subproblem is solved by the PG method [24] and the H -
 778 subproblem is solved by the smoothing projected gradient method [41]. No
 779 active set strategy is adopted.
- 780 4) MU method: this method is a state-of-art method that employs (4.6) and
 781 (4.7) recursively to solve the penalized counterpart of $L_{1/2}$ -NMF problem in
 782 (4.5).
- 783 5) Adaptive HT method: this method is proposed in [35]. It employs the half-
 784 thresholding algorithm and an adaptive strategy for automatically choosing
 785 regularization parameters $\tau_j^k, j = 1, 2, \dots, N$ in k th iteration, and solving the
 786 penalized $L_{1/2}$ sparsity-constrained NMF defined by

$$787 \quad (4.12) \quad \min_{W \geq 0, H \geq 0} \frac{1}{2} \|V - WH\|_F^2 + \sum_{j=1}^N \tau_j^k \|h_j\|_{\frac{1}{2}}.$$

788
 789 We set the maximum CPU time to be 3000 seconds for all the methods, and
 790 the maximum number of iterations for the MU method to be 3000, and the max-
 791 imum number of iterations for the PAO-ASCG-SASCG, PAO-PG-SASPG-O, and
 792 PAO-PG-SPG-O methods to be 1000, and $n_1 = 5$ in Algorithm 3.1. To overcome
 793 the nonconvexity of the original problem (4.1)-(4.3), and the penalized problem (4.5),
 794 we randomly choose 10 initial points for W^1 and H^1 using the Matlab commands
 795 $\text{rand}(L, K)$ and $\text{rand}(K, N)$ for all the methods, respectively. And each column of
 796 H^1 is further rescaled to be sum to one, according to the ASC in (4.3). The MU and
 797 the PAO-ASCG-SASCG methods involve two essential parameters τ and δ , while the
 798 Adaptive HT method only has one parameter δ , and the PAO-PG-SASCG-O meth-
 799 ods only has one parameter τ . In order to estimate an optimal parameter, we first
 800 determine the intervals $[\tau_{\min}, \tau_{\max}]$, and/or $[\delta_{\min}, \delta_{\max}]$ by trying the values at large
 801 steps. We then search the optimal parameters by trying more values in the interval
 802 $[\tau_{\min}, \tau_{\max}]$, and/or $[\delta_{\min}, \delta_{\max}]$.

803 If (W, H) is a solution of NMF, then $(WD, D^{-1}H)$ is also a solution of NMF for
 804 any positive diagonal matrices D . To get rid of this kind of uncertainty, one intuitive
 805 method is to scale each column of W to be the unit ℓ_1 - or ℓ_2 -norm [39, 43], e.g.,

$$806 \quad (4.13) \quad W_{lk} \leftarrow \frac{W_{lk}}{\sqrt{\sum W_{lk}^2}}, \quad H_{kn} \leftarrow H_{kn} \sqrt{\sum W_{lk}^2}.$$

807 Considering the ASC in (4.3), we further let

$$808 \quad (4.14) \quad H_{kn} \leftarrow \frac{H_{kn}}{\sum_k H_{kn}}.$$

809 To evaluate the performance of the computed solution, we use the spectral angle
810 distance (SAD) and the root mean squared error (RMSE) [30, 35, 43] as two
811 benchmark metrics. The SAD is used to evaluate the endmembers, which is defined
812 as

$$813 \quad (4.15) \quad \text{SAD}(w, \hat{w}) = \arccos \left(\frac{w^T \hat{w}}{\|w\| \|\hat{w}\|} \right),$$

814 where w is an estimated endmember, and \hat{w} is the corresponding ground-truth end-
815 member. The RMSE is used to evaluate the performance of the estimated abundance,
816 which is given by

$$817 \quad (4.16) \quad \text{RMSE}(z, \hat{z}) = \left(\frac{1}{N} \|z - \hat{z}\|^2 \right)^{1/2},$$

818 where N is the number of pixels in the image, z is the estimated abundance map
819 (a row vector in the abundance matrix H), and \hat{z} is the corresponding ground-truth
820 abundance map. In general, a smaller SAD and a smaller RMSE correspond to a
821 better hyperspectral unmixing result.

822 We draw in Fig. 1 the corresponding objective value $\frac{1}{2} \|V - WH\|_F^2 + \tau \|H\|_{1,2}^{\frac{1}{2}}$
823 of each iterate point versus the CPU time obtained by the PAO-PG-SASPG-O
824 and the PAO-PG-SPG-O method, using the same optimal parameter $\tau = 1.5 \times 10^6$, and
825 the same initial point on Jasper Ridge data, respectively. We divide the x -axis to be
826 $[0, 200]$ and $[200, 3000]$ in two subfigures to see clear the decrease tendency and the
827 final objective value. We can find from Fig. 1 that our PAO-PG-SASPG-O decreases
828 faster and gets lower objective value than the PAO-PG-SPG-O method. The final
829 objective value obtained by the PAO-PG-SASPG-O method is 2.6494e10, which is
830 much lower than 2.6988e10 that obtained by the PAO-PG-SPG-O method. It is easy
831 to see that the active set strategy helps fasten the computational speed.

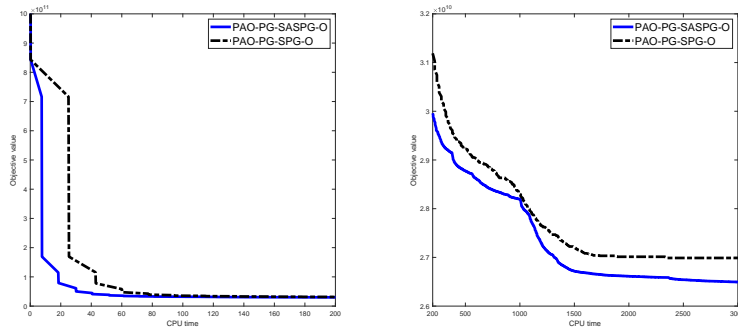


FIG. 1. Convergence curve of objective value versus CPU time using the PAO-PG-SPG-O and the PAO-PG-SASPG-O on the Jasper Ridge data, respectively.

832 For Jasper Ridge, we record in Table 1 the final SAD and RMSE for each end-
833 member corresponding to the computed solution with the smallest sum of SAD and

834 RMSE, among the 10 trials of initial points as well as the choices of parameters. The
 835 lowest SAD and RMSE for each endmember, and the lowest average SAD and RMSE
 836 are indicated in bold face in Table 1. It is easy to see that the computed solution
 837 obtained by the PAO-PG-SASPG-O method proposed in this paper has the lowest
 838 average SAD and RMSE of the four endmembers. Our proposed PAO-ASCG-SASCG
 839 method that solves the penalized version of $L_{1/2}$ -NMF also provides lower average
 840 SAD and RMSE than the MU and the Adaptive HT methods.

841 For Urban, we record in Table 2 the final SAD and the final RMSE for each
 842 endmember. The lowest SAD and RMSE for each endmember, and the lowest average
 843 SAD and RMSE are indicated in bold face in Table 2. Clearly the PAO-ASCG-SASCG
 844 method provides the solution that obtains the lowest average SAD and RMSE than
 845 the other four methods. The PAO-PG-SPG-O and PAO-ASCG-SASCG-O method
 846 for solving the original model (4.1)-(4.3) do not provide satisfying SAD and RMSE.
 847 The reason, we think, is due to the model itself. As pointed out in [43], applying an
 848 identical strength of constraints to all the factors, (that is, in our case, using the same
 849 $p = \frac{1}{2}$ for all the columns of H) does not hold in practice. Therefore, in [43] they
 850 proposed to solve

$$851 \quad (4.17) \quad \min_{W \geq 0, H \geq 0} \frac{1}{2} \|V - WH\|^2 + \tau \sum_{j=1}^N \|h_j\|_{p_j}^{p_j},$$

852 where $p_j \in (0, 1)$, $j = 1, 2, \dots, N$, are estimated from the original data V using
 853 two-steps procedures. If the pixels indeed have very different levels of sparsity as in
 854 Urban, the sum-to-one constraints will make the original model (4.1)-(4.3) deviate a
 855 lot from the true model. The PAO-ASCG-SASCG method, in contrast, because of
 856 the lack of the sum-to-one constraints, has the ability to adjust the sparsity levels of
 857 different pixels to some degree. The Adaptive HT method, which adaptively adjusts
 858 the different regularization parameter for each column of H , also has the effect to
 859 assign different level of sparsity for each pixel. When the pixels have not so much
 860 different levels of sparsity as in Jasper, the PAO-PG-SASPG-O that solves the original
 861 model (4.1)-(4.3) with the sum-to-one constraints provides the best SAD and RMSE.
 862

TABLE 1
 SAD and RMSE on the Jasper Ridge data estimated by our methods and the other methods

	SAD				Avg.
Jasper Ridge ($K = 4$)	#1	#2	#3	#4	#1 ~ #4
MU	0.2070	0.1185	0.3324	0.2939	0.2379
Adaptive HT	0.1451	0.3099	0.1367	0.1515	0.1858
PAO-ASCG-SASCG	0.1241	0.0690	0.1859	0.1645	0.1359
PAO-PG-SPG-O	0.1315	0.0606	0.1132	0.0516	0.0892
PAO-PG-SASPG-O	0.1301	0.0616	0.1019	0.0609	0.0886
	RMSE				Avg.
MU	0.1235	0.0953	0.1773	0.0953	0.1361
Adaptive HT	0.1016	0.1483	0.1761	0.1885	0.1536
PAO-ASCG-SASCG	0.0836	0.0425	0.1244	0.1052	0.0889
PAO-PG-SPG-O	0.0846	0.0581	0.0929	0.0875	0.0808
PAO-PG-SASPG-O	0.0840	0.0578	0.0930	0.0842	0.0798

863 The abundance fractions for Jasper Ridge from the ground-truth, and separated
 864 by the five methods are shown in Fig. 2. We can also see that our proposed PAO-

TABLE 2

SAD and RMSE on the Urban data estimated by our methods and the other methods

Urban ($K = 4$)	SAD				Avg.
	#1	#2	#3	#4	#1 ~ #4
MU	0.1976	0.0318	0.0454	0.1445	0.1048
Adaptive HT	0.0715	0.0393	0.0704	0.3288	0.1275
PAO-ASCG-SASCG	0.0738	0.0525	0.0314	0.0736	0.0578
PAO-PG-SPG-O	0.0900	0.1940	0.0423	0.3424	0.1672
PAO-PG-SASPG-O	0.0925	0.1026	0.0397	0.2153	0.1125
Urban ($K = 4$)	RMSE				Avg.
	#1	#2	#3	#4	#1 ~ #4
MU	0.0989	0.1037	0.0707	0.0995	0.0932
Adaptive HT	0.1165	0.0964	0.0794	0.0895	0.0954
PAO-ASCG-SASCG	0.1101	0.1085	0.0562	0.0548	0.0824
PAO-PG-SPG-O	0.2595	0.2242	0.1281	0.2052	0.2242
PAO-PG-SASPG-O	0.2452	0.1715	0.1435	0.2082	0.1921

865 ASCG-SASCG and PAO-PG-SASPG-O methods provide good estimates of abun-
866 dance. The abundance fractions for Urban from the ground-truth, and separated by
867 the MU, the Adaptive HT, and the PAO-ASCG-SASCG methods are shown in Fig.
868 3. It is easy to see that our proposed PAO-ASCG-SASCG method provide the best
869 estimates of abundance.

870 The numerical results demonstrate that our proposed PAO-PG-SASPG-O method
871 and PAO-ASCG-SASCG method can efficiently solve the original and penalized $L_{1/2}$ -
872 NMF problem, respectively. Moreover, at least one of our methods provides an excel-
873 lent unmixing performance, compared to the popular MU method and the Adaptive
874 HT method.

875 It is worth pointing out that our smoothing active set method can deal with the
876 sum-to-one constraints, but the MU method and the Adaptive HT method can not.
877 Our smoothing active set method is flexible to solve the new model in (4.17) with
878 additional sum-to-one constraints. It is interesting to further investigate how to get
879 good estimation of p_j , $j = 1, 2, \dots, N$, and whether applying our smoothing active
880 set method to this new model can provide even better unmixing results in future.

881 **5. Conclusion remarks.** We develop Algorithm 3.1, a novel smoothing active
882 set method, for solving problem (1.1) where the objective function f may be non-
883 Lipschitz continuous. We approximate f by a continuously differentiable function \tilde{f}
884 and employ Algorithm 2.1 for solving the smooth optimization problem (3.5) until the
885 special updating rule holds in the inner loop of Algorithm 3.1. Algorithm 2.1 is a new
886 active set method for linearly constrained smooth optimization, which ensures that for
887 any positive smoothing parameter μ_k , the iterate x^{k+1} satisfies $\|\nabla_{\Omega}\tilde{f}(x^{k+1}, \mu_k)\| \leq$
888 $\hat{\gamma}\mu_k$. This property is essential for the convergence result of Algorithm 3.1. It is
889 worth noting that convergence results of most existing active set methods for the
890 smooth minimization problem (2.1) are in the sense $\liminf_{k \rightarrow \infty} P_{\Omega}[x^k - \nabla f(x^k)] -$
891 $x^k = 0$, which does not imply $\liminf_{k \rightarrow \infty} \|\nabla_{\Omega}f(x^k)\| = 0$. See inequality (2.35)
892 and Example 1. Our global convergence result, as well as the nice finite identification
893 property, and the local convergence result makes Algorithm 2.1 not only important for
894 approximately solving subproblems in Algorithm 3.1 for non-Lipschitz minimization
895 problem (1.1), but also advanced for smooth problem (2.1).

896 **Acknowledgements.** We are very grateful to Prof. W. W. Hager and the
897 anonymous referees for valuable comments.

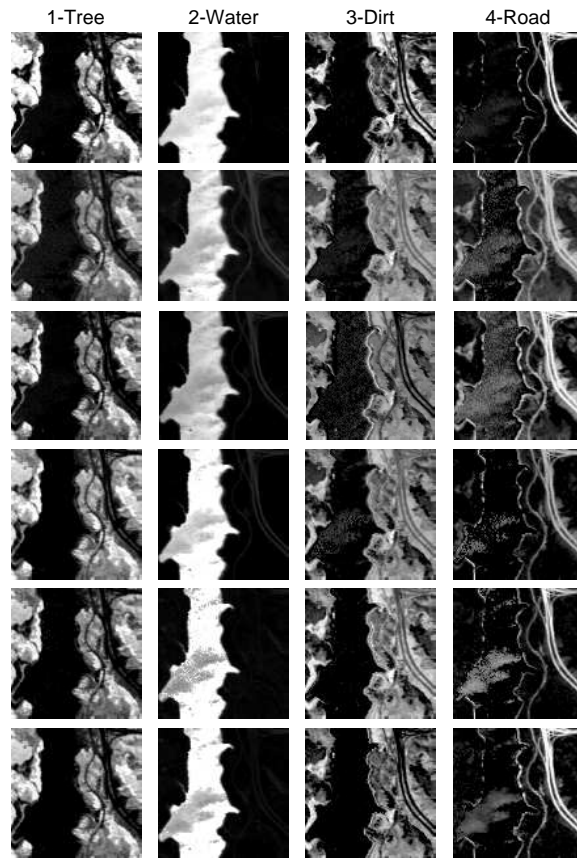


FIG. 2. *Abundance maps from the ground-truth, MU, Adaptive HT, PAO-ASCG-SASCG, PAO-PG-SPG-O, and PAO-PG-SASPG-O (from the first row to the last row sequentially) for four targets in the Jasper Ridge data.*

898

REFERENCES

- 899 [1] N. S. AYBAT AND G. IYENGAR, *A first-order smoothed penalty method for compressed sensing*,
900 SIAM J. Optim., 21 (2011), pp. 287–313.
- 901 [2] M. S. BARZARAA, H. D. SHERALI AND C. M. SHETTY, *Nonlinear Programming: Theory and*
902 *Algorithms*, 2nd ed., John Wiley & Sons, New York, 1993.
- 903 [3] D. P. BERTSEKAS, *Nonlinear Programming*, 2nd ed., Athena Scientific, Belmont, MA, 1999.
- 904 [4] W. BIAN AND X. CHEN, *Linearly constrained non-Lipschitz optimization for image restoration*,
905 SIAM J. Imaging Sci., 8 (2015), pp. 2294–2322.
- 906 [5] W. BIAN AND X. CHEN, *Optimality and complexity for constrained optimization problems with*
907 *nonconvex regularization*, Math. Oper. Res., 42 (2017), pp. 1063–1084.
- 908 [6] J. BOLTE, A. DANILIDIS, A. LEWIS, AND M. SHIOTA, *Clarke subgradients of stratifiable functions*,
909 SIAM J. Optim., 18 (2007), pp. 556–572.
- 910 [7] J. V. BURKE, T. HOHEISEL, AND C. KANZOW, *Gradient consistency for integral-convolution*
911 *smoothing*, Set-Valued Var. Anal., 21 (2013), pp. 359–376.
- 912 [8] P. H. CALAMAI AND J. J. MORÉ, *Projected gradient method for linearly constrained problems*,
913 Math. Program., 39 (1987), pp. 93–116.
- 914 [9] X. CHEN, L. GUO, Z. LU, AND J. J. YE, *An augmented Lagrangian method for non-Lipschitz*

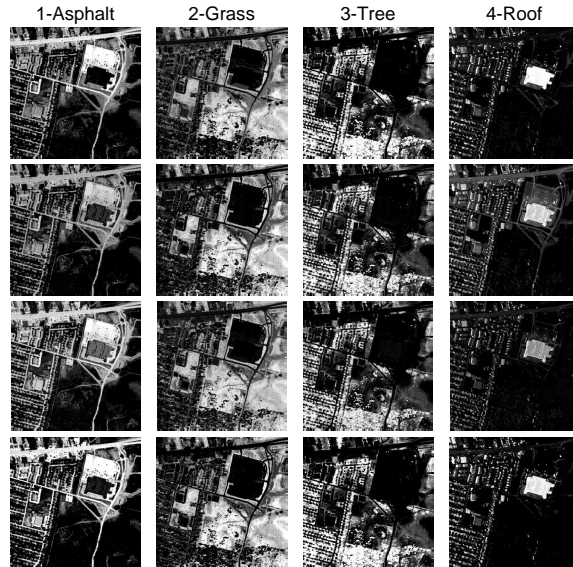


FIG. 3. *Abundance maps from the ground-truth, MU, Adaptive HT and PAO-ASCG (from the first row to the last row sequentially) for four targets in the Urban data.*

- 915 *nonconvex programming*, SIAM J. Numer. Anal., 55 (2017), pp. 168–193.
- 916 [10] X. CHEN, *Smoothing methods for nonsmooth, nonconvex minimization*, Math. Program., 134
917 (2012), pp. 71–99.
- 918 [11] X. CHEN, L. NIU, AND Y. YUAN, *Optimality conditions and smoothing trust region Newton
919 method for non-Lipschitz optimization*, SIAM J. Optim., 23 (2013), pp. 1528–1552.
- 920 [12] X. CHEN AND W. ZHOU, *Smoothing nonlinear conjugate gradient method for image restoration
921 using nonsmooth nonconvex minimization*, SIAM J. Imaging Sci., 3 (2010), pp. 765–790.
- 922 [13] Z. DOSTÁL AND T. KOZUBEK, *An optimal algorithm with superrelaxation for minimization of
923 a quadratic functions subject to separable constraints with applications*, Math. Program.,
924 135 (2012), pp. 195–220.
- 925 [14] F. FACCHINEI, A. FISCHER, AND C. KANZOW, *On the accurate identification of active con-
926 straints*, SIAM J. Optim., 9 (1998), pp. 14–32.
- 927 [15] B. FISTRICH, S. PATERLINI, AND P. WINKER, *Cardinality versus q-norm constraints for index
928 tracking*, Quant. Finance, 14 (2014), pp. 2019–2032.
- 929 [16] M. FUKUSHIMA, *A modified Frank-Wolfe algorithm for solving the traffic assignment problem*,
930 Transport Res. Part B: Meth., 18 (1984), pp. 169–177.
- 931 [17] W. W. HAGER AND H. ZHANG, *A new active set algorithm for box constrained optimization*,
932 SIAM J. Optim., 17 (2006), pp. 526–557.
- 933 [18] W. W. HAGER AND H. ZHANG, *An active set algorithm for nonlinear optimization with poly-
934 hedral constraints*, Sci. China Math., 59 (2016), pp. 1525–1542.
- 935 [19] W. W. HAGER AND H. ZHANG, *An affine scaling method for optimization problems with poly-
936 hedral constraints*, Comput. Optim. Appl., 59 (2014), pp. 163–183.
- 937 [20] D. HEINZ AND C.-I. CHANG, *Fully constrained least squares linear spectral mixture analysis
938 method for material quantification in hyperspectral imagery*, IEEE Trans. Geosci. Remote
939 Sens., 39 (2001), pp. 529–545.
- 940 [21] N. KESKAR AND A. WÄCHTER, *A limited-memory quasi-Newton algorithm for bound-
941 constrained nonsmooth optimization*, Optim. Method Softw., 34 (2019), pp. 150–171.
- 942 [22] J. KIM AND H. PARK, *Fast nonnegative matrix factorization: An active-set-like method and
943 comparisons*, SIAM J. Sci. Comput., 33 (2011), pp. 3261–3281.
- 944 [23] D. LEE AND H. SEUNG, *Algorithms for non-negative matrix factorization*, Adv. Neural Inf.
945 Process. Syst., 13 (2001), pp. 556–562.
- 946 [24] C.-J. LIN, *Projected gradient methods for nonnegative matrix factorization*, Neural Comput.,
947 19 (2007), pp. 2756–2779.

- 948 [25] C.-J. LIN AND J. J. MORE, *Newton's method for large bound-constrained optimization problems*,
 949 SIAM J. Optim., 9 (1999), pp. 1100–1127.
- 950 [26] Y. LIU, S. MA, Y. H. DAI, S. ZHANG, *A smoothing sequential quadratic programming (SSQP)*
 951 *framework for a class of composite L_q minimization over polyhedron*, Math. Program., 158
 952 (2016), pp. 467–500.
- 953 [27] YU. NESTEROV, *Smooth minimization of non-smooth functions*, Math. Program., 103 (2005),
 954 pp. 127–152.
- 955 [28] M. NIKOLOVA, *Minimizers of cost-functions involving nonsmooth data-fidelity terms. Applica-*
 956 *tion to the processing of outliers*, SIAM J. Numer. Anal., 40 (2002), pp. 965–994.
- 957 [29] E. R. PANIER, *An active set method for solving linearly constrained nonsmooth optimization*
 958 *problems*, Math. Program., 37 (1987), pp. 269–292.
- 959 [30] Y. QIAN, S. JIA, J. ZHOU, AND A. ROBLES-KELLY, *Hyperspectral unmixing via $L_{1/2}$ sparsity-*
 960 *constrained nonnegative matrix factorization*, IEEE Trans. Geosci. Remote Sens., 49
 961 (2011), pp. 4282–4297.
- 962 [31] R. T. ROCKAFELLAR AND R. J-B. WETS, *Variational Analysis*, Springer, Germany, 1998.
- 963 [32] A. SHAPIRO, D. DENTCHEVA, AND A. RUSZCZYŃSKI, *Lectures on Stochastic Programming: Mod-*
 964 *eling and Theory*, Society for Industrial and Applied Mathematics, Philadelphia, 2009.
- 965 [33] P. TSENG, I. M. BOMZE, AND W. SCHACHINGER, *A first-order interior-point for linearly con-*
 966 *strained smooth optimization*, Math. Program., 127 (2011), pp. 399–424.
- 967 [34] W. WANG AND M. A. CARREIRA-PERPINÀN, *Projection onto the probability simplex: An effi-*
 968 *cient algorithm with a simple proof, and an application*, arXiv preprint arXiv: 1309.1541,
 969 2013.
- 970 [35] W. WANG AND Y. QIAN, *Adaptive $L_{1/2}$ sparsity-constrained NMF with half-thresholding algo-*
 971 *rithm for hyperspectral unmixing*, IEEE J. Sel Topics Appl. Earth Observ. Remote Sens.,
 972 8 (2015), pp. 2618–2631.
- 973 [36] Z. WEN, W. YIN, D. GOLDFARB, AND Y. ZHANG, *A fast algorithm for sparse reconstruction*
 974 *based on shrinkage subspace optimization and continuation*, SIAM J. Comput., 32 (2010),
 975 pp. 1832–1857.
- 976 [37] Z. WEN, W. YIN, H. ZHANG, AND D. GOLDFARB, *On the convergence of an active set method*
 977 *for ℓ_1 minimization*, Optim. Method Softw., 27 (2012), pp. 1127–1146.
- 978 [38] M. XU, J. J. YE, AND L. ZHANG, *Smoothing SQP methods for solving degenerate nonsmooth*
 979 *constrained optimization problems with applications to bilevel programs*, SIAM J. Optim.,
 980 25 (2015), pp. 1388–1410.
- 981 [39] W. XU, X. LIU, AND Y. GONG, *Document clustering based on nonnegative matrix factorization*,
 982 in Proc. 26th Int. Conf. Res. Develop. Inf. Retr. (SIGIR), (2003), pp. 267–273.
- 983 [40] Z. YANG, G. ZHOU, S. XIE, S. DING, J. M. YANG, AND J. ZHANG, *Blind spectral unmixing*
 984 *based on sparse nonnegative matrix factorization*, IEEE Trans. Image Process., 20 (2011),
 985 pp. 1112–1125.
- 986 [41] C. ZHANG AND X. CHEN, *Smoothing projected gradient method and its application to stochastic*
 987 *linear complementarity problems*, SIAM J. Optim., 20 (2009), pp. 627–649.
- 988 [42] C. ZHANG, L. JING, AND N. XIU, *A new active set method for nonnegative matrix factorization*,
 989 SIAM J. Sci. Comput., 36 (2014), pp. A2633–A2653.
- 990 [43] F. ZHU, Y. WANG, B. FAN, S. XIANG, AND G. MENG, *Spectral unmixing via data-guided sparsity*,
 991 IEEE Trans. Image Process., 23 (2014), pp. 5412–5427.