

Continuous exact relaxation and alternating proximal gradient algorithm for partial sparse and partial group sparse optimization problems

Qingqing Wu, Dingtao Peng, Xian Zhang

Abstract In this paper, we consider a partial sparse and partial group sparse optimization problem, where the loss function is a continuously differentiable function (possibly nonconvex), and the penalty term consists of two parts associated with sparsity and group sparsity. The first part is the ℓ_0 norm of \mathbf{x} , the second part is the $\ell_{2,0}$ norm of \mathbf{y} , i.e. $\lambda_1 \|\mathbf{x}\|_0 + \lambda_2 \|\mathbf{y}\|_{2,0}$, where $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{n+m}$ is the decision variable. We give a continuous relaxation model of the above original problem, where the two parts of the penalty term are relaxed by Capped- ℓ_1 of \mathbf{x} and group Capped- ℓ_1 of \mathbf{y} respectively. Firstly, we define two kinds of stationary points of the relaxation model. Based on the lower bound property of d-stationary points of the relaxation model, we establish the equivalence of solutions of the original problem and the relaxation model, which provides a theoretical basis for solving the original problem via solving the relaxation problem. Secondly, we propose an alternating proximal gradient (APG) algorithm to solve the relaxation model, and prove that the whole sequence of the APG algorithm converges to a critical point under some mild conditions. Finally, numerical experiments on simulated data and multichannel image as well as comparison with some state-of-art algorithms are presented to illustrate the effectiveness and robustness of the proposed algorithm for partial sparse and partial group sparse optimization problem.

Keywords Partial sparse and partial group sparse optimization problem; continuous exact relaxation; stationary point; alternating proximal gradient algorithm; whole sequence convergence

MSC(2010) 90C26 · 90C46

1 Introduction

In the past decade, sparse optimization problems have attracted great attention in variable selection, image restoration, gene expression, and so on [5, 10, 14, 20, 21, 22, 39, 45]. The basic framework of sparse optimization problem is to seek a sparse solution of an underde-

Dingtao Peng ✉

School of Mathematics and Statistics, Guizhou University, Guiyang 550025, China. E-mail: dingtaopeng@126.com

Qingqing Wu

School of Mathematics and Statistics, Guizhou University, Guiyang 550025, China. E-mail: gs.qqwu21@gzu.edu.cn

Xian Zhang

School of Mathematics and Statistics, Guizhou University, Guiyang 550025, China. E-mail: zhangxian05@163.com

30 terminated linear system. The general sparse optimization problem is as follows:

$$\min_{\mathbf{x} \in \mathbb{R}^n} F(\mathbf{x}) = f(\mathbf{x}) + \lambda \|\mathbf{x}\|_0,$$

31 where $f : \mathbb{R}^n \rightarrow \mathbb{R}_+$ is a loss function, $\lambda > 0$, $\|\mathbf{x}\|_0 := \sum_{i=1, x_i \neq 0}^n |x_i|^0$. A vector $\mathbf{x} \in \mathbb{R}^n$ is said
32 to be sparse if $\|\mathbf{x}\|_0 \ll n$, and the sparsity of vector $\mathbf{x} \in \mathbb{R}^n$ is usually provided by its ℓ_0
33 norm.

34 Due to the fact that traditional sparse optimization problems only consider the sparsity
35 of a single item and do not have sufficient ability to handle complex structures such as group
36 sparse structures, Yuan and Lin [39] first use group sparse structures as prior information.
37 Group sparse structure refers to dividing variables into multiple groups, and then considering
38 whether each group as a whole is zero. Let $\mathbf{x} = (\mathbf{x}_{(1)}^\top, \dots, \mathbf{x}_{(J)}^\top)^\top$ with J disjoint groups, where
39 $\mathbf{x}_{(i)} = (x_{(i)1}, \dots, x_{(i)n_i})^\top \in \mathbb{R}^{n_i}$, $n_i > 0$ and $\sum_{i=1}^J n_i = n$. Then the optimization problem
40 with group sparse structure can be formulated as the following group sparse optimization
41 [24, 30, 31]:

$$\min_{\mathbf{x} \in \mathbb{R}^n} F(\mathbf{x}) = f(\mathbf{x}) + \lambda \|\mathbf{x}\|_{2,0},$$

42 where $\|\mathbf{x}\|_{2,0} := \#\{i \mid \|\mathbf{x}_{(i)}\| \neq 0, i = 1, \dots, J\}$ is called $\ell_{2,0}$ norm that counts the number of
43 nonzero groups of \mathbf{x} , in which $\|\mathbf{x}_{(i)}\|$ denotes the ℓ_2 norm of the subvector $\mathbf{x}_{(i)}$. Note that
44 $\|\cdot\|_{2,0}$ is nonconvex, nonsmooth, and even discontinuous, which causes the above problem
45 to be NP-hard. Many researchers consider the relaxation problem of this problem, such as
46 group LASSO model [35], Bayes group LASSO models [9, 33], group SCAD model [23, 31,
47 38], group MCP model [31, 40] and other models [30, 32, 42, 44, 46].

48 When the data consist of two parts such that the first part has a sparse structure and
49 the second part has a certain group sparse structure, it naturally makes sense for us to
50 investigate the following partial sparse and partial group sparse optimization problem:

$$\min_{\mathbf{x} \in \mathbb{R}^n, \mathbf{y} \in \mathbb{R}^m} F(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}, \mathbf{y}) + \lambda_1 \|\mathbf{x}\|_0 + \lambda_2 \|\mathbf{y}\|_{2,0}, \quad (1.1)$$

51 where $f(\mathbf{x}, \mathbf{y})$ is a loss function which we suppose it to be continuously differentiable but not
52 necessarily convex in this paper. In (1.1), $\lambda_1, \lambda_2 > 0$, $\|\mathbf{x}\|_0 = \sum_{i=0, x_i \neq 0}^n |x_i|^0$ is called ℓ_0 norm
53 of \mathbf{x} , $\mathbf{y} = (\mathbf{y}_{(1)}^\top, \dots, \mathbf{y}_{(J)}^\top)^\top \in \mathbb{R}^m$ with J disjoint groups, and $\|\mathbf{y}\|_{2,0} = \#\{j \mid \|\mathbf{y}_{(j)}\| \neq 0, j =$
54 $1, \dots, J\}$ is called $\ell_{2,0}$ norm of \mathbf{y} . Specially, if \mathbf{x} and \mathbf{y} are same, problem (1.1) degrades to
55 the following sparse plus group sparse optimization problem [26]:

$$\min_{\mathbf{x} \in \mathbb{R}^n} F(\mathbf{x}) = f(\mathbf{x}) + \lambda_1 \|\mathbf{x}\|_0 + \lambda_2 \|\mathbf{x}\|_{2,0}.$$

56 Since both $\|\cdot\|$ and $\|\cdot\|_{2,0}$ are nonconvex, nonsmooth and discontinuous, problem (1.1)
57 in general is NP-hard. One popular way is to relax ℓ_0 ($\ell_{2,0}$) norm to ℓ_1 ($\ell_{2,1}$) norm which
58 are convex [35, 44], but the solution obtained by the relaxation problem is biased and does
59 not satisfy oracle property [16, 17]. Therefore, some researchers propose using several classes
60 of folding concave continuous relaxations which are still nonconvex but have some good
61 properties. These nonconvex relaxations includes ℓ_p ($0 < p < 1$) norm, smoothly clipped
62 absolute deviation (SCAD) penalty [17], minimax concave penalty (MCP) [40], Capped- ℓ_1
63 penalty [28, 41] and their corresponding group structure forms, such as $\ell_{p,q}$ group SCAD

64 and group MCP. The nonconvex relaxations have been widely studied in many works, for
 65 example [3, 4, 11, 30, 36, 37, 42]. It has been proved that the solutions obtained by these kinds
 66 of nonconvex optimization have some desired properties: unbiasedness, sparsity, continuity
 67 and oracle property. Specially, reference [25] has shown that Capped- ℓ_1 relaxation is the
 68 tightest difference-of-convex (DC) relaxation for ℓ_0 norm.

69 In this paper, we consider using Capped- ℓ_1 and group Capped- ℓ_1 to relax ℓ_0 norm and
 70 $\ell_{2,0}$ norm in problem (1.1) respectively, that is, we consider the following problem

$$\min_{\mathbf{x} \in \mathbb{R}^n, \mathbf{y} \in \mathbb{R}^m} F(\mathbf{x}, \mathbf{y}) := f(\mathbf{x}, \mathbf{y}) + \lambda_1 \Phi_1(\mathbf{x}) + \lambda_2 \Phi_2(\mathbf{y}), \quad (1.2)$$

71 where

$$\Phi_1(\mathbf{x}) := \sum_{i=1}^n \varphi_1(|x_i|), \quad \Phi_2(\mathbf{y}) := \sum_{j=1}^J \varphi_2(\|\mathbf{y}_{(j)}\|),$$

72 which are Capped- ℓ_1 regularization and group Capped- ℓ_1 regularization respectively, and

$$\varphi_v(t) := \min \left\{ 1, \frac{t}{\alpha_v} \right\} = \frac{t}{\alpha_v} - \max \left\{ 0, \frac{t}{\alpha_v} - 1 \right\} = \begin{cases} \frac{t}{\alpha_v}, & \text{if } 0 \leq t < \alpha_v, \\ 1, & \text{if } t \geq \alpha_v, \end{cases}$$

73 with $\alpha_v > 0$, $v = 1, 2$. The penalty function $\varphi_v : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ can be written in the form of DC
 74 form as $\varphi_v(t) := g_v(t) - h_v(t)$ with $g_v(t) = \frac{t}{\alpha_v}$, $h_v(t) = \max\{0, \frac{t}{\alpha_v} - 1\}$. Therefore, problem
 75 (1.2) can be rewritten as follows:

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^n, \mathbf{y} \in \mathbb{R}^m} F(\mathbf{x}, \mathbf{y}) = & f(\mathbf{x}, \mathbf{y}) + \lambda_1 \sum_{i=1}^n (g_1(|x_i|) - h_1(|x_i|)) \\ & + \lambda_2 \sum_{j=1}^J (g_2(\|\mathbf{y}_{(j)}\|) - h_2(\|\mathbf{y}_{(j)}\|)). \end{aligned} \quad (1.3)$$

76 In recent years, many scholars have studied the relaxation models of sparse or group
 77 sparse optimization problems. For the sparse optimization problem with the linear least
 78 square loss and ℓ_p regularization, the reference [12] established the lower bound property
 79 of nonzero entires of local solutions. When the loss function is convex and the constraint
 80 set is a box, the reference [3] studied the relationship between the original ℓ_0 regularization
 81 problem and the Capped- ℓ_1 relaxation problem. Under certain conditions, the equivalence
 82 of global solutions and the inclusion relationship of local solutions between the two prob-
 83 lems are proved. The authors also proposed a smoothing proximal gradient algorithm for
 84 solving the relaxation problem. The reference [31] considered a class of group sparse opti-
 85 mization problems with nonconvex folding concave continuous relaxations, and researched
 86 the first-order and second-order directional stationary points of the problem. The reference
 87 [30] considered three kinds of group sparse optimization models with linear inequality con-
 88 straints and discussed the relationship between stationary points, local solutions and global
 89 solutions. The reference [42] considered a class of group sparse optimization models with a
 90 general constraint set, and discussed the relationship of local solutions and global solutions
 91 between original problem and relaxation problem.

92 In this paper, inspired by the above works, we study the stationary points of problem
 93 (1.2), the equivalence of solutions between problems (1.1) and (1.2), and provide an efficient
 94 algorithm for solving problem (1.2).

95 This paper is organized as follows. In Section 2, we give some preliminaries that will be
 96 used in this paper. In Section 3, we define two classes of stationary points for the relaxation
 97 model and discuss their characterization, relationship and some properties. In Section 4, we
 98 establish the equivalence of solutions between the original problem (1.1) and the relaxation
 99 model (1.2). In Section 5, we propose an APG algorithm for problem (1.2) and establish the
 100 convergence result of the whole sequence. In Section 6, we test the proposed APG algorithm
 101 through rich numerical experiments on recovering the simulated partial sparse and partial
 102 group sparse signals and some real images. In Section 7, we make a brief conclusion of this
 103 paper.

104 2 Notations and preliminaries

105 In this section, we provide some basic notations, and introduce the preliminaries of
 106 several kinds of stationary points and subdifferentials.

107 **Notations:** For any $n \in \mathbb{N}^+$, $[n] := \{1, \dots, n\}$. For any $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{y} \in \mathbb{R}^m$, $\nabla f(\mathbf{x}, \mathbf{y}) =$
 108 $(\nabla_{\mathbf{x}}^\top f(\mathbf{x}, \mathbf{y}), \nabla_{\mathbf{y}}^\top f(\mathbf{x}, \mathbf{y}))^\top$, where $\nabla_{\mathbf{y}} f(\mathbf{x}, \mathbf{y}) = ([\nabla_{\mathbf{y}}^\top f(\mathbf{x}, \mathbf{y})]_{(1)}, \dots, [\nabla_{\mathbf{y}}^\top f(\mathbf{x}, \mathbf{y})]_{(J)})^\top$, and
 109 $[\nabla_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})]_{(j)} = ([\nabla_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})]_{(j)1}, \dots, [\nabla_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})]_{(j)m_j})^\top$. For convenience, we define the
 110 following index sets

$$\begin{aligned} I_1(\mathbf{x}) &:= \{i : |x_i| = 0, \forall i \in [n]\}, \\ I_2(\mathbf{x}) &:= \{i : 0 < |x_i| < \alpha_1, \forall i \in [n]\}, \\ I_3(\mathbf{x}) &:= \{i : |x_i| = \alpha_1, \forall i \in [n]\}, \\ I_4(\mathbf{x}) &:= \{i : |x_i| > \alpha_1, \forall i \in [n]\}, \\ J_1(\mathbf{y}) &:= \{j : \|\mathbf{y}_{(j)}\| = 0, \forall j \in [J]\}, \\ J_2(\mathbf{y}) &:= \{j : 0 < \|\mathbf{y}_{(j)}\| < \alpha_2, \forall j \in [J]\}, \\ J_3(\mathbf{y}) &:= \{j : \|\mathbf{y}_{(j)}\| = \alpha_2, \forall j \in [J]\}, \\ J_4(\mathbf{y}) &:= \{j : \|\mathbf{y}_{(j)}\| > \alpha_2, \forall j \in [J]\}. \end{aligned}$$

Let $I(\mathbf{x}) := I_2(\mathbf{x}) \cup I_3(\mathbf{x}) \cup I_4(\mathbf{x})$, and $J(\mathbf{y}) := J_2(\mathbf{y}) \cup J_3(\mathbf{y}) \cup J_4(\mathbf{y})$. Denote

$$\ell(x_i) := |x_i|, \quad \rho_j(\mathbf{y}_{(j)}) := \|\mathbf{y}_{(j)}\|,$$

111 then problem (1.3) can be rewritten as follows

$$\begin{aligned} \min_{\mathbf{x} \in \mathbb{R}^n, \mathbf{y} \in \mathbb{R}^m} F(\mathbf{x}, \mathbf{y}) &= f(\mathbf{x}, \mathbf{y}) + \lambda_1 \sum_{i=1}^n (g_1 \circ \ell(x_i) - h_1 \circ \ell(x_i)) \\ &\quad + \lambda_2 \sum_{j=1}^J (g_2 \circ \rho_j(\mathbf{y}_{(j)}) - h_2 \circ \rho_j(\mathbf{y}_{(j)})), \end{aligned} \tag{2.1}$$

112 where "o" denotes the composition of two functions.

113 Next, we introduce several important concepts to characterize optimal conditions of
 114 problem (1.2).

115 **Definition 2.1** [13, 31] Let $h : \mathbb{R}^{n+m} \rightarrow \mathbb{R} \cup \{\infty\}$, for any $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{n+m}$, the directional
 116 derivative of h at $(\hat{\mathbf{x}}, \hat{\mathbf{y}}) \in \mathbb{R}^{n+m}$ is defined as

$$h'((\hat{\mathbf{x}}, \hat{\mathbf{y}}); (\mathbf{x} - \hat{\mathbf{x}}, \mathbf{y} - \hat{\mathbf{y}})) := \lim_{t \downarrow 0} \frac{h((\hat{\mathbf{x}}, \hat{\mathbf{y}}) + t(\mathbf{x} - \hat{\mathbf{x}}, \mathbf{y} - \hat{\mathbf{y}})) - h(\hat{\mathbf{x}}, \hat{\mathbf{y}})}{t}.$$

117 If h is differentiable at $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$, then $h'((\hat{\mathbf{x}}, \hat{\mathbf{y}}); (\mathbf{x} - \hat{\mathbf{x}}, \mathbf{y} - \hat{\mathbf{y}})) = \langle \nabla h(\hat{\mathbf{x}}, \hat{\mathbf{y}}), (\mathbf{x} - \hat{\mathbf{x}}, \mathbf{y} - \hat{\mathbf{y}}) \rangle$.

118 By the definition, for any $(\mathbf{x}, \mathbf{y}), (\hat{\mathbf{x}}, \hat{\mathbf{y}}) \in \mathbb{R}^{n+m}$, we can get

$$\ell'(\hat{x}_i; x_i - \hat{x}_i) = \begin{cases} |x_i|, & \text{if } i \in I_1(\hat{\mathbf{x}}), \\ \text{sgn}(\hat{x}_i)(x_i - \hat{x}_i), & \text{otherwise,} \end{cases} \quad (2.2)$$

119 and

$$\rho'_j(\hat{\mathbf{y}}^{(j)}; \mathbf{y}^{(j)} - \hat{\mathbf{y}}^{(j)}) = \begin{cases} \|\mathbf{y}^{(j)}\|, & \text{if } j \in J_1(\hat{\mathbf{y}}), \\ \frac{\hat{\mathbf{y}}^{(j)\top}(\mathbf{y}^{(j)} - \hat{\mathbf{y}}^{(j)})}{\|\hat{\mathbf{y}}^{(j)}\|}, & \text{otherwise,} \end{cases} \quad (2.3)$$

120 where

$$\text{sgn}(t) = \begin{cases} 1, & \text{if } t > 0, \\ [-1, 1], & \text{if } t = 0, \\ -1, & \text{if } t < 0. \end{cases}$$

121 **Definition 2.2** [13] Let $h : \mathbb{R}^{n+m} \rightarrow \mathbb{R} \cup \{\infty\}$ be locally Lipschitz at $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$, for any
122 $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{n+m}$ near $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$, the generalized directional derivative of h at $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$ is defined as

$$h^\circ((\hat{\mathbf{x}}, \hat{\mathbf{y}}); (\mathbf{x} - \hat{\mathbf{x}}, \mathbf{y} - \hat{\mathbf{y}})) := \limsup_{\substack{(\mathbf{x}, \mathbf{y}) \rightarrow (\hat{\mathbf{x}}, \hat{\mathbf{y}}) \\ t \downarrow 0}} \frac{h((\mathbf{x}, \mathbf{y}) + t(\mathbf{x} - \hat{\mathbf{x}}, \mathbf{y} - \hat{\mathbf{y}})) - h(\mathbf{x}, \mathbf{y})}{t}.$$

123 As we all know, the existence of the generalized directional derivative does not imply
124 the existence of the directional derivative. But if the directional derivative exists, then

$$h'((\hat{\mathbf{x}}, \hat{\mathbf{y}}); (\mathbf{x} - \hat{\mathbf{x}}, \mathbf{y} - \hat{\mathbf{y}})) \leq h^\circ((\hat{\mathbf{x}}, \hat{\mathbf{y}}); (\mathbf{x} - \hat{\mathbf{x}}, \mathbf{y} - \hat{\mathbf{y}})), \quad \forall (\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{n+m}. \quad (2.4)$$

125 Next, we introduce several types of definitions of subdifferential.

126 **Definition 2.3** [34] Let $h : \mathbb{R}^{n+m} \rightarrow \mathbb{R} \cup \{\infty\}$ be a proper convex function, the subdiffer-
127 ential $\partial h(\hat{\mathbf{x}}, \hat{\mathbf{y}})$ of h at $(\hat{\mathbf{x}}, \hat{\mathbf{y}}) \in \text{dom}h$ is the set of $\xi \in \partial h(\hat{\mathbf{x}}, \hat{\mathbf{y}})$, called subgradients of h at
128 $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$, such that

$$h(\mathbf{x}, \mathbf{y}) \geq h(\hat{\mathbf{x}}, \hat{\mathbf{y}}) + \langle \xi, (\mathbf{x} - \hat{\mathbf{x}}, \mathbf{y} - \hat{\mathbf{y}}) \rangle, \quad \forall (\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{n+m}.$$

129 If $(\hat{\mathbf{x}}, \hat{\mathbf{y}}) \notin \text{dom}h$, then $\partial h(\hat{\mathbf{x}}, \hat{\mathbf{y}}) = \emptyset$.

130 **Definition 2.4** [13] Let $h : \mathbb{R}^{n+m} \rightarrow \mathbb{R} \cup \{\infty\}$ be a locally Lipschitz function. The Clarke
131 subdifferential of h at $(\hat{\mathbf{x}}, \hat{\mathbf{y}}) \in \text{dom}h$, written $\partial^C h(\hat{\mathbf{x}}, \hat{\mathbf{y}})$, is defined as

$$\text{con}\{\xi \in \mathbb{R}^{n+m} \mid \langle \xi, (\mathbf{x} - \hat{\mathbf{x}}, \mathbf{y} - \hat{\mathbf{y}}) \rangle \leq h^\circ((\hat{\mathbf{x}}, \hat{\mathbf{y}}); (\mathbf{x} - \hat{\mathbf{x}}, \mathbf{y} - \hat{\mathbf{y}})), \forall (\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{n+m}\},$$

132 where "con" represents the convex hull of a set.

133 The above definition implies that [13, Corollary 2.9.1]

$$h^\circ((\hat{\mathbf{x}}, \hat{\mathbf{y}}); (\mathbf{x} - \hat{\mathbf{x}}, \mathbf{y} - \hat{\mathbf{y}})) = \max_{\xi \in \partial^C h(\hat{\mathbf{x}}, \hat{\mathbf{y}})} \langle \xi, (\mathbf{x} - \hat{\mathbf{x}}, \mathbf{y} - \hat{\mathbf{y}}) \rangle.$$

134 It is known that [13, Proposition 2.3.6] if h is convex, then $h^\circ((\hat{\mathbf{x}}, \hat{\mathbf{y}}); (\mathbf{x} - \hat{\mathbf{x}}, \mathbf{y} - \hat{\mathbf{y}})) =$
135 $h'((\hat{\mathbf{x}}, \hat{\mathbf{y}}); (\mathbf{x} - \hat{\mathbf{x}}, \mathbf{y} - \hat{\mathbf{y}}))$ and $\partial^C h(\mathbf{x}, \mathbf{y}) = \partial h(\mathbf{x}, \mathbf{y})$; if h is continuously differentiable, then
136 $\partial^C h(\mathbf{x}, \mathbf{y}) = \{\nabla h(\mathbf{x}, \mathbf{y})\}$.

137 Since the penalty terms in (1.2) are known as Capped- ℓ_1 functions, we can gain that
138 the objective function F is nonconvex and lower semicontinuous. We now give the definition
139 of limiting subdifferential.

140 **Definition 2.5** [34] Let $h : \mathbb{R}^{n+m} \rightarrow \mathbb{R} \cup \{\infty\}$ be a proper lower semicontinuous function.

141 (i) The Fréchet subdifferential of h at $(\hat{\mathbf{x}}, \hat{\mathbf{y}}) \in \text{dom}h$, written $\widehat{\partial}h(\hat{\mathbf{x}}, \hat{\mathbf{y}})$, is defined as

$$\{\xi \in \mathbb{R}^{n+m} \mid h(\mathbf{x}, \mathbf{y}) \geq h(\hat{\mathbf{x}}, \hat{\mathbf{y}}) + \langle \xi, (\mathbf{x} - \hat{\mathbf{x}}, \mathbf{y} - \hat{\mathbf{y}}) \rangle + o(\|(\mathbf{x} - \hat{\mathbf{x}}, \mathbf{y} - \hat{\mathbf{y}})\|), \forall (\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{n+m}\},$$

142 If $(\hat{\mathbf{x}}, \hat{\mathbf{y}}) \notin \text{dom}h := \{(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{n+m} \mid h(\mathbf{x}, \mathbf{y}) < \infty\}$, then $\widehat{\partial}h(\hat{\mathbf{x}}, \hat{\mathbf{y}}) = \emptyset$.

143 (ii) The limiting subdifferential of h at $(\hat{\mathbf{x}}, \hat{\mathbf{y}}) \in \text{dom}h$, written $\partial h(\hat{\mathbf{x}}, \hat{\mathbf{y}})$, is defined as

$$\{\xi \in \mathbb{R}^{n+m} \mid \exists (\mathbf{x}^k, \mathbf{y}^k) \rightarrow (\hat{\mathbf{x}}, \hat{\mathbf{y}}), h(\mathbf{x}^k, \mathbf{y}^k) \rightarrow h(\hat{\mathbf{x}}, \hat{\mathbf{y}}), \xi^k \in \widehat{\partial}h(\mathbf{x}^k, \mathbf{y}^k) \text{ such that } \xi^k \rightarrow \xi\},$$

144 If $(\hat{\mathbf{x}}, \hat{\mathbf{y}}) \notin \text{dom}h$, then $\partial h(\hat{\mathbf{x}}, \hat{\mathbf{y}}) = \emptyset$.

145 From [34], it is known that if h is locally Lipschitz, then $\partial^C h(\hat{\mathbf{x}}, \hat{\mathbf{y}}) = \text{cl}(\text{con}(\partial h(\hat{\mathbf{x}}, \hat{\mathbf{y}})))$
 146 which is the closed convex hull of $\partial h(\hat{\mathbf{x}}, \hat{\mathbf{y}})$. If h is a convex function, then the Fréchet
 147 subdifferential, limit subdifferential and Clarke subdifferential of h at (\mathbf{x}, \mathbf{y}) are all consistent
 148 with the classical subdifferential of convex function.

149 3 Directional stationary points and critical points of problem (1.2)

150 The optimality conditions of optimization problems are often characterized by stationary
 151 points. In this section, we give the characterization of the d(irectional)-stationary points and
 152 the critical points of problem (1.2), and analyze their properties. Then we investigate the
 153 relationship between the two types of stationary points.

154 Based on the DC expression (1.3) of problem (1.2), we give the definition of critical
 155 point of problem (1.2).

156 **Definition 3.1** [29, 34] [critical point] $(\hat{\mathbf{x}}, \hat{\mathbf{y}}) \in \mathbb{R}^{n+m}$ is called a critical point of problem
 157 (1.2), if

$$\begin{aligned} \mathbf{0} \in & \nabla f(\hat{\mathbf{x}}, \hat{\mathbf{y}}) + \lambda_1 \partial \left(\sum_{i=1}^n (g_1 \circ \ell)(\hat{x}_i) \right) - \lambda_1 \partial \left(\sum_{i=1}^n (h_1 \circ \ell)(\hat{x}_i) \right) \\ & + \lambda_2 \partial \left(\sum_{j=1}^J (g_2 \circ \rho_j)(\hat{\mathbf{y}}_{(j)}) \right) - \lambda_2 \partial \left(\sum_{j=1}^J (h_2 \circ \rho_j)(\hat{\mathbf{y}}_{(j)}) \right). \end{aligned}$$

158 The set of critical points of problem (1.2) is denoted by $\text{crit}F$.

159 Based on this definition, [34, Proposition 10.5] and [42, Theorem 3.4], we give the char-
 160 acterization of critical point of problem (1.2) as follows.

161 **Theorem 3.2** Let $(\hat{\mathbf{x}}, \hat{\mathbf{y}}) \in \mathbb{R}^{n+m}$ be a critical point of problem (1.2), then

$$\begin{aligned} \mathbf{0} \in & \nabla f(\hat{\mathbf{x}}, \hat{\mathbf{y}}) + \lambda_1 \left\{ \partial(g_1 \circ \ell)(\hat{x}_1) \times \cdots \times \partial(g_1 \circ \ell)(\hat{x}_n) - \partial(h_1 \circ \ell)(\hat{x}_1) \times \cdots \times \partial(h_1 \circ \ell)(\hat{x}_n) \right\} \\ & + \lambda_2 \left\{ \partial(g_2 \circ \rho_1)(\hat{\mathbf{y}}_{(1)}) \times \cdots \times \partial(g_2 \circ \rho_J)(\hat{\mathbf{y}}_{(J)}) - \partial(h_2 \circ \rho_1)(\hat{\mathbf{y}}_{(1)}) \times \cdots \times \partial(h_2 \circ \rho_J)(\hat{\mathbf{y}}_{(J)}) \right\}, \end{aligned}$$

162 where

$$\begin{aligned}\partial(g_1 \circ \ell)(\hat{x}_i) &= \begin{cases} [-\frac{1}{\alpha_1}, \frac{1}{\alpha_1}], & \text{if } i \in I_1(\hat{\mathbf{x}}), \\ \{\frac{1}{\alpha_1} \text{sgn}(\hat{x}_i)\}, & \text{otherwise,} \end{cases} \\ \partial(h_1 \circ \ell)(\hat{x}_i) &= \begin{cases} 0, & \text{if } i \in I_1(\hat{\mathbf{x}}) \cup I_2(\hat{\mathbf{x}}), \\ \text{con}\{0, \frac{1}{\alpha_2} \text{sgn}(\hat{x}_i)\}, & \text{if } i \in I_3(\hat{\mathbf{x}}), \\ \frac{1}{\alpha_2} \text{sgn}(\hat{x}_i), & \text{if } i \in I_4(\hat{\mathbf{x}}), \end{cases} \\ \partial(g_2 \circ \rho_j)(\hat{\mathbf{y}}_{(j)}) &= \begin{cases} \frac{1}{\alpha_2} B^{m_j}, & \text{if } j \in J_1(\hat{\mathbf{y}}), \\ \{\frac{\hat{\mathbf{y}}_{(j)}}{\alpha_2 \|\hat{\mathbf{y}}_{(j)}\|}\}, & \text{otherwise,} \end{cases} \\ \partial(h_2 \circ \rho_j)(\hat{\mathbf{y}}_{(j)}) &= \begin{cases} \mathbf{0}, & \text{if } j \in J_1(\hat{\mathbf{y}}) \cup J_2(\hat{\mathbf{y}}), \\ \text{con}\{\mathbf{0}, \frac{\hat{\mathbf{y}}_{(j)}}{\alpha_2 \|\hat{\mathbf{y}}_{(j)}\|}\}, & \text{if } j \in J_3(\hat{\mathbf{y}}), \\ \{\frac{\hat{\mathbf{y}}_{(j)}}{\alpha_2 \|\hat{\mathbf{y}}_{(j)}\|}\}, & \text{if } j \in J_4(\hat{\mathbf{y}}), \end{cases}\end{aligned}$$

163 in which "con" denotes the convex hull of a set and B^{m_j} denotes the closed unit ball in \mathbb{R}^{m_j} .

164 Proof According to [34, Proposition 10.5], we get

$$\begin{aligned}\partial \left(\sum_{i=1}^n (g_1 \circ \ell)(\hat{x}_i) \right) &= \left\{ \partial(g_1 \circ \ell)(\hat{x}_1) \times \cdots \times \partial(g_1 \circ \ell)(\hat{x}_n) \right\} \\ \partial \left(\sum_{i=1}^n (h_1 \circ \ell)(\hat{x}_i) \right) &= \left\{ \partial(h_1 \circ \ell)(\hat{x}_1) \times \cdots \times \partial(h_1 \circ \ell)(\hat{x}_n) \right\} \\ \partial \left(\sum_{j=1}^J (g_2 \circ \rho_j)(\hat{\mathbf{y}}_{(j)}) \right) &= \left\{ \partial(g_2 \circ \rho_1)(\hat{\mathbf{y}}_{(1)}) \times \cdots \times \partial(g_2 \circ \rho_J)(\hat{\mathbf{y}}_{(J)}) \right\} \\ \partial \left(\sum_{j=1}^J (h_2 \circ \rho_j)(\hat{\mathbf{y}}_{(j)}) \right) &= \left\{ \partial(h_2 \circ \rho_1)(\hat{\mathbf{y}}_{(1)}) \times \cdots \times \partial(h_2 \circ \rho_J)(\hat{\mathbf{y}}_{(J)}) \right\}.\end{aligned}$$

165 By the definition of critical point and the direct calculation, we get that

$$\begin{aligned}\partial(g_1 \circ \ell)(\hat{x}_i) &= \begin{cases} [-\frac{1}{\alpha_1}, \frac{1}{\alpha_1}], & \text{if } i \in I_1(\hat{\mathbf{x}}), \\ \{\frac{1}{\alpha_1} \text{sgn}(\hat{x}_i)\}, & \text{otherwise,} \end{cases} \\ \partial(h_1 \circ \ell)(\hat{x}_i) &= \begin{cases} 0, & \text{if } i \in I_1(\hat{\mathbf{x}}) \cup I_2(\hat{\mathbf{x}}), \\ \text{con}\{0, \frac{1}{\alpha_2} \text{sgn}(\hat{x}_i)\}, & \text{if } i \in I_3(\hat{\mathbf{x}}), \\ \{\frac{1}{\alpha_2} \text{sgn}(\hat{x}_i)\}, & \text{if } i \in I_4(\hat{\mathbf{x}}). \end{cases}\end{aligned}$$

166 Similar to [42, Theorem 3.4], we can get

$$\begin{aligned}\partial(g_2 \circ \rho_j)(\hat{\mathbf{y}}_{(j)}) &= \begin{cases} \frac{1}{\alpha_2} B^{m_j}, & \text{if } j \in J_1(\hat{\mathbf{y}}), \\ \{\frac{\hat{\mathbf{y}}_{(j)}}{\alpha_2 \|\hat{\mathbf{y}}_{(j)}\|}\}, & \text{otherwise,} \end{cases} \\ \partial(h_2 \circ \rho_j)(\hat{\mathbf{y}}_{(j)}) &= \begin{cases} \mathbf{0}, & \text{if } j \in J_1(\hat{\mathbf{y}}) \cup J_2(\hat{\mathbf{y}}), \\ \text{con}\{\mathbf{0}, \frac{\hat{\mathbf{y}}_{(j)}}{\alpha_2 \|\hat{\mathbf{y}}_{(j)}\|}\}, & \text{if } j \in J_3(\hat{\mathbf{y}}), \\ \{\frac{\hat{\mathbf{y}}_{(j)}}{\alpha_2 \|\hat{\mathbf{y}}_{(j)}\|}\}, & \text{if } j \in J_4(\hat{\mathbf{y}}). \end{cases}\end{aligned}$$

167 The proof is thus finished. \square

168 Now we give the definition of d-stationary point of problem (1.2).

169 **Definition 3.3** [42] [*d-stationary point*] $(\hat{\mathbf{x}}, \hat{\mathbf{y}}) \in \mathbb{R}^{n+m}$ is called a *d-stationary point* of
170 problem (1.2), if

$$F'((\hat{\mathbf{x}}, \hat{\mathbf{y}}); (\mathbf{x} - \hat{\mathbf{x}}, \mathbf{y} - \hat{\mathbf{y}})) \geq 0, \quad \forall (\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{n+m}.$$

171 The following theorem gives the characterization of d-stationary point of problem (1.2).

172 **Theorem 3.4** Let $(\hat{\mathbf{x}}, \hat{\mathbf{y}}) \in \mathbb{R}^{n+m}$ be a *d-stationary point* of problem (1.2), then

$$\begin{aligned} & \langle \nabla f(\hat{\mathbf{x}}, \hat{\mathbf{y}}), (\mathbf{x} - \hat{\mathbf{x}}, \mathbf{y} - \hat{\mathbf{y}}) \rangle + \lambda_1 \left(\sum_{i=1}^n (g_1 \circ \ell)'(\hat{x}_i; x_i - \hat{x}_i) \right) - \lambda_1 \left(\sum_{i=1}^n (h_1 \circ \ell)'(\hat{x}_i; x_i - \hat{x}_i) \right) \\ & + \lambda_2 \left(\sum_{j=1}^J (g_2 \circ \rho_j)'(\hat{\mathbf{y}}^{(j)}; \mathbf{y}^{(j)} - \hat{\mathbf{y}}^{(j)}) \right) - \lambda_2 \left(\sum_{j=1}^J (h_2 \circ \rho_j)'(\hat{\mathbf{y}}^{(j)}; \mathbf{y}^{(j)} - \hat{\mathbf{y}}^{(j)}) \right) \geq 0 \end{aligned}$$

173 for any $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{n+m}$, where

$$\begin{aligned} (g_1 \circ \ell)'(\hat{x}_i; x_i - \hat{x}_i) &= \begin{cases} \frac{|x_i|}{\alpha_1}, & \text{if } i \in I_1(\hat{\mathbf{x}}), \\ \frac{\text{sgn}(\hat{x}_i)(x_i - \hat{x}_i)}{\alpha_1}, & \text{otherwise,} \end{cases} \\ (h_1 \circ \ell)'(\hat{x}_i; x_i - \hat{x}_i) &= \begin{cases} 0, & \text{if } i \in I_1(\hat{\mathbf{x}}) \cup I_2(\hat{\mathbf{x}}), \\ \max\{0, \frac{\text{sgn}(\hat{x}_i)(x_i - \hat{x}_i)}{\alpha_1}\}, & \text{if } i \in I_3(\hat{\mathbf{x}}), \\ \frac{\text{sgn}(\hat{x}_i)(x_i - \hat{x}_i)}{\alpha_1}, & \text{if } i \in I_4(\hat{\mathbf{x}}). \end{cases} \\ (g_2 \circ \rho_j)'(\hat{\mathbf{y}}^{(j)}; \mathbf{y}^{(j)} - \hat{\mathbf{y}}^{(j)}) &= \begin{cases} \frac{\|\mathbf{y}^{(j)}\|}{\alpha_2}, & \text{if } j \in J_1(\hat{\mathbf{y}}), \\ \frac{\hat{\mathbf{y}}_{(j)}^{\top}(\mathbf{y}^{(j)} - \hat{\mathbf{y}}^{(j)})}{\alpha_2 \|\hat{\mathbf{y}}^{(j)}\|}, & \text{otherwise.} \end{cases} \\ (h_2 \circ \rho_j)'(\hat{\mathbf{y}}^{(j)}; \mathbf{y}^{(j)} - \hat{\mathbf{y}}^{(j)}) &= \begin{cases} 0, & \text{if } j \in J_1(\hat{\mathbf{y}}) \cup J_2(\hat{\mathbf{y}}), \\ \max\{0, \frac{\hat{\mathbf{y}}_{(j)}^{\top}(\mathbf{y}^{(j)} - \hat{\mathbf{y}}^{(j)})}{\alpha_2 \|\hat{\mathbf{y}}^{(j)}\|}\}, & \text{if } j \in J_3(\hat{\mathbf{y}}), \\ \frac{\hat{\mathbf{y}}_{(j)}^{\top}(\mathbf{y}^{(j)} - \hat{\mathbf{y}}^{(j)})}{\alpha_2 \|\hat{\mathbf{y}}^{(j)}\|}, & \text{if } j \in J_4(\hat{\mathbf{y}}). \end{cases} \end{aligned} \quad (3.1)$$

174 *Proof* From the definition of d-stationary point, the DC form (1.3) and the analysis similar
175 to [42, Theorem 3.2], we can directly obtain the conclusion. \square

176 The following theorem provides the relationship between d-stationary point and critical
177 point.

178 **Theorem 3.5** Let $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$ be a *d-stationary point* of problem (1.2), then $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$ is a *critical*
179 *point* of problem (1.2).

180 *Proof* Since $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$ is a d-stationary point of problem (1.2), according to inequality (2.4), we
181 have

$$\begin{aligned} 0 &\leq F'((\hat{\mathbf{x}}, \hat{\mathbf{y}}); (\mathbf{x} - \hat{\mathbf{x}}, \mathbf{y} - \hat{\mathbf{y}})) \\ &\leq F^\circ((\hat{\mathbf{x}}, \hat{\mathbf{y}}); (\mathbf{x} - \hat{\mathbf{x}}, \mathbf{y} - \hat{\mathbf{y}})) \\ &= \max_{\xi \in \partial^C F(\hat{\mathbf{x}}, \hat{\mathbf{y}})} \langle \xi, (\mathbf{x} - \hat{\mathbf{x}}, \mathbf{y} - \hat{\mathbf{y}}) \rangle. \end{aligned}$$

182 Therefore, according to the operational properties of Clarke differential [13, Propostion 2.3.3,
183 Corollary 2.3.3.2], we obtain that

$$\begin{aligned} \mathbf{0} &\in \partial^C F(\hat{\mathbf{x}}, \hat{\mathbf{y}}) \\ &\subseteq \partial^C f(\hat{\mathbf{x}}, \hat{\mathbf{y}}) + \lambda_1 \partial^C \left(\sum_{i=1}^n g_1(|\hat{x}_i|) - h_1(|\hat{x}_i|) \right) + \lambda_2 \partial^C \left(\sum_{j=1}^J g_2(\|\hat{\mathbf{y}}_{(j)}\|) - h_2(\|\hat{\mathbf{y}}_{(j)}\|) \right) \\ &\subseteq \partial^C f(\hat{\mathbf{x}}, \hat{\mathbf{y}}) + \lambda_1 \partial^C \left(\sum_{i=1}^n g_1(|\hat{x}_i|) \right) - \lambda_1 \partial^C \left(\sum_{i=1}^n h_1(|\hat{x}_i|) \right) \\ &\quad + \lambda_2 \partial^C \left(\sum_{j=1}^J g_2(\|\hat{\mathbf{y}}_{(j)}\|) \right) - \lambda_2 \partial^C \left(\sum_{j=1}^J h_2(\|\hat{\mathbf{y}}_{(j)}\|) \right). \end{aligned}$$

184 Since f is continuously differentiable, g_ν and h_ν ($\nu = 1, 2$) are all convex functions, according
185 to [13, Propostion 2.3.6(b)], we get

$$\begin{aligned} \partial^C f(\hat{\mathbf{x}}, \hat{\mathbf{y}}) &= \partial f(\hat{\mathbf{x}}, \hat{\mathbf{y}}) = \{\nabla f(\hat{\mathbf{x}}, \hat{\mathbf{y}})\} \\ \partial^C \left(\sum_{i=1}^n g_1(|\hat{x}_i|) \right) &= \partial \left(\sum_{i=1}^n g_1(|\hat{x}_i|) \right) \\ \partial^C \left(\sum_{i=1}^n h_1(|\hat{x}_i|) \right) &= \partial \left(\sum_{i=1}^n h_1(|\hat{x}_i|) \right) \end{aligned}$$

186 and

$$\begin{aligned} \partial^C \left(\sum_{j=1}^J g_2(\|\hat{\mathbf{y}}_{(j)}\|) \right) &= \partial \left(\sum_{j=1}^J g_2(\|\hat{\mathbf{y}}_{(j)}\|) \right) \\ \partial^C \left(\sum_{j=1}^J h_2(\|\hat{\mathbf{y}}_{(j)}\|) \right) &= \partial \left(\sum_{j=1}^J h_2(\|\hat{\mathbf{y}}_{(j)}\|) \right), \end{aligned}$$

187 then

$$\begin{aligned} \mathbf{0} &\in \nabla f(\hat{\mathbf{x}}, \hat{\mathbf{y}}) + \lambda_1 \partial \left(\sum_{i=1}^n g_1(|\hat{x}_i|) \right) - \lambda_1 \partial \left(\sum_{i=1}^n h_1(|\hat{x}_i|) \right) \\ &\quad + \lambda_2 \partial \left(\sum_{j=1}^J g_2(\|\hat{\mathbf{y}}_{(j)}\|) \right) - \lambda_2 \partial \left(\sum_{j=1}^J h_2(\|\hat{\mathbf{y}}_{(j)}\|) \right). \end{aligned}$$

188 From Definition 3.1, the above inequality implies that $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$ is a critical point of problem
189 (1.2). \square

190 **Remark 3.6** From the proof of Lemma 3.5, we have that if $\mathbf{0} \in \partial^C F(\hat{\mathbf{x}}, \hat{\mathbf{y}})$, then $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$
191 is a critical point of problem (1.2).

192 The following lemma characterize the property of gradient of f at the d-stationary point
193 of problem (1.2).

194 **Lemma 3.7** Let $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$ be a d -stationary point of problem (1.2), the following statements
195 hold:

- 196 (i) $|\nabla_{\mathbf{x}} f(\hat{\mathbf{x}}, \hat{\mathbf{y}})]_i| = \frac{\lambda_1}{\alpha_1}, \forall i \in I_2(\hat{\mathbf{x}}); \quad [\nabla_{\mathbf{x}} f(\hat{\mathbf{x}}, \hat{\mathbf{y}})]_i = 0, \forall i \in I_4(\hat{\mathbf{x}}).$
197 (ii) $\|[\nabla_{\mathbf{y}} f(\hat{\mathbf{x}}, \hat{\mathbf{y}})]_{(j)}\| = \frac{\lambda_2}{\alpha_2}, \forall j \in J_2(\hat{\mathbf{y}}); \quad \|[\nabla_{\mathbf{y}} f(\hat{\mathbf{x}}, \hat{\mathbf{y}})]_{(j)}\| = 0, \forall i \in J_4(\hat{\mathbf{y}}).$

198 *Proof* (i). From Theorem 3.4, for any $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{n+m}$, we have

$$\begin{aligned} 0 &\leq F'((\hat{\mathbf{x}}, \hat{\mathbf{y}}); (\mathbf{x} - \hat{\mathbf{x}}, \mathbf{y} - \hat{\mathbf{y}})) \\ &= \langle \nabla f(\hat{\mathbf{x}}, \hat{\mathbf{y}}), (\mathbf{x} - \hat{\mathbf{x}}, \mathbf{y} - \hat{\mathbf{y}}) \rangle + \lambda_1 \sum_{i=1}^n (g_1 \circ \ell)'(\hat{x}_i; x_i - \hat{x}_i) - \lambda_1 \sum_{i=1}^n (h_1 \circ \ell)'(\hat{x}_i; x_i - \hat{x}_i) \\ &\quad + \lambda_2 \sum_{j=1}^J (g_2 \circ \rho_j)'(\hat{\mathbf{y}}_{(j)}; \mathbf{y}_{(j)} - \hat{\mathbf{y}}_{(j)}) - \lambda_2 \sum_{j=1}^J (h_2 \circ \rho_j)'(\hat{\mathbf{y}}_{(j)}; \mathbf{y}_{(j)} - \hat{\mathbf{y}}_{(j)}). \end{aligned} \quad (3.2)$$

199 According to the arbitrariness of $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{n+m}$, let $\mathbf{y} = \hat{\mathbf{y}}$, then

$$\begin{aligned} 0 &\leq F'((\hat{\mathbf{x}}, \hat{\mathbf{y}}); (\mathbf{x} - \hat{\mathbf{x}}, \mathbf{0})) \\ &= \sum_{i=1}^n [\nabla_{\mathbf{x}} f(\hat{\mathbf{x}}, \hat{\mathbf{y}})]_i (x_i - \hat{x}_i) + \lambda_1 \sum_{i=1}^n (g_1 \circ \ell)'(\hat{x}_i; x_i - \hat{x}_i) - \lambda_1 \sum_{i=1}^n (h_1 \circ \ell)'(\hat{x}_i; x_i - \hat{x}_i). \end{aligned}$$

200 From (3.1), we have

$$\begin{aligned} 0 &\leq F'((\hat{\mathbf{x}}, \hat{\mathbf{y}}); (\mathbf{x} - \hat{\mathbf{x}}, \mathbf{0})) \\ &= \sum_{i=1}^n [\nabla_{\mathbf{x}} f(\hat{\mathbf{x}}, \hat{\mathbf{y}})]_i (x_i - \hat{x}_i) + \frac{\lambda_1}{\alpha_1} \left(\sum_{i \in I_1(\hat{\mathbf{x}})} |x_i| + \sum_{i \in [n] \setminus I_1(\hat{\mathbf{x}})} \operatorname{sgn}(\hat{x}_i) (x_i - \hat{x}_i) \right. \\ &\quad \left. - \sum_{i \in I_3(\hat{\mathbf{x}})} \max\{0, \operatorname{sgn}(\hat{x}_i) (x_i - \hat{x}_i)\} - \sum_{i \in I_4(\hat{\mathbf{x}})} \operatorname{sgn}(\hat{x}_i) (x_i - \hat{x}_i) \right) \\ &= \sum_{i=1}^n [\nabla_{\mathbf{x}} f(\hat{\mathbf{x}}, \hat{\mathbf{y}})]_i (x_i - \hat{x}_i) + \frac{\lambda_1}{\alpha_1} \left(\sum_{i \in I_1(\hat{\mathbf{x}})} |x_i| + \sum_{i \in I_2(\hat{\mathbf{x}}) \cup I_3(\hat{\mathbf{x}})} \operatorname{sgn}(\hat{x}_i) (x_i - \hat{x}_i) \right. \\ &\quad \left. - \sum_{i \in I_3(\hat{\mathbf{x}})} \max\{0, \operatorname{sgn}(\hat{x}_i) (x_i - \hat{x}_i)\} \right). \end{aligned} \quad (3.3)$$

Let

$$\tilde{x}_i^1 = \begin{cases} \hat{x}_i, & \text{if } i \in [n] \setminus I_2(\hat{\mathbf{x}}), \\ \hat{x}_i - ([\nabla_{\mathbf{x}} f(\hat{\mathbf{x}}, \hat{\mathbf{y}})]_i + \frac{\lambda_1}{\alpha_1} \operatorname{sgn}(\hat{x}_i)), & \text{if } i \in I_2(\hat{\mathbf{x}}), \end{cases}$$

201 then from (3.3), we obtain that

$$\begin{aligned} 0 &\leq F'((\hat{\mathbf{x}}, \hat{\mathbf{y}}); (\tilde{\mathbf{x}}^1 - \hat{\mathbf{x}}, \mathbf{0})) \\ &= \sum_{i \in I_2(\hat{\mathbf{x}})} ([\nabla_{\mathbf{x}} f(\hat{\mathbf{x}}, \hat{\mathbf{y}})]_i + \frac{\lambda_1}{\alpha_1} \operatorname{sgn}(\hat{x}_i)) (\tilde{x}_i^1 - \hat{x}_i), \\ &= - \sum_{i \in I_2(\hat{\mathbf{x}})} \left([\nabla_{\mathbf{x}} f(\hat{\mathbf{x}}, \hat{\mathbf{y}})]_i + \frac{\lambda_1}{\alpha_1} \operatorname{sgn}(\hat{x}_i) \right)^2. \end{aligned} \quad (3.4)$$

202 From inequality (3.4), we obtain

$$[\nabla_{\mathbf{x}} f(\hat{\mathbf{x}}, \hat{\mathbf{y}})]_i + \frac{\lambda_1}{\alpha_1} \operatorname{sgn}(\hat{x}_i) = 0, \quad \forall i \in I_2(\hat{\mathbf{x}}).$$

203 Thus, $|\llbracket \nabla_{\mathbf{x}} f(\hat{\mathbf{x}}, \hat{\mathbf{y}}) \rrbracket_i| = \frac{\lambda_1}{\alpha_1}, \forall i \in I_2(\hat{\mathbf{x}})$.

Let

$$\tilde{x}_i^2 = \begin{cases} \hat{x}_i, & \text{if } i \in [n] \setminus I_4(\hat{\mathbf{x}}), \\ \hat{x}_i - \llbracket \nabla_{\mathbf{x}} f(\hat{\mathbf{x}}, \hat{\mathbf{y}}) \rrbracket_i, & \text{if } i \in I_4(\hat{\mathbf{x}}), \end{cases}$$

204 then from (3.3), we obtain that

$$\begin{aligned} 0 &\leq F'((\hat{\mathbf{x}}, \hat{\mathbf{y}}); (\tilde{\mathbf{x}}^2 - \hat{\mathbf{x}}, \mathbf{0})) \\ &= \sum_{i \in I_4(\hat{\mathbf{x}})} (\llbracket \nabla_{\mathbf{x}} f(\hat{\mathbf{x}}, \hat{\mathbf{y}}) \rrbracket_i) (\tilde{x}_i^2 - \hat{x}_i) = - \sum_{i \in I_4(\hat{\mathbf{x}})} (\llbracket \nabla_{\mathbf{x}} f(\hat{\mathbf{x}}, \hat{\mathbf{y}}) \rrbracket_i)^2, \end{aligned} \quad (3.5)$$

205 From inequality (3.5), we obtain

$$\llbracket \nabla_{\mathbf{x}} f(\hat{\mathbf{x}}, \hat{\mathbf{y}}) \rrbracket_i = 0, \quad \forall i \in I_4(\hat{\mathbf{x}}).$$

206 (ii). The proof is similar to that of (i). By the arbitrariness of $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{n+m}$ in (3.2),
207 take $\mathbf{x} = \hat{\mathbf{x}}$, then

$$\begin{aligned} 0 &\leq F'((\hat{\mathbf{x}}, \hat{\mathbf{y}}); (\mathbf{0}, \mathbf{y} - \hat{\mathbf{y}})) \\ &= \sum_{j=1}^J \llbracket \nabla_{\mathbf{y}} f(\hat{\mathbf{x}}, \hat{\mathbf{y}}) \rrbracket_{(j)} (\mathbf{y}_{(j)} - \hat{\mathbf{y}}_{(j)}) + \lambda_2 \sum_{j=1}^J (g_2 \circ \rho_j)'(\hat{\mathbf{y}}_{(j)}; \mathbf{y}_{(j)} - \hat{\mathbf{y}}_{(j)}) \\ &\quad - \lambda_2 \sum_{j=1}^J (h_2 \circ \rho_j)'(\hat{\mathbf{y}}_{(j)}; \mathbf{y}_{(j)} - \hat{\mathbf{y}}_{(j)}). \end{aligned}$$

208 From (3.1), we have

$$\begin{aligned} 0 &\leq F'((\hat{\mathbf{x}}, \hat{\mathbf{y}}); (\mathbf{0}, \mathbf{y} - \hat{\mathbf{y}})) \\ &= \sum_{j=1}^J \llbracket \nabla_{\mathbf{y}} f(\hat{\mathbf{x}}, \hat{\mathbf{y}}) \rrbracket_{(j)}^\top (\mathbf{y}_{(j)} - \hat{\mathbf{y}}_{(j)}) + \frac{\lambda_2}{\alpha_2} \left(\sum_{j \in J_1(\hat{\mathbf{x}})} \|\mathbf{y}_{(j)}\| + \sum_{j \in [m] \setminus J_1(\hat{\mathbf{y}})} \frac{\hat{\mathbf{y}}_{(j)}^\top (\mathbf{y}_{(j)} - \hat{\mathbf{y}}_{(j)})}{\|\hat{\mathbf{y}}_{(j)}\|} \right. \\ &\quad \left. - \sum_{j \in J_3(\hat{\mathbf{y}})} \max \left\{ \mathbf{0}, \frac{\hat{\mathbf{y}}_{(j)}^\top (\mathbf{y}_{(j)} - \hat{\mathbf{y}}_{(j)})}{\|\hat{\mathbf{y}}_{(j)}\|} \right\} - \sum_{j \in J_4(\hat{\mathbf{y}})} \frac{\hat{\mathbf{y}}_{(j)}^\top (\mathbf{y}_{(j)} - \hat{\mathbf{y}}_{(j)})}{\|\hat{\mathbf{y}}_{(j)}\|} \right) \quad (3.6) \\ &= \sum_{j=1}^J \llbracket \nabla_{\mathbf{y}} f(\hat{\mathbf{x}}, \hat{\mathbf{y}}) \rrbracket_{(j)}^\top (\mathbf{y}_{(j)} - \hat{\mathbf{y}}_{(j)}) + \frac{\lambda_2}{\alpha_2} \left(\sum_{j \in J_1(\hat{\mathbf{x}})} \|\mathbf{y}_{(j)}\| + \sum_{j \in J_2(\hat{\mathbf{y}}) \cup J_3(\hat{\mathbf{y}})} \frac{\hat{\mathbf{y}}_{(j)}^\top (\mathbf{y}_{(j)} - \hat{\mathbf{y}}_{(j)})}{\|\hat{\mathbf{y}}_{(j)}\|} \right. \\ &\quad \left. - \sum_{j \in J_3(\hat{\mathbf{y}})} \max \left\{ \mathbf{0}, \frac{\hat{\mathbf{y}}_{(j)}^\top (\mathbf{y}_{(j)} - \hat{\mathbf{y}}_{(j)})}{\|\hat{\mathbf{y}}_{(j)}\|} \right\} \right). \end{aligned}$$

Let

$$\tilde{\mathbf{y}}_{(j)}^1 = \begin{cases} \hat{\mathbf{y}}_{(j)}, & \text{if } j \in [m] \setminus J_2(\hat{\mathbf{y}}), \\ \hat{\mathbf{y}}_{(j)} - \left(\llbracket \nabla_{\mathbf{y}} f(\hat{\mathbf{x}}, \hat{\mathbf{y}}) \rrbracket_{(j)} + \frac{\lambda_2}{\alpha_2} \frac{\hat{\mathbf{y}}_{(j)}}{\|\hat{\mathbf{y}}_{(j)}\|} \right), & \text{if } j \in J_2(\hat{\mathbf{y}}), \end{cases}$$

209 then from (3.6), we obtain that

$$\begin{aligned} 0 &\leq F'((\hat{\mathbf{x}}, \hat{\mathbf{y}}); (\mathbf{0}, \tilde{\mathbf{y}}^1 - \hat{\mathbf{y}})) \\ &= \sum_{j \in J_2(\hat{\mathbf{y}})} \left(\llbracket \nabla_{\mathbf{y}} f(\hat{\mathbf{x}}, \hat{\mathbf{y}}) \rrbracket_{(j)} + \frac{\lambda_2}{\alpha_2} \frac{\hat{\mathbf{y}}_{(j)}}{\|\hat{\mathbf{y}}_{(j)}\|} \right)^\top (\tilde{\mathbf{y}}_{(j)}^1 - \hat{\mathbf{y}}_{(j)}) \\ &= - \sum_{j \in J_2(\hat{\mathbf{y}})} \left\| \llbracket \nabla_{\mathbf{y}} f(\hat{\mathbf{x}}, \hat{\mathbf{y}}) \rrbracket_{(j)} + \frac{\lambda_2}{\alpha_2} \frac{\hat{\mathbf{y}}_{(j)}}{\|\hat{\mathbf{y}}_{(j)}\|} \right\|^2. \end{aligned} \quad (3.7)$$

210 From inequality (3.7), we obtain

$$[\nabla_{\mathbf{y}} f(\hat{\mathbf{x}}, \hat{\mathbf{y}})]_{(j)} + \frac{\lambda_2}{\alpha_2} \frac{\hat{\mathbf{y}}_{(j)}}{\|\hat{\mathbf{y}}_{(j)}\|} = 0, \text{ i.e., } [\nabla_{\mathbf{y}} f(\hat{\mathbf{x}}, \hat{\mathbf{y}})]_{(j)} = -\frac{\lambda_2}{\alpha_2} \frac{\hat{\mathbf{y}}_{(j)}}{\|\hat{\mathbf{y}}_{(j)}\|}, \forall j \in J_2(\hat{\mathbf{y}}).$$

211 Take ℓ_2 norm on both sides of the above equality, then we get

$$\|[\nabla_{\mathbf{y}} f(\hat{\mathbf{x}}, \hat{\mathbf{y}})]_{(j)}\| = \frac{\lambda_2}{\alpha_2}, \forall j \in J_2(\hat{\mathbf{y}}).$$

Let

$$\tilde{\mathbf{y}}_{(j)}^2 = \begin{cases} \hat{\mathbf{y}}_{(j)}, & \text{if } j \in [m] \setminus J_4(\hat{\mathbf{y}}), \\ \hat{\mathbf{y}}_{(j)} - [\nabla_{\mathbf{y}} f(\hat{\mathbf{x}}, \hat{\mathbf{y}})]_{(j)}, & \text{if } j \in J_4(\hat{\mathbf{y}}), \end{cases}$$

212 then from (3.6), we obtain that

$$\begin{aligned} 0 &\leq F'((\hat{\mathbf{x}}, \hat{\mathbf{y}}); (\mathbf{0}, \tilde{\mathbf{y}}^2 - \hat{\mathbf{y}})) \\ &= \sum_{j \in J_4(\hat{\mathbf{y}})} \left([\nabla_{\mathbf{y}} f(\hat{\mathbf{x}}, \hat{\mathbf{y}})]_{(j)} \right)^\top (\tilde{\mathbf{y}}_{(j)}^2 - \hat{\mathbf{y}}_{(j)}) \\ &= - \sum_{j \in J_4(\hat{\mathbf{y}})} \left\| [\nabla_{\mathbf{y}} f(\hat{\mathbf{x}}, \hat{\mathbf{y}})]_{(j)} \right\|^2. \end{aligned} \quad (3.8)$$

213 From inequality (3.8), we obtain

$$\|[\nabla_{\mathbf{y}} f(\hat{\mathbf{x}}, \hat{\mathbf{y}})]_{(j)}\| = 0, \forall j \in J_4(\hat{\mathbf{y}}).$$

214 The proof is thus complete. \square

215 The following theorem gives the lower bound property of the d-stationary points of
216 problem (1.2).

217 **Theorem 3.8** Let $(\hat{\mathbf{x}}, \hat{\mathbf{y}}) \in \mathbb{R}^{n+m}$ be a d-stationary point of problem (1.2). Suppose

218 $\|[\nabla f(\hat{\mathbf{x}}, \hat{\mathbf{y}})]_{I_2(\hat{\mathbf{x}}) \cup J_2(\hat{\mathbf{y}})}\| < \min \left\{ \frac{\lambda_1}{\alpha_1}, \frac{\lambda_2}{\alpha_2} \right\}$, then the following statements hold:

219 (i) $I_2(\hat{\mathbf{x}}) = \emptyset$, that is, if $\hat{x}_i \neq 0$, then $|\hat{x}_i| \geq \alpha_1$;

220 (ii) $J_2(\hat{\mathbf{y}}) = \emptyset$, that is, if $\hat{\mathbf{y}}_{(j)} \neq \mathbf{0}$, then $\|\hat{\mathbf{y}}_{(j)}\| \geq \alpha_2$.

221 *Proof* (i) Assume, on the contrary, that $I_2(\hat{\mathbf{x}}) \neq \emptyset$. Let $i_0 \in I_2(\hat{\mathbf{x}})$, then from Lemma 3.7, we
222 have

$$\frac{\lambda_1}{\alpha_1} = \|[\nabla_{\mathbf{x}} f(\hat{\mathbf{x}}, \hat{\mathbf{y}})]_{i_0}\| \leq \|[\nabla f(\hat{\mathbf{x}}, \hat{\mathbf{y}})]_{I_2(\hat{\mathbf{x}}) \cup J_2(\hat{\mathbf{y}})}\| < \frac{\lambda_1}{\alpha_1},$$

223 which is a contradiction, and implies that $I_2(\hat{\mathbf{x}}) = \emptyset$.

224 (ii) Assume, on the contrary, that $J_2(\hat{\mathbf{y}}) \neq \emptyset$. Let $j_0 \in J_2(\hat{\mathbf{y}})$, then from Lemma 3.7, we
225 have

$$\frac{\lambda_2}{\alpha_2} = \|[\nabla_{\mathbf{y}} f(\hat{\mathbf{x}}, \hat{\mathbf{y}})]_{(j_0)}\| \leq \|[\nabla f(\hat{\mathbf{x}}, \hat{\mathbf{y}})]_{I_2(\hat{\mathbf{x}}) \cup J_2(\hat{\mathbf{y}})}\| < \frac{\lambda_2}{\alpha_2},$$

226 which is a contradiction, and implies that $J_2(\hat{\mathbf{y}}) = \emptyset$. \square

227 **Remark 3.9** (1) If $\|[\nabla f(\hat{\mathbf{x}}, \hat{\mathbf{y}})]\| < \min \left\{ \frac{\lambda_1}{\alpha_1}, \frac{\lambda_2}{\alpha_2} \right\}$, then $\|[\nabla f(\hat{\mathbf{x}}, \hat{\mathbf{y}})]_{I_2(\hat{\mathbf{x}}) \cup J_2(\hat{\mathbf{y}})}\| < \min \left\{ \frac{\lambda_1}{\alpha_1}, \frac{\lambda_2}{\alpha_2} \right\}$.

228 (2) If f is locally Lipschitz at $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$ with modulus $L < \min \left\{ \frac{\lambda_1}{\alpha_1}, \frac{\lambda_2}{\alpha_2} \right\}$, then $\|[\nabla f(\hat{\mathbf{x}}, \hat{\mathbf{y}})]\| <$

229 $\min \left\{ \frac{\lambda_1}{\alpha_1}, \frac{\lambda_2}{\alpha_2} \right\}$.

230 (3) If $f : \mathbb{R}^{n+m} \rightarrow \mathbb{R}$ is convex, then f is locally Lipschitz on \mathbb{R}^{n+m} .

231 4 Equivalence of problem (1.1) and problem (1.2)

232 In this section, we investigate the relationship between the original problem (1.1) and
233 the relaxation problem (1.2) by considering the global solutions and local solutions of them.

234 **Theorem 4.1** *Suppose $\|\nabla f(\mathbf{x}, \mathbf{y})\| < \min\{\frac{\lambda_1}{\alpha_1}, \frac{\lambda_2}{\alpha_2}\}$ holds on \mathbb{R}^{n+m} , then the following
235 statements hold.*

236 (i) *The global optimal solution sets and optimal value of problem (1.1) are same as those
237 of problem (1.2) respectively;*

238 (ii) *If $(\hat{\mathbf{x}}, \hat{\mathbf{y}}) \in \mathbb{R}^{n+m}$ is a local minimizer of problem (1.2), then $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$ is also a local
239 minimizer of problem (1.1), and the objective function value of problems (1.1) and (1.2) at
240 $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$ are same.*

241 *Proof* (i). (a) Let $(\hat{\mathbf{x}}, \hat{\mathbf{y}}) \in \mathbb{R}^{n+m}$ be a global optimal solution of problem (1.2), then $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$
242 is also a d-stationary point of problem (1.2). From (1.2) and Theorem 3.8, we obtain

$$\varphi_1(|\hat{x}_i|) = \begin{cases} 0, & \text{if } i \in I_1(\hat{\mathbf{x}}) \cup I_2(\hat{\mathbf{x}}), \\ 1, & \text{if } i \in I_3(\hat{\mathbf{x}}) \cup I_4(\hat{\mathbf{x}}), \end{cases} \quad \text{and} \quad \varphi_2(\|\hat{\mathbf{y}}_{(j)}\|) = \begin{cases} 0, & \text{if } j \in J_1(\hat{\mathbf{y}}) \cup J_2(\hat{\mathbf{y}}), \\ 1, & \text{if } j \in J_3(\hat{\mathbf{y}}) \cup J_4(\hat{\mathbf{y}}), \end{cases}$$

243 then

$$\Phi_1(\hat{\mathbf{x}}) = \sum_{i \in I_3(\hat{\mathbf{x}}) \cup I_4(\hat{\mathbf{x}})} \varphi_1(|\hat{x}_i|) = \|\hat{\mathbf{x}}\|_0, \quad \Phi_2(\hat{\mathbf{y}}) = \sum_{j \in J_3(\hat{\mathbf{y}}) \cup J_4(\hat{\mathbf{y}})} \varphi_2(\|\hat{\mathbf{y}}_{(j)}\|) = \|\hat{\mathbf{y}}\|_{2,0}. \quad (4.1)$$

244 For any $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{n+m}$, since $\varphi_v(t) = \min\left\{1, \frac{t}{\alpha_v}\right\} \leq 1$, ($v = 1, 2$), then

$$\begin{aligned} \Phi_1(\mathbf{x}) &= \sum_{i \in [n] \setminus I_1(\mathbf{x})} \varphi_1(|x_i|) \leq \sum_{i \in [n] \setminus I_1(\mathbf{x})} 1 = \|\mathbf{x}\|_0, \\ \Phi_2(\mathbf{y}) &= \sum_{j \in [m] \setminus J_1(\hat{\mathbf{y}})} \varphi_2(\|\hat{\mathbf{y}}_{(j)}\|) \leq \sum_{j \in [m] \setminus J_1(\hat{\mathbf{y}})} 1 = \|\mathbf{y}\|_{2,0}. \end{aligned}$$

245 Thus, we have

$$\begin{aligned} f(\hat{\mathbf{x}}, \hat{\mathbf{y}}) + \lambda_1 \|\hat{\mathbf{x}}\|_0 + \lambda_2 \|\hat{\mathbf{y}}\|_{2,0} &= f(\hat{\mathbf{x}}, \hat{\mathbf{y}}) + \lambda_1 \Phi_1(\hat{\mathbf{x}}) + \lambda_2 \Phi_2(\hat{\mathbf{y}}) \\ &\leq f(\mathbf{x}, \mathbf{y}) + \lambda_1 \Phi_1(\mathbf{x}) + \lambda_2 \Phi_2(\mathbf{y}) \\ &\leq f(\mathbf{x}, \mathbf{y}) + \lambda_1 \|\mathbf{x}\|_0 + \lambda_2 \|\mathbf{y}\|_{2,0}. \end{aligned}$$

246 Therefore, $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$ is a global solution of problem (1.1), and (4.1) implies that optimal value
247 of problems (1.1) and (1.2) at $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$ are same.

248 (b) On the other hand, let $(\hat{\mathbf{x}}, \hat{\mathbf{y}}) \in \mathbb{R}^{n+m}$ be a global minimizer of problem (1.1).
249 Assume, on the contrary, that $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$ is not a global minimizer of problem (1.2), then

$$\Phi_1(\hat{\mathbf{x}}) \leq \|\hat{\mathbf{x}}\|_0 \quad \text{and} \quad \Phi_2(\hat{\mathbf{y}}) \leq \|\hat{\mathbf{y}}\|_{2,0}.$$

250 Let $(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \in \mathbb{R}^{n+m}$ be a global minimizer of problem (1.2), then

$$f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) + \lambda_1 \Phi_1(\bar{\mathbf{x}}) + \lambda_2 \Phi_2(\bar{\mathbf{y}}) < f(\hat{\mathbf{x}}, \hat{\mathbf{y}}) + \lambda_1 \Phi_1(\hat{\mathbf{x}}) + \lambda_2 \Phi_2(\hat{\mathbf{y}}).$$

251 What's more, from (i)(a), we know that

$$\Phi_1(\bar{\mathbf{x}}) = \|\bar{\mathbf{x}}\|_0 \quad \text{and} \quad \Phi_2(\bar{\mathbf{y}}) = \|\bar{\mathbf{y}}\|_{2,0}.$$

Thus, we have

$$\begin{aligned} f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) + \lambda_1 \|\bar{\mathbf{x}}\|_0 + \lambda_2 \|\bar{\mathbf{y}}\|_{2,0} &= f(\bar{\mathbf{x}}, \bar{\mathbf{y}}) + \lambda_1 \Phi_1(\bar{\mathbf{x}}) + \lambda_2 \Phi_2(\bar{\mathbf{y}}) \\ &< f(\hat{\mathbf{x}}, \hat{\mathbf{y}}) + \lambda_1 \Phi_1(\hat{\mathbf{x}}) + \lambda_2 \Phi_2(\hat{\mathbf{y}}) \\ &\leq f(\hat{\mathbf{x}}, \hat{\mathbf{y}}) + \lambda_1 \|\hat{\mathbf{x}}\|_0 + \lambda_2 \|\hat{\mathbf{y}}\|_{2,0}, \end{aligned}$$

which contradicts that $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$ is a global minimizer of problem (1.1). Therefore, $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$ must be a global minimizer of problem (1.2).

(ii). Let $(\hat{\mathbf{x}}, \hat{\mathbf{y}}) \in \mathbb{R}^{n+m}$ be a local minimizer of problem (1.2), then there exists a neighborhood W of $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$ such that

$$f(\hat{\mathbf{x}}, \hat{\mathbf{y}}) + \lambda_1 \Phi_1(\hat{\mathbf{x}}) + \lambda_2 \Phi_2(\hat{\mathbf{y}}) \leq f(\mathbf{x}, \mathbf{y}) + \lambda_1 \Phi_1(\mathbf{x}) + \lambda_2 \Phi_2(\mathbf{y}), \quad \forall (\mathbf{x}, \mathbf{y}) \in W,$$

It is easy to know that $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$ is also a d-stationary point of problem (1.2). From Theorem 3.8 and (4.1), we have

$$\Phi_1(\hat{\mathbf{x}}) = \|\hat{\mathbf{x}}\|_0 \quad \text{and} \quad \Phi_2(\hat{\mathbf{y}}) = \|\hat{\mathbf{y}}\|_{2,0}. \quad (4.2)$$

Hence, we have

$$\begin{aligned} f(\hat{\mathbf{x}}, \hat{\mathbf{y}}) + \lambda_1 \|\hat{\mathbf{x}}\|_0 + \lambda_2 \|\hat{\mathbf{y}}\|_{2,0} &= f(\hat{\mathbf{x}}, \hat{\mathbf{y}}) + \lambda_1 \Phi_1(\hat{\mathbf{x}}) + \lambda_2 \Phi_2(\hat{\mathbf{y}}) \\ &\leq f(\mathbf{x}, \mathbf{y}) + \lambda_1 \Phi_1(\mathbf{x}) + \lambda_2 \Phi_2(\mathbf{y}) \\ &\leq f(\mathbf{x}, \mathbf{y}) + \lambda_1 \|\mathbf{x}\|_0 + \lambda_2 \|\mathbf{y}\|_{2,0}, \quad \forall (\mathbf{x}, \mathbf{y}) \in W. \end{aligned}$$

Therefore, $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$ is a local minimizer of problem (1.1), and (4.2) implies that the objective function value of problems (1.1) and (1.2) at $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$ are equal. \square

Remark 4.2 (1) The result in Theorem (4.1) reveals that problems (1.1) and (1.2) have some equivalence, which provides a theoretical basis for solving problem (1.1) via solving problem (1.2).

(2) From Remark 3.9, we know that the hypothesis of Theorem (4.1) is easy to satisfy.

5 Alternating proximal gradient algorithm for problem (1.2)

In this section, we propose an APG algorithm to solve problem (1.2), and discuss the convergence of the sequence generated by the APG algorithm.

5.1 Scheme of APG algorithm

Noting that the objective function F in (1.2) has two parts of variables, the alternating minimization may be the suitable way to solve problem (1.2), which transforms problem (1.2) into two subproblems.

Take the initial point $(\mathbf{x}^0, \mathbf{y}^0) \in \mathbb{R}^{n+m}$, and let the sequence $\{(\mathbf{x}^{k+1}, \mathbf{y}^{k+1})\}_{k \in \mathbb{N}}$ be generated through the following subproblems:

$$\begin{cases} \mathbf{x}^{k+1} \in \arg \min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}^k, \mathbf{y}^k) + \langle \mathbf{x} - \mathbf{x}^k, \nabla_{\mathbf{x}} f(\mathbf{x}^k, \mathbf{y}^k) \rangle + \frac{1}{2t_1} \|\mathbf{x} - \mathbf{x}^k\|^2 + \lambda_1 \Phi_1(\mathbf{x}), & (5.1a) \\ \mathbf{y}^{k+1} \in \arg \min_{\mathbf{y} \in \mathbb{R}^m} f(\mathbf{x}^{k+1}, \mathbf{y}^k) + \langle \mathbf{y} - \mathbf{y}^k, \nabla_{\mathbf{y}} f(\mathbf{x}^{k+1}, \mathbf{y}^k) \rangle + \frac{1}{2t_2} \|\mathbf{y} - \mathbf{y}^k\|^2 + \lambda_2 \Phi_2(\mathbf{y}). & (5.1b) \end{cases}$$

274 One can note that the two subproblems in (5.1) are both nonconvex and nonsmooth
 275 since $\Phi_1(\mathbf{x})$ and $\Phi_2(\mathbf{y})$ are both nonconvex and nonsmooth. Fortunately, in the following
 276 part, we can provide their closed form solutions, which is very important for the efficiency
 277 of the APG algorithm.

278 The subproblem (5.1a) solves \mathbf{x} with the fixed \mathbf{y}^k . It can be explicitly reexpressed as
 279 the following form:

$$\mathbf{x}^{k+1} \in \arg \min_{\mathbf{x} \in \mathbb{R}^n} \left\{ \frac{1}{2t_1} \|\mathbf{x} - (\mathbf{x}^k - t_1 \nabla_{\mathbf{x}} f(\mathbf{x}^k, \mathbf{y}^k))\|^2 + \lambda_1 \Phi_1(\mathbf{x}) \right\}. \quad (5.2)$$

280 Denote $\mathbf{v}^k := \mathbf{x}^k - t_1 \nabla_{\mathbf{x}} f(\mathbf{x}^k, \mathbf{y}^k)$, then (5.2) can be rewritten as

$$\mathbf{x}^{k+1} \in \arg \min_{\mathbf{x} \in \mathbb{R}^n} \left\{ \frac{1}{2t_1} \|\mathbf{x} - \mathbf{v}^k\|^2 + \lambda_1 \Phi_1(\mathbf{x}) \right\}. \quad (5.3)$$

281 Note that $\Phi_1(\mathbf{x})$ is separable in the component of \mathbf{x} , then problem (5.3) is also separable.
 282 That is,

$$\mathbf{x}^{k+1} \in \text{Prox}_{t_1 \lambda_1 \Phi_1}(\mathbf{v}^k) = \text{Prox}_{t_1 \lambda_1 \varphi_1}(v_1^k) \times \cdots \times \text{Prox}_{t_1 \lambda_1 \varphi_1}(v_n^k), \quad (5.4)$$

283 where the proximal operator $\text{Prox}_{t_1 \lambda_1 \varphi_1}(\cdot)$ is the optimal solution of the following problem

$$\text{Prox}_{t_1 \lambda_1 \varphi_1}(v) = \arg \min_{x \in \mathbb{R}} \left\{ \frac{1}{2t_1} (x - v)^2 + \lambda_1 \varphi_1(x) \right\}, \quad \forall v \in \mathbb{R}. \quad (5.5)$$

284 The solution of (5.5) is known to have the following closed form [3, 19, 30, 43]

$$\begin{aligned} \text{Prox}_{t_1 \lambda_1 \varphi_1}(v) &= \begin{cases} 0, & |v| \leq \frac{\lambda_1 t_1}{\alpha_1}, \\ \text{sgn}(v)(|v| - \frac{\lambda_1 t_1}{\alpha_1}), & \frac{\lambda_1 t_1}{\alpha_1} < |v| < \alpha_1 + \frac{\lambda_1 t_1}{2\alpha_1}, \\ \text{sgn}(v)(\alpha_1 \pm \frac{\lambda_1 t_1}{2\alpha_1}), & |v| = \alpha_1 + \frac{\lambda_1 t_1}{2\alpha_1}, \\ v, & |v| > \alpha_1 + \frac{\lambda_1 t_1}{2\alpha_1}. \end{cases} \\ &= \begin{cases} \text{sgn}(v)(|v| - \frac{\lambda_1 t_1}{\alpha_1})_+, & |v| \leq \alpha_1 + \frac{\lambda_1 t_1}{2\alpha_1}, \\ v, & |v| \geq \alpha_1 + \frac{\lambda_1 t_1}{2\alpha_1}. \end{cases} \end{aligned} \quad (5.6)$$

285 which means that $\text{Prox}_{t_1 \lambda_1 \varphi_1}(v)$ has two values when $|v| = \alpha_1 + \frac{\lambda_1 t_1}{2\alpha_1}$.

286 The subproblem (5.1b) solves \mathbf{y} with the fixed \mathbf{x}^{k+1} . It can be explicitly reexpressed as

$$\mathbf{y}^{k+1} \in \arg \min_{\mathbf{y} \in \mathbb{R}^m} \left\{ \frac{1}{2t_2} \|\mathbf{y} - (\mathbf{y}^k - t_2 \nabla_{\mathbf{y}} f(\mathbf{x}^{k+1}, \mathbf{y}^k))\|^2 + \lambda_2 \Phi_2(\mathbf{y}) \right\}. \quad (5.7)$$

287 Denote $\mathbf{u}^k := \mathbf{y}^k - t_2 \nabla_{\mathbf{y}} f(\mathbf{x}^{k+1}, \mathbf{y}^k)$, then (5.7) can be simplified as

$$\mathbf{y}^{k+1} \in \arg \min_{\mathbf{y} \in \mathbb{R}^m} \left\{ \frac{1}{2t_2} \|\mathbf{y} - \mathbf{u}^k\|^2 + \lambda_2 \Phi_2(\mathbf{y}) \right\}. \quad (5.8)$$

288 Note that $\Phi_2(\mathbf{y})$ is separable in the group of \mathbf{y} , then problem (5.8) is also group separable.

289 That is, the solution of (5.8) have the following closed form

$$\mathbf{y}^{k+1} \in \text{Prox}_{t_2 \lambda_2 \Phi_2}(\mathbf{u}^k) = [\text{Prox}_{t_2 \lambda_2 \Phi_2}(\mathbf{u}^k)]_{(1)} \times \cdots \times [\text{Prox}_{t_2 \lambda_2 \Phi_2}(\mathbf{u}^k)]_{(J)} \quad (5.9)$$

290 with

$$[\text{Prox}_{t_2 \lambda_2 \Phi_2}(\mathbf{u})]_{(j)} = \begin{cases} (\|\mathbf{u}_{(j)}\| - \frac{\lambda_2 t_2}{\alpha_2})_+ + \frac{\mathbf{u}_{(j)}}{\|\mathbf{u}_{(j)}\|}, & \|\mathbf{u}_{(j)}\| \leq \alpha_2 + \frac{\lambda_2 t_2}{2\alpha_2}, \\ \mathbf{u}_{(j)}, & \|\mathbf{u}_{(j)}\| \geq \alpha_2 + \frac{\lambda_2 t_2}{2\alpha_2}, \end{cases}$$

291 for $j = 1, \dots, J$, which can be obtained by the similar way to (5.6) or [42].

292 From (5.4) and (5.9), we give the scheme of the APG algorithm for solving problem
293 (1.2) as below.

Algorithm 1 APG algorithm

– **Initialize:** For given $\alpha_1 > 0, \alpha_2 > 0, \lambda_1 > 0, \lambda_2 > 0, t_1 > 0, t_2 > 0, \text{xtol} > 0$, take $(\mathbf{x}^0, \mathbf{y}^0) \in \mathbb{R}^{n+m}$, and set $k = 0$.

– **Step1.** Compute

$$\begin{cases} \mathbf{x}^{k+1} \in \arg \min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}^k, \mathbf{y}^k) + \langle \mathbf{x} - \mathbf{x}^k, \nabla_{\mathbf{x}} f(\mathbf{x}^k, \mathbf{y}^k) \rangle + \frac{1}{2t_1} \|\mathbf{x} - \mathbf{x}^k\|^2 + \lambda_1 \Phi_1(\mathbf{x}), \\ \mathbf{y}^{k+1} \in \arg \min_{\mathbf{y} \in \mathbb{R}^m} f(\mathbf{x}^{k+1}, \mathbf{y}^k) + \langle \mathbf{y} - \mathbf{y}^k, \nabla_{\mathbf{y}} f(\mathbf{x}^{k+1}, \mathbf{y}^k) \rangle + \frac{1}{2t_2} \|\mathbf{y} - \mathbf{y}^k\|^2 + \lambda_2 \Phi_2(\mathbf{y}). \end{cases}$$

The calculation process is as follows:

I. Compute $\mathbf{x}^{k+1} \in \text{Prox}_{t_1 \lambda_1 \Phi_1}(\mathbf{x}^k - t_1 \nabla_{\mathbf{x}} f(\mathbf{x}^k, \mathbf{y}^k))$ according to (5.4);

II. Let $\mathbf{u}^k = \mathbf{y}^k - t_2 \nabla_{\mathbf{y}} f(\mathbf{x}^{k+1}, \mathbf{y}^k)$, then divide \mathbf{u}^k into J groups according to the given group of \mathbf{y} ;

III. Compute $\mathbf{y}^{k+1} \in \text{Prox}_{t_2 \lambda_2 \Phi_2}(\mathbf{u}^k)$ according to (5.9).

– **Step2.** Let $\mathbf{z}^{k+1} := (\mathbf{x}^{k+1}, \mathbf{y}^{k+1})$, if $\frac{\|\mathbf{z}^{k+1} - \mathbf{z}^k\|}{\max\{1, \|\mathbf{z}^{k+1}\|\}} \leq \text{xtol}$, terminate.

Otherwise, let $k := k + 1$ then return to **Step1**.

– **Output:** $(\mathbf{x}^k, \mathbf{y}^k)$

294 5.2 Convergence analysis

295 Before the convergence analysis of the APG algorithm, we give some basic assumptions.

296 **Assumption 5.1** (i) $\inf\{F(\mathbf{x}, \mathbf{y}) = f(\mathbf{x}, \mathbf{y}) + \lambda_1 \Phi_1(\mathbf{x}) + \lambda_2 \Phi_2(\mathbf{y})\} > -\infty$.

297 (ii) $f(\mathbf{x}, \mathbf{y}) \rightarrow \infty$ as $\|(\mathbf{x}, \mathbf{y})\| \rightarrow \infty$.

298 (iii) $\nabla_{\mathbf{x}} f(\cdot, \cdot)$ is Lipschitz continuous with modulus L_1 , that is

$$\|\nabla_{\mathbf{x}} f(\mathbf{x}^1, \mathbf{y}^1) - \nabla_{\mathbf{x}} f(\mathbf{x}^2, \mathbf{y}^2)\| \leq L_1(\|\mathbf{x}_1 - \mathbf{x}_2\| + \|\mathbf{y}_1 - \mathbf{y}_2\|), \quad \forall (\mathbf{x}^1, \mathbf{y}^1), (\mathbf{x}^2, \mathbf{y}^2) \in \mathbb{R}^{n+m}.$$

299 Meanwhile, for any \mathbf{x} , $\nabla_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})$ is Lipschitz continuous about \mathbf{y} with modulus L_2 .

300 (v) The parameters satisfy

$$0 < t_1 < \frac{1}{L_1}, \quad 0 < t_2 < \frac{1}{L_2}.$$

301 It is easy to check that there are many loss functions satisfy Assumption 5.1, for example,
302 ℓ_2 loss and logistic loss.

303 Next, we investigate the convergence of the proposed APG algorithm under Assumption
304 5.1.

305 **Lemma 5.2** Let $\{(\mathbf{x}^k, \mathbf{y}^k)\}_{k \in \mathbb{N}}$ be the sequence generated by the APG algorithm. Suppose
306 Assumption 5.1 holds, then

$$\rho (\|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2 + \|\mathbf{y}^{k+1} - \mathbf{y}^k\|^2) \leq F(\mathbf{x}^k, \mathbf{y}^k) - F(\mathbf{x}^{k+1}, \mathbf{y}^{k+1}), \quad (5.10)$$

307 where $\rho = \min\left\{\frac{1}{2t_1} - \frac{L_1}{2}, \frac{1}{2t_2} - \frac{L_2}{2}\right\} > 0$, which implies that the sequence $\{F(\mathbf{x}^k, \mathbf{y}^k)\}$ is
308 nonincreasing.

309 *Proof* From Step 1 in the APG algorithm, we know that

$$\begin{aligned} \lambda_1 \Phi_1(\mathbf{x}^k) + f(\mathbf{x}^k, \mathbf{y}^k) &\geq f(\mathbf{x}^k, \mathbf{y}^k) + \langle \nabla_{\mathbf{x}} f(\mathbf{x}^k, \mathbf{y}^k), \mathbf{x}^{k+1} - \mathbf{x}^k \rangle \\ &\quad + \frac{1}{2t_1} \|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2 + \lambda_1 \Phi_1(\mathbf{x}^{k+1}), \end{aligned} \quad (5.11)$$

310 and that

$$\begin{aligned} \lambda_2 \Phi_2(\mathbf{y}^k) + f(\mathbf{x}^{k+1}, \mathbf{y}^k) &\geq f(\mathbf{x}^{k+1}, \mathbf{y}^k) + \langle \nabla_{\mathbf{y}} f(\mathbf{x}^{k+1}, \mathbf{y}^k), \mathbf{y}^{k+1} - \mathbf{y}^k \rangle \\ &\quad + \frac{1}{2t_2} \|\mathbf{y}^{k+1} - \mathbf{y}^k\|^2 + \lambda_2 \Phi_2(\mathbf{y}^{k+1}). \end{aligned} \quad (5.12)$$

311 Summing (5.11) and (5.12), we obtain that

$$\begin{aligned} \lambda_1 \Phi_1(\mathbf{x}^k) + \lambda_2 \Phi_2(\mathbf{y}^k) &\geq \lambda_1 \Phi_1(\mathbf{x}^{k+1}) + \lambda_2 \Phi_2(\mathbf{y}^{k+1}) \\ &\quad + \langle \nabla_{\mathbf{x}} f(\mathbf{x}^k, \mathbf{y}^k), \mathbf{x}^{k+1} - \mathbf{x}^k \rangle + \frac{1}{2t_1} \|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2 \\ &\quad + \langle \nabla_{\mathbf{y}} f(\mathbf{x}^{k+1}, \mathbf{y}^k), \mathbf{y}^{k+1} - \mathbf{y}^k \rangle + \frac{1}{2t_2} \|\mathbf{y}^{k+1} - \mathbf{y}^k\|^2. \end{aligned} \quad (5.13)$$

312 From the Lipschitz continuity of $\nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y})$ and $\nabla_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})$ (Assumption 5.1 (iii)), we can
313 obtain

$$\begin{aligned} f(\mathbf{x}^k, \mathbf{y}^k) &\geq f(\mathbf{x}^{k+1}, \mathbf{y}^k) - \langle \nabla_{\mathbf{x}} f(\mathbf{x}^k, \mathbf{y}^k), \mathbf{x}^{k+1} - \mathbf{x}^k \rangle - \frac{L_1}{2} \|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2, \\ f(\mathbf{x}^{k+1}, \mathbf{y}^k) &\geq f(\mathbf{x}^{k+1}, \mathbf{y}^{k+1}) - \langle \nabla_{\mathbf{y}} f(\mathbf{x}^{k+1}, \mathbf{y}^k), \mathbf{y}^{k+1} - \mathbf{y}^k \rangle - \frac{L_2}{2} \|\mathbf{y}^{k+1} - \mathbf{y}^k\|^2. \end{aligned}$$

314 The above two inequalities yield that

$$\begin{aligned} f(\mathbf{x}^k, \mathbf{y}^k) &\geq f(\mathbf{x}^{k+1}, \mathbf{y}^{k+1}) - \langle \nabla_{\mathbf{x}} f(\mathbf{x}^k, \mathbf{y}^k), \mathbf{x}^{k+1} - \mathbf{x}^k \rangle - \frac{L_1}{2} \|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2, \\ &\quad - \langle \nabla_{\mathbf{y}} f(\mathbf{x}^{k+1}, \mathbf{y}^k), \mathbf{y}^{k+1} - \mathbf{y}^k \rangle - \frac{L_2}{2} \|\mathbf{y}^{k+1} - \mathbf{y}^k\|^2. \end{aligned} \quad (5.14)$$

315 Summing (5.13) and (5.14), we have

$$F(\mathbf{x}^k, \mathbf{y}^k) \geq F(\mathbf{x}^{k+1}, \mathbf{y}^{k+1}) + \left(\frac{1}{2t_1} - \frac{L_1}{2} \right) \|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2 + \left(\frac{1}{2t_2} - \frac{L_2}{2} \right) \|\mathbf{y}^{k+1} - \mathbf{y}^k\|^2.$$

316 By Assumption 5.1 (iii), we get $\frac{1}{2t_1} - \frac{L_1}{2} > 0$, $\frac{1}{2t_2} - \frac{L_2}{2} > 0$. Let $\rho = \min \left\{ \frac{1}{2t_1} - \frac{L_1}{2}, \frac{1}{2t_2} - \frac{L_2}{2} \right\}$,
317 then we obtain

$$\rho (\|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2 + \|\mathbf{y}^{k+1} - \mathbf{y}^k\|^2) \leq F(\mathbf{x}^k, \mathbf{y}^k) - F(\mathbf{x}^{k+1}, \mathbf{y}^{k+1}).$$

318 This completes the proof. \square

319 **Theorem 5.3** Suppose Assumption 5.1 holds. Let $\{\mathbf{z}^k := (\mathbf{x}^k, \mathbf{y}^k)\}$ be generated by the
320 APG Algorithm, then the following statements hold.

- 321 (i) $\{\mathbf{z}^k\}$ is bounded and $\{F(\mathbf{z}^k)\}$ is convergent;
322 (ii) $\sum_{k=0}^{\infty} \|\mathbf{z}^{k+1} - \mathbf{z}^k\|^2 < \infty$, $\lim_{k \rightarrow \infty} \|\mathbf{x}^{k+1} - \mathbf{x}^k\| = 0$ and $\lim_{k \rightarrow \infty} \|\mathbf{y}^{k+1} - \mathbf{y}^k\| = 0$.

323 *Proof* (i). From Lemma 5.2 and Assumption 5.1 (i), it follows that $\{F(\mathbf{x}^k, \mathbf{y}^k)\}$ is nonincreas-
 324 ing and F is bounded from below, and hence $\{F(\mathbf{x}^k, \mathbf{y}^k)\}$ is convergent. From $\{(\mathbf{x}^k, \mathbf{y}^k)\} \subset$
 325 $\{(\mathbf{x}, \mathbf{y}) : F(\mathbf{x}, \mathbf{y}) \leq F(\mathbf{x}^0, \mathbf{y}^0)\}$ which is bounded due to Assumption 5.1 (ii) and $\Phi_1(\mathbf{x}) \geq 0$
 326 as well as $\Phi_2(\mathbf{y}) \geq 0$, it follows that $\{\mathbf{x}^k, \mathbf{y}^k\}$ is bounded.

327 (ii). From (5.10) and (i), we have

$$\rho \|\mathbf{z}^{k+1} - \mathbf{z}^k\|^2 = \rho(\|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2 + \|\mathbf{y}^{k+1} - \mathbf{y}^k\|^2) \leq F(\mathbf{x}^k, \mathbf{y}^k) - F(\mathbf{x}^{k+1}, \mathbf{y}^{k+1}).$$

328 Summing both sides of the above inequality from 0 to N , we get

$$\begin{aligned} \sum_{k=0}^N \rho \|\mathbf{z}^{k+1} - \mathbf{z}^k\|^2 &= \sum_{k=0}^N \rho(\|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2 + \|\mathbf{y}^{k+1} - \mathbf{y}^k\|^2) \\ &\leq \sum_{k=0}^N (F(\mathbf{x}^k, \mathbf{y}^k) - F(\mathbf{x}^{k+1}, \mathbf{y}^{k+1})) \\ &= F(\mathbf{x}^0, \mathbf{y}^0) - F(\mathbf{x}^{N+1}, \mathbf{y}^{N+1}), \end{aligned}$$

329 Letting $N \rightarrow \infty$, we obtain

$$\sum_{k=0}^{\infty} \rho \|\mathbf{z}^{k+1} - \mathbf{z}^k\|^2 = \sum_{k=0}^{\infty} \rho(\|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2 + \|\mathbf{y}^{k+1} - \mathbf{y}^k\|^2) \leq F(\mathbf{x}^0, \mathbf{y}^0) - F(\mathbf{x}^*, \mathbf{y}^*) < \infty,$$

330 Then

$$\lim_{k \rightarrow \infty} \|\mathbf{z}^{k+1} - \mathbf{z}^k\|^2 = \lim_{k \rightarrow \infty} (\|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2 + \|\mathbf{y}^{k+1} - \mathbf{y}^k\|^2) = 0.$$

331 Thus, $\lim_{k \rightarrow \infty} \|\mathbf{x}^{k+1} - \mathbf{x}^k\| = 0$ and $\lim_{k \rightarrow \infty} \|\mathbf{y}^{k+1} - \mathbf{y}^k\| = 0$. \square

332 In order to prove a global convergence of the whole sequence $\{(\mathbf{x}^k, \mathbf{y}^k)\}$, we first prove
 333 the following results.

Lemma 5.4 *Suppose Assumption 5.1 holds, and $\{(\mathbf{x}^k, \mathbf{y}^k)\}_{k \in \mathbb{N}}$ is generated by the APG algorithm with the initial point $(\mathbf{x}^0, \mathbf{y}^0)$. Let*

$$q_{\mathbf{x}}^{k+1} := \nabla_{\mathbf{x}} f(\mathbf{x}^{k+1}, \mathbf{y}^{k+1}) - \nabla_{\mathbf{x}} f(\mathbf{x}^k, \mathbf{y}^k) - \frac{1}{t_1}(\mathbf{x}^{k+1} - \mathbf{x}^k),$$

$$q_{\mathbf{y}}^{k+1} := \nabla_{\mathbf{y}} f(\mathbf{x}^{k+1}, \mathbf{y}^{k+1}) - \nabla_{\mathbf{y}} f(\mathbf{x}^{k+1}, \mathbf{y}^k) - \frac{1}{t_2}(\mathbf{y}^{k+1} - \mathbf{y}^k),$$

334 then

$$q_{\mathbf{x}}^{k+1} \in \partial_{\mathbf{x}} F(\mathbf{x}^{k+1}, \mathbf{y}^{k+1}), \quad q_{\mathbf{y}}^{k+1} \in \partial_{\mathbf{y}} F(\mathbf{x}^{k+1}, \mathbf{y}^{k+1})$$

335 and

$$\|(q_{\mathbf{x}}^{k+1}, q_{\mathbf{y}}^{k+1})\| \leq \alpha \|(\mathbf{x}^{k+1}, \mathbf{y}^{k+1})\|,$$

336 where $\alpha^2 = \max \left\{ 2(L_1 + \frac{1}{t_1})^2, 2L_1^2 + (L_2 + \frac{1}{t_2})^2 \right\}$.

337 *Proof* It follows from (5.2) that

$$0 \in \frac{1}{t_1}(\mathbf{x}^{k+1} - \mathbf{x}^k) + \nabla_{\mathbf{x}}f(\mathbf{x}^k, \mathbf{y}^k) + \lambda_1 \partial \Phi_1(\mathbf{x}^{k+1}), \quad (5.15)$$

338 Adding $\nabla_{\mathbf{x}}f(\mathbf{x}^{k+1}, \mathbf{y}^{k+1})$ to both sides of (5.15) and rearranging terms, we obtain

$$\begin{aligned} \nabla_{\mathbf{x}}f(\mathbf{x}^{k+1}, \mathbf{y}^{k+1}) - \nabla_{\mathbf{x}}f(\mathbf{x}^k, \mathbf{y}^k) - \frac{1}{t_1}(\mathbf{x}^{k+1} - \mathbf{x}^k) &\in \nabla_{\mathbf{x}}f(\mathbf{x}^{k+1}, \mathbf{y}^{k+1}) + \lambda_1 \partial \Phi_1(\mathbf{x}^{k+1}) \\ &= \partial_{\mathbf{x}}F(\mathbf{x}^{k+1}, \mathbf{y}^{k+1}). \end{aligned} \quad (5.16)$$

339 Similarly, it follows from (5.7) that

$$0 \in \frac{1}{t_2}(\mathbf{y}^{k+1} - \mathbf{y}^k) + \nabla_{\mathbf{y}}f(\mathbf{x}^{k+1}, \mathbf{y}^k) + \lambda_2 \partial \Phi_2(\mathbf{y}^{k+1}). \quad (5.17)$$

340 Adding $\nabla_{\mathbf{y}}f(\mathbf{x}^{k+1}, \mathbf{y}^{k+1})$ to both sides of (5.17) and rearranging terms, we obtain

$$\begin{aligned} \nabla_{\mathbf{y}}f(\mathbf{x}^{k+1}, \mathbf{y}^{k+1}) - \nabla_{\mathbf{y}}f(\mathbf{x}^{k+1}, \mathbf{y}^k) - \frac{1}{t_2}(\mathbf{y}^{k+1} - \mathbf{y}^k) &\in \nabla_{\mathbf{y}}f(\mathbf{x}^{k+1}, \mathbf{y}^{k+1}) + \lambda_2 \partial \Phi_2(\mathbf{y}^{k+1}) \\ &= \partial_{\mathbf{y}}F(\mathbf{x}^{k+1}, \mathbf{y}^{k+1}). \end{aligned} \quad (5.18)$$

341 Combining (5.16) and (5.18), we obtain

$$q_{\mathbf{x}}^{k+1} \in \partial_{\mathbf{x}}F(\mathbf{x}^{k+1}, \mathbf{y}^{k+1}), \quad q_{\mathbf{y}}^{k+1} \in \partial_{\mathbf{y}}F(\mathbf{x}^{k+1}, \mathbf{y}^{k+1}), \quad (5.19)$$

342 which then implies $(q_{\mathbf{x}}^{k+1}, q_{\mathbf{y}}^{k+1}) \in \partial F(\mathbf{x}^{k+1}, \mathbf{y}^{k+1})$.

343 From (5.19) and Assumption 5.1 (iii), we have

$$\begin{aligned} \|q_{\mathbf{x}}^{k+1}\| &= \|\nabla_{\mathbf{x}}f(\mathbf{x}^{k+1}, \mathbf{y}^{k+1}) - \nabla_{\mathbf{x}}f(\mathbf{x}^k, \mathbf{y}^k) - \frac{1}{t_1}(\mathbf{x}^{k+1} - \mathbf{x}^k)\| \\ &\leq \|\nabla_{\mathbf{x}}f(\mathbf{x}^{k+1}, \mathbf{y}^{k+1}) - \nabla_{\mathbf{x}}f(\mathbf{x}^k, \mathbf{y}^k)\| + \frac{1}{t_1}\|\mathbf{x}^{k+1} - \mathbf{x}^k\| \\ &\leq \left(L_1 + \frac{1}{t_1}\right)\|\mathbf{x}^{k+1} - \mathbf{x}^k\| + L_1\|\mathbf{y}^{k+1} - \mathbf{y}^k\|. \end{aligned}$$

344 Similarly, we have

$$\begin{aligned} \|q_{\mathbf{y}}^{k+1}\| &= \|\nabla_{\mathbf{y}}f(\mathbf{x}^{k+1}, \mathbf{y}^{k+1}) - \nabla_{\mathbf{y}}f(\mathbf{x}^{k+1}, \mathbf{y}^k) - \frac{1}{t_2}(\mathbf{y}^{k+1} - \mathbf{y}^k)\| \\ &\leq \|\nabla_{\mathbf{y}}f(\mathbf{x}^{k+1}, \mathbf{y}^{k+1}) - \nabla_{\mathbf{y}}f(\mathbf{x}^{k+1}, \mathbf{y}^k)\| + \frac{1}{t_2}\|\mathbf{y}^{k+1} - \mathbf{y}^k\| \\ &\leq \left(L_2 + \frac{1}{t_2}\right)\|\mathbf{y}^{k+1} - \mathbf{y}^k\|. \end{aligned}$$

345 then

$$\begin{aligned} \|(q_{\mathbf{x}}^{k+1}, q_{\mathbf{y}}^{k+1})\|^2 &= \|q_{\mathbf{x}}^{k+1}\|^2 + \|q_{\mathbf{y}}^{k+1}\|^2 \\ &\leq 2\left(L_1 + \frac{1}{t_1}\right)^2\|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2 + \left(2L_1^2 + \left(L_2 + \frac{1}{t_2}\right)^2\right)\|\mathbf{y}^{k+1} - \mathbf{y}^k\|^2. \end{aligned}$$

346 Let $\alpha^2 = \max\{2\left(L_1 + \frac{1}{t_1}\right)^2, 2L_1^2 + \left(L_2 + \frac{1}{t_2}\right)^2\}$, then we have

$$\|(q_{\mathbf{x}}^{k+1}, q_{\mathbf{y}}^{k+1})\|^2 \leq \alpha^2(\|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2 + \|\mathbf{y}^{k+1} - \mathbf{y}^k\|^2).$$

347 As a consequence, we get

$$\|(q_{\mathbf{x}}^{k+1}, q_{\mathbf{y}}^{k+1})\| \leq \alpha\|\mathbf{x}^{k+1} - \mathbf{x}^k, \mathbf{y}^{k+1} - \mathbf{y}^k\|.$$

348 The proof is thus complete. \square

349 To analyze the convergence of the generated sequence of the APG algorithm, we discuss
 350 some properties of the limit point sets of the sequence at first. For convenience, we denote
 351 $\mathbf{z}^k = (\mathbf{x}^k, \mathbf{y}^k)$ and $F(\mathbf{z}^k) = F(\mathbf{x}^k, \mathbf{y}^k)$. The set of all limit points of $\{\mathbf{z}^k\}$ is denoted by $\Gamma(\mathbf{z}^0)$,
 352 i.e.,

$$\Gamma(\mathbf{z}^0) = \{\bar{\mathbf{z}} = (\bar{\mathbf{x}}, \bar{\mathbf{y}}) \in \mathbb{R}^{n+m} \mid \exists \{k_j\} \in \mathbb{N}, \text{ s.t. } \mathbf{z}^{k_j} \rightarrow \bar{\mathbf{z}}, j \rightarrow \infty\}.$$

353 **Theorem 5.5** *Suppose Assumption 5.1 holds. Let $\{\mathbf{z}^k\}$ be generated by the APG algo-*
 354 *rithm with the initial point $\mathbf{z}^0 = (\mathbf{x}^0, \mathbf{y}^0)$, then the following statements hold.*

355 (i) $\Gamma(\mathbf{z}^0)$ is a nonempty and compact set, and the objective value of F is finite and
 356 constant on $\Gamma(\mathbf{z}^0)$.

357 (ii) $\Gamma(\mathbf{z}^0) \subset \text{crit}F$.

358 (iii) $\lim_{k \rightarrow \infty} \text{dist}(\mathbf{z}^k, \Gamma(\mathbf{z}^0)) = 0$.

359 *Proof* (i). From the boundedness of $\{\mathbf{z}^k\}$, it follows that $\Gamma(\mathbf{z}^0)$ is nonempty. Note that $\Gamma(\mathbf{z}^0)$
 360 can be represented as an intersection of compact sets, i.e.,

$$\Gamma(\mathbf{z}^0) = \bigcap_{s \in \mathbb{N}} \overline{\bigcup_{k \geq s} \{\mathbf{z}^k\}}.$$

361 Since the intersection of bounded closed sets is still bounded and closed, $\Gamma(\mathbf{z}^0)$ is also a
 362 compact set. For any $\bar{\mathbf{z}} = (\bar{\mathbf{x}}, \bar{\mathbf{y}}) \in \Gamma(\mathbf{z}^0)$, $\exists \{k_j\} \subset \mathbb{N}$ such that

$$\lim_{j \rightarrow \infty} \mathbf{z}^{k_j} = \bar{\mathbf{z}}.$$

363 By the continuity of F , we have

$$\lim_{j \rightarrow \infty} F(\mathbf{z}^{k_j}) = F(\bar{\mathbf{z}}).$$

364 From Theorem 5.3 (i), we have $F(\mathbf{z}^k) \rightarrow F^*$ ($k \rightarrow \infty$). Then, for arbitrary subsequence
 365 $F(\mathbf{z}^{k_j})$, it holds

$$\lim_{j \rightarrow \infty} F(\mathbf{z}^{k_j}) = F(\bar{\mathbf{z}}) = F^*. \quad (5.20)$$

366 That is, the value of F on $\Gamma(\mathbf{z}^0)$ is a constant.

367 (ii) From Theorem 5.3 (ii), we have

$$\lim_{j \rightarrow \infty} \|\mathbf{x}^{k_j+1} - \mathbf{x}^{k_j}\| = 0, \quad \lim_{j \rightarrow \infty} \|\mathbf{y}^{k_j+1} - \mathbf{y}^{k_j}\| = 0,$$

368 then

$$\lim_{j \rightarrow \infty} \mathbf{x}^{k_j+1} = \lim_{j \rightarrow \infty} \mathbf{x}^{k_j} = \bar{\mathbf{x}}, \quad \lim_{j \rightarrow \infty} \mathbf{y}^{k_j+1} = \lim_{j \rightarrow \infty} \mathbf{y}^{k_j} = \bar{\mathbf{y}}.$$

369 From Lemma 5.4, we have

$$\|(q_{\mathbf{x}}^{k_j+1}, q_{\mathbf{y}}^{k_j+1})\| \leq \alpha \|(\mathbf{x}^{k_j+1} - \mathbf{x}^{k_j}, \mathbf{y}^{k_j+1} - \mathbf{y}^{k_j})\|.$$

370 Let $j \rightarrow \infty$, then

$$\lim_{j \rightarrow \infty} \|(q_{\mathbf{x}}^{k_j+1}, q_{\mathbf{y}}^{k_j+1})\| = 0, \quad \text{i.e.,} \quad (q_{\mathbf{x}}^{k_j+1}, q_{\mathbf{y}}^{k_j+1}) \rightarrow (0, 0), \text{ for } j \rightarrow \infty.$$

371 From Lemma 5.4, we know $(q_{\mathbf{x}}^{k_j+1}, q_{\mathbf{y}}^{k_j+1}) \in \partial F(\mathbf{x}^{k_j+1}, \mathbf{y}^{k_j+1}) \subset \partial^C F(\mathbf{x}^{k_j+1}, \mathbf{y}^{k_j+1})$. Further,
 372 by the closedness of the mapping $\partial^C F(\cdot)$ [13, Propostion 2.1.5(b)], we obtain

$$(0, 0) \in \partial^C F(\bar{\mathbf{x}}, \bar{\mathbf{y}}).$$

373 From Remark 3.6, this implies that $\bar{\mathbf{z}} = (\bar{\mathbf{x}}, \bar{\mathbf{y}})$ is a critical point of problem (1.2), and
 374 $\Gamma(\mathbf{z}^0) \subset \text{crit}F$.

375 (iii) This conclusion follows from the definition of $\Gamma(\mathbf{z}^0)$. \square

376 In order to give the global convergence of the whole sequence $\{(\mathbf{x}^k, \mathbf{y}^k)\}$, we first intro-
 377 duce the Kurdyka-Łojasiewicz (KL) property of F . The KL property was used to analyze
 378 smooth problems, then Bolte, Daniilidis and Lewis [6] used KL property to analyze nons-
 379 mooth problems. Since then, lots of researchers have done much research on this basis, for
 380 example, [1, 2, 7, 27]. Now, it is well-known that the KL property have played the important
 381 roles in the convergence analysis of proximal algorithms. Let's recall the KL property.

382 Let $\eta \in (0, +\infty]$, we denote by Ψ_η the class of all concave and continuous functions
 383 $\psi : [0, \eta) \rightarrow [0, \infty)$ such that

- 384 (i) $\psi(0) = 0$;
- 385 (ii) ψ is continuously differentiable on $(0, \eta)$;
- 386 (iii) $\psi'(s) > 0$ for all $s \in (0, \eta)$.

387 **Definition 5.6** [2, 7, 27] [KL property] Let $h : \mathbb{R}^{n+m} \rightarrow \mathbb{R} \cup \{\infty\}$ be a proper lower
 388 semicontinuous function.

389 (i) h is said to have the KL property at $\bar{\mathbf{w}} \in \text{dom}\partial h := \{\mathbf{w} \in \mathbb{R}^{n+m} \mid \partial h(\mathbf{w}) \neq \emptyset\}$, if there
 390 exist $\eta \in (0, +\infty]$, a neighborhood Ω of $\bar{\mathbf{w}}$ and a function $\psi \in \Psi_\eta$, such that for all

$$\mathbf{w} \in \Omega \cap [h(\bar{\mathbf{w}}) < h(\mathbf{w}) < h(\bar{\mathbf{w}}) + \eta],$$

391 the following inequality holds

$$\psi'(h(\mathbf{w}) - h(\bar{\mathbf{w}}))\text{dist}(0, \partial h(\mathbf{w})) \geq 1.$$

392 (ii) If h satisfies the KL property at each point of $\text{dom}\partial h$, then h is called a KL function.

393 **Lemma 5.7** [7, 27] [Uniformized KL property] Let Ω be a compact set and $h : \mathbb{R}^{n+m} \rightarrow$
 394 $\mathbb{R} \cup \{\infty\}$ be a proper lower semicontinuous function. Assume that h is constant on Ω and
 395 satisfies the KL property at each point of Ω . Then there exist $\epsilon > 0, \eta > 0$ and $\psi \in \Psi_\eta$ such
 396 that for all $\bar{\mathbf{w}}$ in Ω and all

$$\mathbf{w} \in \{\mathbf{w} \in \mathbb{R}^{n+m} : \text{dist}(\mathbf{w}, \Omega) < \epsilon\} \cap [h(\bar{\mathbf{w}}) < h(\mathbf{w}) < h(\bar{\mathbf{w}}) + \eta],$$

397 one has,

$$\psi'(h(\mathbf{w}) - h(\bar{\mathbf{w}}))\text{dist}(0, \partial h(\mathbf{w})) \geq 1.$$

398 The KL functions have a wide range including semi-algebraic, subanalytic and log-exp
 399 and so on [7]. It is easy to check that our objective function F in (1.2) meets the KL property.

400 Now we can give the global convergence of the whole sequence $\{(\mathbf{x}^k, \mathbf{y}^k)\}$ under the
 401 condition of KL function.

402 **Theorem 5.8** Suppose Assumption 5.1 holds and F is a KL function. Let $\{\mathbf{z}^k = (\mathbf{x}^k, \mathbf{y}^k)\}$
 403 be generated by the APG algrithom. Then the following statements hold.

- 404 (i) $\sum_{k=0}^{\infty} \|\mathbf{z}^{k+1} - \mathbf{z}^k\| < \infty$;
- 405 (ii) The sequence $\{\mathbf{z}^k\}_{k \in \mathbb{N}}$ converges to a critical point $\mathbf{z}^* = (\mathbf{x}^*, \mathbf{y}^*)$ of problem (1.2).

406 *Proof* (i). Firstly, we suppose that $F(\mathbf{z}^k) \neq F(\bar{\mathbf{z}})$ for all $k \in \mathbb{N}$; Otherwise, the algorithm will
407 terminate.

408 On the one hand, it follows from (5.20) that $\lim_{k \rightarrow \infty} F(\mathbf{z}^k) = F^* = F(\bar{\mathbf{z}})$. Then, for any
409 $\eta > 0$, there exists $k_0 > 0$, such that for any $k > k_0$, it holds

$$F(\bar{\mathbf{z}}) < F(\mathbf{z}^k) < F(\bar{\mathbf{z}}) + \eta,$$

410 that is,

$$\mathbf{z}^k \in [F(\bar{\mathbf{z}}) < F(\mathbf{z}) < F(\bar{\mathbf{z}}) + \eta], \quad \forall k > k_0.$$

411 On the other hand, by Theorem 5.5 (iii), we have $\lim_{k \rightarrow \infty} \text{dist}(\mathbf{z}^k, \Gamma(\mathbf{z}^0)) = 0$. Therefore, for
412 any $\epsilon > 0$, there exists $k_1 > 0$, such that for any $k > k_1$, it holds

$$\text{dist}(\mathbf{z}^k, \Gamma(\mathbf{z}^0)) < \epsilon.$$

413 Let $k_2 = \max\{k_0, k_1\}$, then for any $k > k_2$, we have

$$\mathbf{z}^k \in \{\mathbf{z} \mid \text{dist}(\mathbf{z}, \Gamma(\mathbf{z}^0)) < \epsilon\} \cap [F(\bar{\mathbf{z}}) < F(\mathbf{z}) < F(\bar{\mathbf{z}}) + \eta], \quad \forall k > k_2.$$

414 Since the value of F on $\Gamma(\mathbf{z}^0)$ is a constant, by the uniformized KL property (Lemma
415 5.7), there exists $\psi \in \Psi_\eta$, such that

$$\psi'(F(\mathbf{z}^k) - F(\bar{\mathbf{z}})) \text{dist}(0, \partial F(\mathbf{z}^k)) \geq 1. \quad (5.21)$$

416 By Lemma 5.4, we have $(q_{\mathbf{x}}^{k+1}, q_{\mathbf{y}}^{k+1}) \in \partial F(\mathbf{z}^{k+1})$ and $\|(q_{\mathbf{x}}^{k+1}, q_{\mathbf{y}}^{k+1})\| \leq \alpha \|(\mathbf{x}^{k+1} - \mathbf{x}^k, \mathbf{y}^{k+1} - \mathbf{y}^k)\|$, then
417

$$\text{dist}(0, \partial F(\mathbf{z}^k)) \leq \|(q_{\mathbf{x}}^{k+1}, q_{\mathbf{y}}^{k+1})\| \leq \alpha \|(\mathbf{x}^{k+1} - \mathbf{x}^k, \mathbf{y}^{k+1} - \mathbf{y}^k)\| = \alpha \|\mathbf{z}^{k+1} - \mathbf{z}^k\|.$$

418 Substitute the above inequality into (5.21), then we obtain

$$\psi'(F(\mathbf{z}^k) - F(\bar{\mathbf{z}})) \geq \frac{1}{\text{dist}(0, \partial F(\mathbf{z}^k))} \geq \frac{1}{\alpha \|\mathbf{z}^{k+1} - \mathbf{z}^k\|}.$$

419 Since ψ is a concave function, we have

$$\psi(F(\mathbf{z}^{k+1}) - F(\bar{\mathbf{z}})) \leq \psi(F(\mathbf{z}^k) - F(\bar{\mathbf{z}})) + \psi'(F(\mathbf{z}^k) - F(\bar{\mathbf{z}}))(F(\mathbf{z}^{k+1}) - F(\mathbf{z}^k)).$$

420 Due to the above two inequalities and the sufficient descending property of function F given
421 by Lemma 5.2, we get

$$\begin{aligned} & \psi(F(\mathbf{z}^k) - F(\bar{\mathbf{z}})) - \psi(F(\mathbf{z}^{k+1}) - F(\bar{\mathbf{z}})) \\ & \geq \psi'(F(\mathbf{z}^k) - F(\bar{\mathbf{z}}))(F(\mathbf{z}^k) - F(\mathbf{z}^{k+1})) \\ & \geq \frac{\rho \|\mathbf{z}^{k+1} - \mathbf{z}^k\|^2}{\alpha \|\mathbf{z}^k - \mathbf{z}^{k-1}\|}. \end{aligned}$$

422 Denote by $C = \alpha/\rho$ and $\Delta_k = \psi(F(\mathbf{z}^k) - F(\bar{\mathbf{z}}))$, then Δ_k is monotonically non-increasing
423 with respect to k , and $\bar{\Delta} = \lim_{k \rightarrow \infty} \Delta_k$ makes sense. Therefore, the above inequality can be
424 rewritten as

$$\Delta_k - \Delta_{k+1} \geq \frac{\|\mathbf{z}^{k+1} - \mathbf{z}^k\|^2}{C \|\mathbf{z}^k - \mathbf{z}^{k-1}\|}.$$

425 By the inequality $4ab \leq (a+b)^2$, then

$$\begin{aligned} \|\mathbf{z}^{k+1} - \mathbf{z}^k\|^2 &\leq C\|\mathbf{z}^k - \mathbf{z}^{k-1}\|(\Delta_k - \Delta_{k+1}) \\ &\leq \left(\frac{\|\mathbf{z}^k - \mathbf{z}^{k-1}\| + C(\Delta_k - \Delta_{k+1})}{2} \right)^2, \end{aligned}$$

426 hence,

$$2\|\mathbf{z}^{k+1} - \mathbf{z}^k\| \leq \|\mathbf{z}^k - \mathbf{z}^{k-1}\| + C(\Delta_k - \Delta_{k+1}).$$

427 Summing the left and right sides of the above inequality respecting to k , we obtain

$$\begin{aligned} 2 \sum_{k=k_2+1}^K \|\mathbf{z}^{k+1} - \mathbf{z}^k\| &\leq C(\Delta_{k_2+1} - \Delta_{K+1}) + \sum_{k=k_2+1}^K \|\mathbf{z}^k - \mathbf{z}^{k-1}\| \\ &= C(\Delta_{k_2+1} - \Delta_{K+1}) + \|\mathbf{z}^{k_2+1} - \mathbf{z}^{k_2}\| \\ &\quad - \|\mathbf{z}^{K+1} - \mathbf{z}^K\| + \sum_{k=k_2+1}^K \|\mathbf{z}^{k+1} - \mathbf{z}^k\|, \end{aligned}$$

428 then

$$\sum_{k=k_2+1}^K \|\mathbf{z}^{k+1} - \mathbf{z}^k\| \leq C(\Delta_{k_2+1} - \Delta_{K+1}) + \|\mathbf{z}^{k_2+1} - \mathbf{z}^{k_2}\| - \|\mathbf{z}^{K+1} - \mathbf{z}^K\|.$$

429 Letting $K \rightarrow \infty$, we obtain

$$\sum_{k=k_2+1}^{\infty} \|\mathbf{z}^{k+1} - \mathbf{z}^k\| \leq C(\Delta_{k_2+1} - \bar{\Delta}) + \|\mathbf{z}^{k_2+1} - \mathbf{z}^{k_2}\| < \infty.$$

430 Therefore,

$$\sum_{k=0}^{\infty} \|\mathbf{z}^{k+1} - \mathbf{z}^k\| = \sum_{k=0}^{k_2} \|\mathbf{z}^{k+1} - \mathbf{z}^k\| + \sum_{k=k_2+1}^{\infty} \|\mathbf{z}^{k+1} - \mathbf{z}^k\| < \infty.$$

431 (ii) For any $p > q \geq k_2$, we have

$$\|\mathbf{z}^p - \mathbf{z}^q\| = \left\| \sum_{k=q}^{p-1} (\mathbf{z}^{k+1} - \mathbf{z}^k) \right\| \leq \sum_{k=q}^{p-1} \|\mathbf{z}^{k+1} - \mathbf{z}^k\| < \sum_{k=q}^{\infty} \|\mathbf{z}^{k+1} - \mathbf{z}^k\|.$$

432 Then $\|\mathbf{z}^p - \mathbf{z}^q\| \rightarrow 0$ as $q \rightarrow \infty$, which indicates that $\{\mathbf{z}^k\}$ is a Cauchy sequence, and hence is
433 a convergent sequence. It then follows from Theorem 5.5 that the limit point $\mathbf{z}^* = (\mathbf{x}^*, \mathbf{y}^*)$
434 of $\{\mathbf{z}^k\}$ is a critical point of problem (1.2). The proof is thus complete. \square

435 6 Numerical experiments

436 In this section, we conduct numerical experiments on the relaxation problem (1.2) to
437 test the APG algorithm.

438 All the numerical experiments are implemented in MATLAB R2018b and on a Lenovo
439 PC (Intel(R) Core(TM) i5-9500, 3.00GHz, 8.00GB of RAM).

440 6.1 Simulated Data

441 In this simulation experiment part, the APG algorithm is applied to solve the following
442 model.

443 **Example 6.1** We consider the least square loss $f(\mathbf{x}, \mathbf{y}) = \frac{1}{2} \|\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y} - \mathbf{c}\|^2$, that is

$$\min_{\mathbf{x} \in \mathbb{R}^n, \mathbf{y} \in \mathbb{R}^m} \frac{1}{2} \|\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y} - \mathbf{c}\|^2 + \lambda_1 \Phi_1(\mathbf{x}) + \lambda_2 \Phi_2(\mathbf{y}), \quad (6.1)$$

444 where $A \in \mathbb{R}^{p \times n}$, $B \in \mathbb{R}^{p \times m}$,

$$\Phi_1(\mathbf{x}) := \sum_{i=1}^n \varphi_1(|x_i|), \quad \Phi_2(\mathbf{y}) := \sum_{j=1}^J \varphi_2(\|\mathbf{y}_{(j)}\|)$$

445 and φ_i ($i = 1, 2$) is defined in (1.2).

446 For this model, the data are generated as follows. We first use MATLAB codes `randn(p, n)`
447 and `randn(p, m)` to randomly generate the i.i.d. Gaussian matrices $A \in \mathbb{R}^{p \times n}$ and $B \in \mathbb{R}^{p \times m}$.
448 Then we generate a sparse solution $\mathbf{x}_0 \in \mathbb{R}^n$ and a group sparse solution $\mathbf{y}_0 \in \mathbb{R}^m$ as the
449 real solution. Let kkx be the number of non-zero entries of \mathbf{x}_0 , then the sparsity level of \mathbf{x}_0
450 is kkx/n . Meanwhile, $\mathbf{y}_0 \in \mathbb{R}^m$ is randomly divided into J groups. The kkj non-zero groups
451 are randomly selected from these J groups, and the remaining $J - kkj$ groups are all set to
452 be zero vectors, so the group sparsity level of \mathbf{y}_0 is kkj/J .

453 For the given positive integers p, n, m, J, kkx, kkj , the real solution $\mathbf{z}_0 = (\mathbf{x}_0, \mathbf{y}_0)$ are
454 generated by the following codes:

```
455   xo =zeros(n,1); Indx =randperm(n); xo(Indx(1:kkx)) =randn(kkx,1);
456   avgsz =floor(m/J); idy =[]; gidy =[gidy; j*ones(avgsz,1)], j =1:J;
457   qqy =randperm(J); suppy =sort(qqy(1:kkj)); yo =zeros(m,1);
458   idy =find(gidy ==suppy(k)), yo(idy) =randn(avgsz,1), k =1:kkj;
459   zo =[xo;yo];
```

460 The observed data $\mathbf{c} \in \mathbb{R}^p$ is generated by

$$\mathbf{c} = \mathbf{A} * \mathbf{x}_0 + \mathbf{B} * \mathbf{y}_0 + \sigma * \text{randn}(p, 1),$$

461 where σ is the standard deviation of additive Gaussian noise.

462 The parameters and initial values in the APG algorithm are given as follows: $\mathbf{z}^0 =$
463 $(\mathbf{x}^0, \mathbf{y}^0) = \mathbf{0}_{n+m}$, $t_1 = 0.7$, $t_2 = 0.9$. In each iteration of the APG algorithm, we sort
464 $\mathbf{E}_x = \{|x_i^k|\}_{i \in [n]}$ and $\mathbf{E}_y = \{\|\mathbf{y}_{(j)}^k\|\}_{j \in [J]}$ in ascending order, we take $\text{crix} = \mathbf{E}_x_{n-kkx}$,
465 $\text{criy} = \mathbf{E}_y_{(m-kkj)}$, $\alpha_1 = 1.2 * \text{crix}$, $\alpha_2 = 1.8 * \text{criy}$, $\lambda_1 = \text{crix} * \alpha_1 / t_1$, and $\lambda_2 = \text{criy} * \alpha_2 / t_2$.
466 Let $\mathbf{z}^* = (\mathbf{x}^*, \mathbf{y}^*) \in \mathbb{R}^{n+m}$ denote the solution produced by the APG algorithm.

467 In this example, for each set of given numbers $\{p, n, m, J = m/4, kkx, kkj, \sigma\}$, we run
468 100 instances and use three indicators to evaluate the experimental effect of the proposed
469 APG algorithm: average relative error ($\text{Rel-err} := \frac{\|\mathbf{z}^* - \mathbf{z}_0\|}{\max\{1, \|\mathbf{z}_0\|\}}$), average CPU time and suc-
470 cessful rate (Suc-rat) where $\text{Rel-err} < 10^{-2}$ is regarded success. Set $\text{xtol} = 10^{-4}$. The ex-
471 perimental results are shown in Table 1, where we consider two cases: noiseless $\sigma = 0$ and
472 noised $\sigma = 10^{-3}$.

473 In Figure 1, the scatter plots of real and numerical solutions for $p = 2000$, $n = 3000$,
474 $m = 4000$, $kkx = 50$, $kkj = 20$, $J = 1000$ are displayed.

Table 1: Average numerical results of the APG algorithm

Problem						$\sigma = 0$			$\sigma = 10^{-3}$		
p	n	m	k _{kx}	k _{k_y}	J	Time	Rel-err	Suc-rat	Time	Rel-err	Suc-rat
800	1200	1600	5	5	400	0.08	1.58e-4	100%	0.22	1.69e-3	100%
800	1200	1600	40	20	400	0.11	2.88e-4	100%	0.25	1.59e-3	100%
800	1200	1600	80	40	400	0.19	6.19e-4	100%	0.55	1.65e-3	100%
1000	1500	2000	5	5	500	0.12	1.36e-4	100%	0.12	2.05e-3	100%
1000	1500	2000	100	50	500	0.31	5.74e-4	100%	0.53	1.73e-3	100%
2000	3000	4000	5	5	1000	0.49	1.39e-4	100%	0.43	2.25e-3	100%
2000	3000	4000	50	20	1000	0.54	1.79e-4	100%	1.41	1.55e-3	100%
2000	3000	4000	200	100	1000	1.14	5.77e-4	100%	1.19	1.63e-3	100%

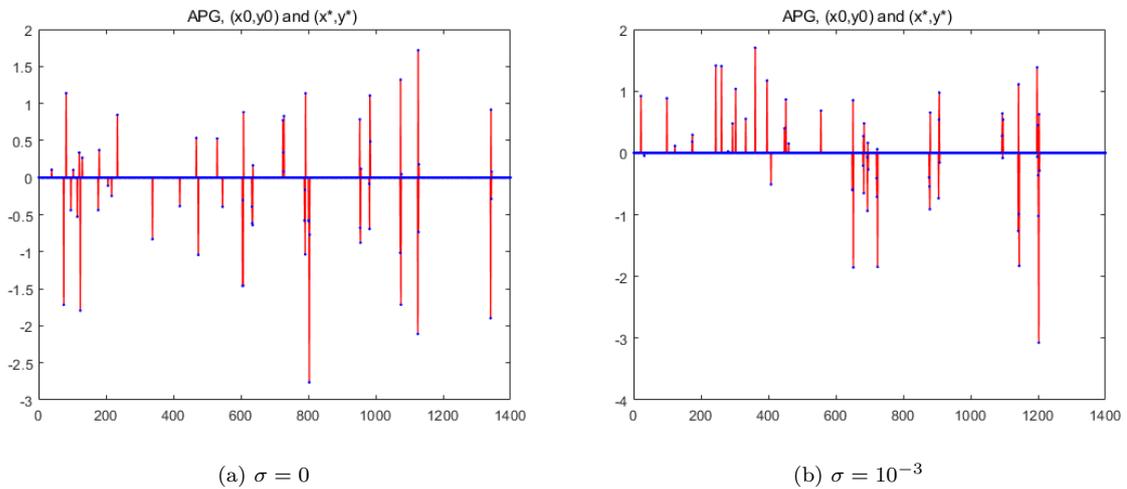


Figure 1. Visual numerical results

475 From Table 1 and Figure 1, we can see that the proposed APG algorithm can quickly
 476 obtain the true solution with high success rate.

477 Next, we compare our APG algorithm with several state-of-art algorithms: PGM-GSO
 478 algorithm [22] for solving $\ell_2\text{-}\ell_{p,q}$ model: $\min \|\mathbf{Ax} - \mathbf{b}\|_2^2 + \lambda \|\mathbf{x}\|_{p,q}^q$ ($p \geq 1, 0 \leq q \leq 1$), IRLS-th
 479 algorithm [18] for solving $\ell_{2,q}$ model: $\min \|\mathbf{Ax} - \mathbf{b}\|_2^2 + \lambda \|\mathbf{x}\|_{2,q}^q$ ($0 < q < 1$), GCD algorithm
 480 [8] for solving group MCP model, and SPGL1 algorithm [15] for solving group lasso model:
 481 $\min \|\mathbf{x}\|_{2,1}$ s.t. $\|\mathbf{Ax} - \mathbf{b}\|_2 \leq \delta$. One can refer to the references for their implementation
 482 details. In order for these algorithms to be used to solve problem (6.1), we group all partial
 483 sparse and partial group sparse data into groups as follows:

484 `gidxy=[]; Jx=floor(n/avgsz); gidxy=[gidxy;i*ones(avgsz,1)], i=1:(Jx+J);`

485 The above grouping way is applied to PGM-GSO, GCD, SPGL1 and IRLS-th. We run 100
 486 times for each instant and record the average CPU time and the average relative error, as
 487 shown in Table 2 and Table 3.

Table 2: Comparison of five algorithms for problem (6.1) with $\sigma = 0$

Problem					APG		PGM-GSO		SPG11		GCD		IRLS-th	
p	n	m	kkx	kky	Time	Rel-err	Time	Rel-err	Time	Rel-err	Time	Rel-err	Time	Rel-err
400	600	800	5	5	0.03	3.59e-4	0.04	2.05e-4	0.01	3.18e-4	1.17	3.41e-2	0.13	1.99e-3
400	600	800	25	15	0.05	8.15e-4	0.08	7.16e-4	0.06	3.50e-4	3.27	4.21e-2	0.42	3.87e-3
800	1200	1600	5	5	0.14	3.10e-4	0.16	1.25e-4	0.03	2.87e-4	3.38	3.91e-2	0.59	6.28e-4
800	1200	1600	20	10	0.17	3.48e-4	0.18	1.92e-4	0.05	3.91e-4	4.10	4.10e-2	0.75	1.17e-2
800	1200	1600	40	20	0.24	5.65e-4	0.24	4.87e-4	0.11	6.03e-4	6.79	2.97e-2	1.04	1.73e-2
1000	1500	2000	5	5	0.25	2.54e-4	0.27	1.35e-4	0.06	2.90e-4	4.06	1.53e-2	1.01	4.06e-4
1000	1500	2000	60	30	0.40	7.49e-4	0.42	5.91e-4	0.26	7.59e-4	13.07	2.89e-2	2.06	2.61e-2
2000	3000	4000	5	5	0.82	3.00e-4	1.32	1.09e-4	0.18	1.58e-4	12.96	1.85e-2	4.94	5.58e-4
2000	3000	4000	50	20	0.95	3.69e-4	1.46	1.95e-4	0.34	1.48e-4	26.84	2.20e-2	7.08	7.62e-4
2000	3000	4000	100	50	1.34	5.35e-4	1.81	4.67e-4	0.83	3.00e-4	43.37	1.68e-2	9.76	2.79e-3
4000	6000	8000	5	5	2.92	2.57e-4	8.63	9.79e-5	0.75	2.01e-4	50.47	8.99e-17	30.52	3.54e-4
4000	6000	8000	60	30	3.42	3.17e-4	8.89	1.40e-4	1.40	1.02e-4	93.57	6.65e-3	43.74	3.02e-3
4000	6000	8000	200	100	5.40	6.73e-4	10.57	4.12e-4	3.41	2.26e-4	144.43	1.15e-2	83.90	1.29e-2
6000	9000	12000	5	5	6.61	2.43e-4	93.44	8.95e-5	1.47	5.10e-5	166.46	1.12e-16	283.06	2.94e-4
6000	9000	12000	80	40	7.54	2.41e-4	104.10	1.77e-4	2.81	4.40e-5	186.57	6.03e-3	456.53	2.82e-3
6000	9000	12000	300	150	10.90	5.20e-4	99.21	4.35e-4	6.17	1.70e-4	312.56	7.78e-3	544.52	1.21e-2

Table 3: Comparison of five algorithms for problem (6.1) with $\sigma = 0.01$

Problem					APG		PGM-GSO		SPG11		GCD		IRLS-th	
p	n	m	kkx	kky	Time	Rel-err	Time	Rel-err	Time	Rel-err	Time	Rel-err	Time	Rel-err
400	600	800	5	5	0.06	3.69e-2	0.05	2.02e-2	0.03	4.17e-2	1.06	1.99e-2	0.26	6.53e-2
400	600	800	25	15	0.07	3.72e-2	0.08	2.40e-2	0.02	1.00e-1	1.83	3.03e-2	0.37	4.55e-2
800	1200	1600	5	5	1.57	3.27e-2	0.18	1.87e-2	0.04	2.95e-2	2.03	1.95e-2	0.85	6.25e-2
800	1200	1600	20	10	1.62	3.66e-2	0.20	2.03e-2	0.04	3.74e-2	3.50	2.43e-2	1.14	4.93e-2
800	1200	1600	40	20	0.28	3.43e-2	0.30	1.18e-2	0.05	3.21e-2	2.94	1.19e-2	2.08	6.89e-2
1000	1500	2000	5	5	0.26	4.29e-2	0.28	1.84e-2	0.05	3.05e-2	4.69	2.37e-2	2.03	8.72e-2
1000	1500	2000	60	30	2.99	3.92e-2	0.50	2.45e-2	0.12	6.09e-2	12.79	2.83e-2	2.75	4.18e-2
2000	3000	4000	5	5	0.89	2.75e-2	1.38	1.28e-2	0.18	3.00e-2	12.43	1.44e-2	9.79	9.93e-2
2000	3000	4000	50	20	10.06	3.14e-2	1.57	2.06e-2	0.23	3.78e-2	26.96	2.51e-2	11.14	5.13e-2
2000	3000	4000	100	50	10.02	3.25e-2	2.01	2.29e-2	0.41	5.29e-2	39.29	2.87e-2	14.35	4.41e-2
4000	6000	8000	5	5	2.94	3.45e-2	8.23	1.61e-2	0.62	2.83e-2	43.36	1.64e-2	49.48	1.11e-1
4000	6000	8000	60	30	36.14	2.68e-2	8.59	1.84e-2	0.85	3.55e-2	89.79	2.25e-2	64.84	5.54e-2
4000	6000	8000	200	100	35.95	2.17e-2	9.89	2.17e-2	1.49	5.08e-2	124.66	2.75e-2	86.07	5.25e-2
6000	9000	12000	5	5	7.68	3.96e-2	297.53	2.16e-2	3.31	3.93e-2	73.93	1.95e-2	3810.94	1.66e-1
6000	9000	12000	80	40	78.77	2.90e-2	555.80	1.97e-2	13.59	3.74e-2	178.94	2.71e-2	4951.63	6.65e-2
6000	9000	12000	300	150	80.17	2.89e-2	967.80	2.10e-2	10.82	4.77e-2	319.92	2.83e-2	12533.64	4.79e-2

488 From Tables 2 and 3, we can observe that in the absence of noise, the average relative
489 errors of APG are similar to SPG11 and PGM-GSO, but smaller than IRLS-th and GCD in
490 most cases; Meanwhile, the average CPU time of APG is less than PGM-GSO, GCD and

491 IRLS-th but more than SPG11. In the presence of noise, the average relative errors of APG
 492 are similar to the other four algorithms; Meanwhile, the average CPU time of APG is more
 493 than SPG11, but less than SPG11, GCD and IRLS-th; It is funny that, for the small scale
 494 instances, the average CPU time of APG is more than PGM-GSO but for the large scale
 495 instances, the average CPU time of APG is less than PGM-GSO. The results indicate that
 496 our APG algorithm is competitive with the four state-of-art algorithms in solving problem
 497 (6.1).

498 6.2 Multichannel image reconstruction

499 In this section, we consider recovering three-channel images from compressive and noisy
 500 measurement. In our experiments, the PSNR (peak signal to noise ratio) is defined by

$$\text{PSNR} = 10 \cdot \log \frac{V^2}{\text{MSE}},$$

501 in which V and $\text{MSE} = \frac{\|\mathbf{z} - \mathbf{z}_0\|^2}{n+m}$ (mean squared error) are the maximum absolute value and
 502 the mean squared error of the reconstruction respectively.

503 The example is taken from [24, 30, 26, 42]. The observed data \mathbf{c} is generated by $\mathbf{c} =$
 504 $\mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{y} + \sigma * \text{randn}(p, 1)$, where \mathbf{A} , \mathbf{B} are random Gaussian matrices, σ is a positive
 505 scalar, \mathbf{x} with sparse structure and \mathbf{y} with group sparse structure are the target coefficients.
 506 For this experiment: $n = 48 * 48 * 1$, $m = 48 * 48 * 2$, $p = m/2$, $J = m/4$, $kkx = 152$, $kk y = 172$.
 507 We still compare experimental results among APG, PGM-GSO, SPG11, GCD and IRLS-th.
 508 The PSNR and CPU time are presented in Table 4, while the original image and the recovered
 509 images for $\sigma = 0.1$ are presented in Figure 2.

Table 4: Numerical results for the three-channel image

σ	algorithm	APG	PGM-GSO	SPG11	GCD	IRLS-th
$\sigma = 0$	CPU time(s)	1.67	3.97	7.70	33.83	17.06
	PSNR	80.11	72.97	61.57	33.27	37.74
$\sigma = 1e - 3$	CPU time(s)	3.11	3.72	1.10	33.56	20.90
	PSNR	64.55	60.16	43.21	60.71	37.75
$\sigma = 1e - 2$	CPU time(s)	5.67	4.38	0.45	29.72	26.77
	PSNR	39.04	37.29	34.68	38.02	33.85
$\sigma = 1e - 1$	CPU time(s)	6.61	5.03	0.38	14.44	77.55
	PSNR	29.29	23.50	24.05	26.72	21.18

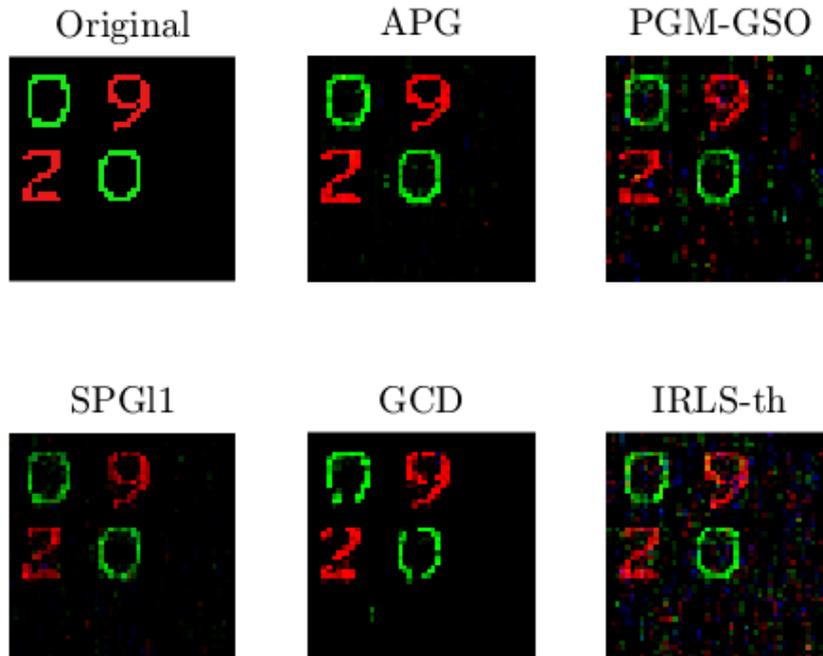


Figure 2. Original image and recovered images by five algorithms for $\sigma = 0.1$

510 From Table 4 and Figure 2, we can see that APG performs better than PGM-GSO,
 511 SPG11, GCD and IRLS-th in restoring the PSNR value of the image. Although APG is not
 512 superior to SPG11 and PGM-GSO in CPU time, it takes less time than GCD and IRLS-th.
 513 The results indicate that our model and APG algorithm are also competitive with the four
 514 state-of-art algorithms in multichannel image reconstruction.

515 7 Conclusion

516 In this paper, we initially studied the partial sparse and partial group sparse optimiza-
 517 tion problem. Firstly, we give the Capped- ℓ_1 relaxation and group Capped- ℓ_1 relaxation
 518 problem of the original problem. Secondly, we introduced d-stationary point and critical
 519 point for the relaxation problem, and prove that any d-stationary point is a critical point.
 520 Under some mild assumptions, we gave the lower bound properties of d-stationary points of
 521 the relaxation problem, based on which, we proved the equivalence of the original problem
 522 and the relaxation problem. This result provides a theoretical basis for solving the original
 523 problem via solving the relaxation problem. Then, we proposed an APG algorithm for the
 524 relaxation problem, and proved that the whole sequence generated by the APG algorithm
 525 converges to a critical point of the relaxation problem. Finally, the rich numerical experi-
 526 ments show that the partial sparse and partial group sparse model and the APG algorithm
 527 have good performance and some practical value.

528 **Acknowledgements** This work is supported by the National Natural Science Foundation of China (12261020),
529 the Guizhou Provincial Science and Technology Program (ZK[2021]009), the Foundation for Selected Excellent Project of
530 Guizhou Province for High-level Talents Back from Overseas ([2018]03), and the Research Foundation for Postgraduates
531 of Guizhou Province (YJSCXJH[2020]085).

532 References

- 533 1. Attouch, H., Bolte, J.: On the convergence of the proximal algorithm for nonsmooth functions involving analytic
534 features. *Mathematical Programming*, 2009, 116(1-2): 5-16.
- 535 2. Attouch, H., Bolte, J., Redont, P., Soubeyran, A.: Proximal alternating minimization and projection methods
536 for nonconvex problems: An approach based on the Kurdyka-Lojasiewicz inequality. *Mathematics of Operations
537 Research*, 2010, 35(2): 438-457.
- 538 3. Bian, W., Chen, X.: A smoothing proximal gradient algorithm for nonsmooth convex regression with cardinality
539 penalty. *SIAM Journal on Numerical Analysis*, 2020, 58(1): 858-883.
- 540 4. Bian, W., Chen, X.: Optimality and complexity for constrained optimization problems with nonconvex regulariza-
541 tion. *Mathematics of Operations Research*, 2017, 42(4): 1063-1084.
- 542 5. Blumensath, T.: Compressed sensing with nonlinear observations and related nonlinear optimization problems.
543 *IEEE Transactions on Information Theory*, 2013, 59(6): 3466-3474.
- 544 6. Bolte, J., Daniilidis, A., Lewis, A.: The Lojasiewicz inequality for nonsmooth subanalytic functions with applications
545 to subgradient dynamical systems. *SIAM Journal on Optimization*, 2007, 17(4): 1205-1223.
- 546 7. Bolte, J., Sabach, S., Teboulle, M.: Proximal alternating linearized minimization for nonconvex and nonsmooth
547 problems. *Mathematical Programming*, 2014, 146(1-2): 459-494.
- 548 8. Breheny, P., Huang, J.: Group descent algorithms for nonconvex penalized linear and logistic regression models
549 with grouped predictors. *Statistics and Computing*, 2015, 25(2): 173-187.
- 550 9. Chandran, M.: Analysis of Bayesian group-Lasso in regression models. University of Florida, 2011.
- 551 10. Chartrand, R.: Exact reconstruction of sparse signals via nonconvex minimization. *IEEE Signal Processing Letters*,
552 2007, 14(10): 707-710.
- 553 11. Chen, X., Pan, L.L., Xiu, N.: Solution sets of three sparse optimization problems for multivariate regression. *Journal
554 of Global Optimization*, 2022, <https://doi.org/10.1007/s10898-021-01124-w>.
- 555 12. Chen, X., Xu, F., Ye, Y.: Lower bound theory of nonzero entries in solutions of ℓ_2 - ℓ_p minimization. *SIAM Journal
556 on Scientific Computing*, 2010, 32(5): 2832-2852.
- 557 13. Clarke, F.H.: Optimization and Nonsmooth Analysis. *SIAM Journal on Control and Optimization*, 1990.
- 558 14. Elad, M., Figueiredo, M.A.T., Ma, Y.: On the role of sparse and redundant representations in image processing.
559 *Proceedings of the IEEE*, 2010, 98(6): 972-982.
- 560 15. Van den Berg, E., Friedlander, M.P.: Probing the pareto frontier for basis pursuit solutions. *SIAM Journal on
561 Scientific Computing*, 2009, 31(2): 890-912.
- 562 16. Fan, J., Li, R.: Statistical challenges with high dimensionality: feature selection in knowledge discovery. *Proceedings
563 of the International Congress of Mathematicians*, 2006, 3: 595-622.
- 564 17. Fan, J., Li, R.: Variable selection via nonconvex penalized likelihood and its oracle properties. *Journal of the
565 American Statistical Association*, 2001, 96(456): 1348-1360.
- 566 18. Feng, X., Yan, S., Wu, C.: The $\ell_{2,q}$ regularized group sparse optimization: lower bound theory, recovery bound and
567 algorithms. *Applied and Computational Harmonic Analysis*, 2020, 49(2): 381-414.
- 568 19. Gong, P., Zhang, C., Lu, Z., Huang, J., Ye, J.: A general iterative shrinkage and thresholding algorithm for non-
569 convex regularized optimization problems. *Proceedings of the 30th International Conference on International Con-
570 ference on Machine Learning (ICML'13)*, **28(2)**, 37-45 (2013)
- 571 20. Huang, J., Ma, S., Xie, H., Zhang, C.H.: A group bridge approach for variable selection. *Biometrika*, 2009, 96(2):
572 339-355.
- 573 21. Huang, J., Zhang, T.: The benefit of group sparsity. *The Annals of Statistics*, 2010, 38(4): 1978-2004.
- 574 22. Hu, Y., Li, C., Meng, K.: Group sparse optimization via $\ell_{p,q}$ regularization. *Journal of Machine Learning Research*,
575 2017, 18(30): 1-52.
- 576 23. Jiang, D.: Concave selection in generalized linear models. University of Iowa, 2012.
- 577 24. Jiao, Y., Jin, B., Lu, X.: Group sparse recovery via the $\ell_0(\ell_2)$ penalty: theory and algorithm. *IEEE Transactions
578 on Signal Processing*, 2017, 65(4): 998-1012.

- 579 25. Le Thi, H.A., Pham Dinh, T., Le, H.M., Vo, X.T.: DC approximation approaches for sparse optimization. *European*
580 *Journal of Operational Research*, 2015, 244(1): 26-46.
- 581 26. Li, W., Bian, W., Toh, K.C.: DC algorithms for a class of sparse group ℓ_0 regularized optimization problems. *SIAM*
582 *Journal on Optimization*, 2022, 32(3): 1614-1641.
- 583 27. Nikolova, M., Tan, P.: Alternating structure-adapted proximal gradient descent for nonconvex nonsmooth block-
584 regularized problems. *SIAM Journal on Optimization*, 2019, 29(3): 2053-2078.
- 585 28. Ong, C.S., An, L.T.H.: Learning sparse classifiers with difference of convex functions algorithms. *Optimization*
586 *Methods and Software*, 2013, 28(4): 830-854.
- 587 29. Pang, J.S., Razaviyayn, M., Alvarado, A.: Computing B-stationary points of nonsmooth DC programs. *Mathematics*
588 *of Operations Research*, 2017, 42(1): 95-118.
- 589 30. Pan, L., Chen, X.: Group sparse optimization for images recovery using capped folded concave functions. *SIAM*
590 *Journal on Imaging Sciences*, 2021, 14(1): 1-25.
- 591 31. Peng, D., Chen, X.: Computation of second-order directional stationary points for group sparse optimization.
592 *Optimization Methods and Software*, 2020, 35(2): 348-376.
- 593 32. Phan, D.N., Le Thi, H.A.: Group variable selection via $\ell_{p,0}$ regularization and application to optimal scoring. *Neural*
594 *Networks*, 2019, 118: 220-234.
- 595 33. Raman, S., Fuchs, T.J., Wild, P.J.: The Bayesian group-Lasso for analyzing contingency tables. *Proceedings of the*
596 *26th Annual International Conference on Machine Learning*, 2009, 881-888.
- 597 34. Rockafellar, R.T., Wets, R.J.B.: *Variational Analysis*. Springer Science & Business Media, 2009.
- 598 35. Simon, N., Friedman, J., Hastie T., Tibshirani R.: A sparse-group Lasso. *Journal of computational and graphical*
599 *statistics*, 2013, 22(2): 231-245.
- 600 36. Soubies, E., Blanc-Féraud, L., Aubert, G.: A continuous exact ℓ_0 penalty (Capped- ℓ_0) for least squares regularized
601 problem. *SIAM Journal on Imaging Sciences*, 2015, 8(3): 1574-1606.
- 602 37. Soubies, E., Blanc-Féraud, L., Aubert, G.: A unified view of exact continuous penalties for $\ell_2 - \ell_0$ minimization.
603 *SIAM Journal on Optimization*, 2017, 27(3): 2034-2060.
- 604 38. Wang, L., Chen, G., Li, H.: Group SCAD regression analysis for microarray time course gene expression data.
605 *Bioinformatics*, 2007, 23(12): 1486-1494.
- 606 39. Yuan, M., Lin, Y.: Model selection and estimation in regression with grouped variables. *Journal of the Royal*
607 *Statistical Society: Series B (Statistical Methodology)*, 2006, 68(1): 49-67.
- 608 40. Zhang, C.H.: Nearly unbiased variable selection under minimax concave penalty. *The Annals of Statistics*, 2010,
609 38(2): 894-942.
- 610 41. Zhang, T.: Analysis of multi-stage convex relaxation for sparse regularization. *Journal of Machine Learning Re-*
611 *search*, 2010, 11(35): 1081-1107.
- 612 42. Zhang, X., Peng, D.: Solving constrained nonsmooth group sparse optimization via group Capped- ℓ_1 relaxation
613 and group smoothing proximal gradient algorithm. *Computational Optimization and Applications*, 2022, 83(3):
614 801-844.
- 615 43. Zhang, X., Peng, D., Su, Y.: A singular value shrinkage thresholding algorithm for folded concave penalized low-rank
616 matrix optimization problems. *Journal of Global Optimization*, 2023, <https://doi.org/10.1007/s10898-023-01322-8>.
- 617 44. Zhang, Y., Zhang, N., Sun, D.: An efficient Hessian based algorithm for solving large-scale sparse group Lasso
618 problems. *Mathematical Programming*, 2020, 179(1): 223-263.
- 619 45. Zhao, P., Rocha, G., Yu, B.: The composite absolute penalties family for grouped and hierarchical variable selection.
620 *The Annals of Statistics*, 2009, 37(6A): 3468-3497.
- 621 46. Zhou, Y., Han, J., Yuan, X.: Inverse sparse group Lasso model for robust object tracking. *IEEE Transactions on*
622 *Multimedia*, 2017, 19(8): 1798-1810.