

Noname manuscript No.
(will be inserted by the editor)

1 Computation of second-order directional stationary points for group 2 sparse optimization

3 Dingtao Peng · Xiaojun Chen

4
5 **Abstract** We consider a nonconvex and nonsmooth group sparse optimization problem
6 where the penalty function is the sum of compositions of a folded concave function and
7 the ℓ_2 vector norm for each group variable. We show that under some mild conditions a
8 first-order directional stationary point is a strict local minimizer that fulfils the first-order
9 growth condition, and a second-order directional stationary point is a strong local minimizer
10 that fulfils the second-order growth condition. In order to compute second-order directional
11 stationary points, we construct a twice continuously differentiable smoothing problem and
12 show that any accumulation point of the sequence of second-order stationary points of the
13 smoothing problem is a second-order directional stationary point of the original problem. We
14 give numerical examples to illustrate how to compute a second-order directional stationary
15 point by the smoothing method.

16 **Keywords** Group sparse optimization; nonconvex and nonsmooth optimization; composite
17 folded concave penalty; directional stationary point; smoothing method

18 **MSC(2010)** 90C26 · 90C46

19 1 Introduction

20 Let $\mathbf{x} = (\mathbf{x}_1^\top, \dots, \mathbf{x}_K^\top)^\top \in \mathbb{R}^n$ with $\mathbf{x}_i = (x_{i(1)}, \dots, x_{i(d_i)})^\top \in \mathbb{R}^{d_i}$, $d_i \geq 1$, $\sum_{i=1}^K d_i = n$. We
21 consider the following optimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) := \mathcal{L}(\mathbf{x}) + \sum_{i=1}^K \varphi(\|\mathbf{x}_i\|), \quad (1.1)$$

22 where $\mathcal{L} : \mathbb{R}^n \rightarrow \mathbb{R}$ is a twice continuously differentiable function, and $\varphi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$
23 is a concave penalty function satisfying the following properties: (i) φ is locally Lipschitz
24 continuous and non-decreasing on $[0, \infty)$ with $\varphi(0) = 0$ and $\varphi(t) > 0$ for $t > 0$; (ii) $\varphi'(0+) >$
25 0 . Throughout this paper, $\|\cdot\|$ denotes the ℓ_2 vector norm.

The paper is dedicated to Professor Ya-Xiang Yuan on the occasion of his 60th birthday.

Dingtao Peng, School of Mathematics and Statistics, Guizhou University, Guiyang 550025, China.

E-mail: dingtaopeng@126.com

Xiaojun Chen , Department of Applied Mathematics, The Hong Kong Polytechnic University, Hong Kong, China.

E-mail: maxjchen@polyu.edu.hk

· This paper is partially supported by the Hong Kong Research Grant Council PolyU, PolyU153000/17P, NSFC (11861020), the Growth Project of Education Department of Guizhou Province for Young Talents in Science and Technology ([2018]121), and the Foundation for Selected Excellent Project of Guizhou Province for High-level Talents Back from Overseas ([2018]03)

In practice, many loss functions are twice continuously differentiable, for example, square loss function $\mathcal{L}(\mathbf{x}) = \frac{1}{2m} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2$, exponential loss function $\mathcal{L}(\mathbf{x}) = \frac{1}{m} \sum_{j=1}^m \exp(-b_j(\mathbf{a}_j^\top \mathbf{x}))$, and logistic loss function

$$\mathcal{L}(\mathbf{x}) = -\frac{1}{m} \sum_{j=1}^m \{b_j \log(1 + \exp(-\mathbf{a}_j^\top \mathbf{x})) + (1 - b_j) \log(1 + \exp(\mathbf{a}_j^\top \mathbf{x}))\},$$

26 where $\mathbf{b} \in \mathbb{R}^m$, $A = (\mathbf{a}_1, \dots, \mathbf{a}_m)^\top \in \mathbb{R}^{m \times n}$.

27 Problem (1.1) is called group sparse optimization due to the group structure in its vari-
 28 able. When $K = n$ and $d_1 = \dots = d_n = 1$, problem (1.1) reduces to the standard sparse
 29 optimization which is aimed to find a sparse solution to minimize the function $\mathcal{L}(\mathbf{x})$. Sparse
 30 optimization has attracted considerable attention in signal processing, machine learning and
 31 statistics in recent years. To yield a sparse solution, a penalty term is often used. Tibshirani
 32 [28] suggested using the ℓ_1 penalty to obtain a sparse vector of regression coefficients in linear
 33 regression problem, which results in a convex optimization problem, called Lasso, and can
 34 be solved by many efficient algorithms. However, Fan and Li [12,13] pointed out that the
 35 solution of the ℓ_1 penalized optimization does not possess some good statistical properties
 36 such as unbiasedness and oracle property. Fan and Li [12,13] then proposed a folded concave
 37 penalty and showed that there exists a local solution with the desired statistical properties
 38 for the resulting non-convex optimization. Till now, many specific folded concave penalty
 39 functions are widely used in signal reconstruction, image restoration, and variable selection,
 40 for example, logarithm penalty [12], fraction penalty [25], hard thresholding penalty (HT-
 41 P) [6,20], capped ℓ_1 penalty (CapL1) [36], minimax concave penalty (MCP) [35], smoothly
 42 clipped absolute deviation (SCAD) [12].

43 Although there exist some local minimizers with good statistical properties for a folded
 44 concave penalized optimization, how to find such local minimizers has not been addressed
 45 satisfactorily. Fan, Xue and Zou [14] proposed a local linear approximation algorithm to
 46 obtain an oracle solution with an initial point being sufficiently close to the true solution. In
 47 [23], the authors developed a concept of subspace second-order optimality which is related
 48 to subspace optimality in [3,4,9,10], and showed that under some conditions the station-
 49 ary point of subspace second-order optimality can be an oracle solution with high prob-
 50 ability. In 1985, Yuan [33] studied convergence of trust region algorithms to a first-order
 51 d(irectional)-stationary point of nonsmooth optimization. Recently, [1,26] adopted a first-
 52 order d(irectional)-stationary point for optimality, and showed that a first-order d-stationary
 53 point must be one of other stationary points using the first-order information of the objec-
 54 tive function. Moreover, [7,27] proposed the concept of second-order directional derivatives
 55 and the concept of second-order d(irectional)-stationary points, and showed that under some
 56 mild conditions second-order d-stationary points can fulfil the second-order growth condition.
 57 However, how to compute second-order directional derivatives and second-order d-stationary
 58 points is unknown for problem (1.1).

59 Group sparse problem was studied by many authors, e.g., see [11,15,16,17,18,19,24,
 60 30,32,34,37]. It has wide applications in statistics, machine learning, and computational
 61 biology such as joint covariate selection [16,17,34,37], multi-task learning [19,32], and gene
 62 finding [15,24]. Most of the literatures use group ℓ_1 penalty which yields group Lasso model.
 63 Huang and Zhang [17] showed that group Lasso is superior to standard Lasso for strongly
 64 group-sparse signals. In consideration of the good performance of folded concave penalties
 65 comparing to ℓ_1 penalty for standard sparse optimization, some authors used group folded
 66 concave penalties such as group SCAD [5,22,29], group MCP [5,22,29], $\ell_q(\ell_p)$ ($0 \leq q \leq 1 \leq$

p) [15] and $\ell_0(\ell_2)$ [19] for group sparse problems. However, these works only used first-order information of objective functions which is weaker than second-order information.

In this paper, we will provide a deep analysis of the second-order directional stationarity for folded concave penalized group sparse optimization. Our main contributions are presented as follows.

In Section 2, by virtue of an explicit formula for computing the directional derivative of the objective function, we show that under some mild conditions a first-order d-stationary point of problem (1.1) is a strict local minimizer that fulfils the first-order growth condition.

In Section 3, we provide an explicit formula for computing the second-order directional derivative, and show that under some mild conditions a second-order d-stationary point of problem (1.1) is a strong local minimizer that fulfils the second-order growth condition. Moreover, we establish lower bounds of the ℓ_2 vector norm of nonzero groups of second-order d-stationary points of problem (1.1). These lower bounds are important for theoretical analysis and numerical algorithms.

In Section 4, we construct a twice continuously differentiable smoothing approximation for the nonsmooth objective function in problem (1.1), and show that any accumulation point of the sequence of second-order stationary points of the smoothing problem is a second-order d-stationary point of the original problem. This result provides a theoretic basis for computing second-order d-stationary points of problem (1.1) using the gradient and Hessian of the smoothing function.

Notations. For any $\hat{\mathbf{x}} \in \mathbb{R}^n$ and the groups $\hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_K$, denote

$$\begin{aligned} I(\hat{\mathbf{x}}) &:= \{i \in \{1, \dots, K\} : \|\hat{\mathbf{x}}_i\| \neq \mathbf{0}\}, & J_i(\hat{\mathbf{x}}) &:= \{j \in \{1, \dots, d_i\} : \hat{x}_{i(j)} \neq 0\} \text{ for } i \in I(\hat{\mathbf{x}}), \\ & i \notin I(\hat{\mathbf{x}}) \text{ if } i \in \{1, \dots, K\} \setminus I(\hat{\mathbf{x}}), & j \notin J_i(\hat{\mathbf{x}}) & \text{ if } i \in I(\hat{\mathbf{x}}) \text{ and } j \in \{1, \dots, d_i\} \setminus J_i(\hat{\mathbf{x}}), \\ [\nabla \mathcal{L}(\hat{\mathbf{x}})]_i &:= ([\nabla \mathcal{L}(\hat{\mathbf{x}})]_{i(1)}, \dots, [\nabla \mathcal{L}(\hat{\mathbf{x}})]_{i(d_i)})^\top, & \nabla \mathcal{L}(\hat{\mathbf{x}}) &:= ([\nabla \mathcal{L}(\hat{\mathbf{x}})]_1^\top, \dots, [\nabla \mathcal{L}(\hat{\mathbf{x}})]_K^\top)^\top, \end{aligned}$$

where $\hat{x}_{i(j)} \in \mathbb{R}$ denotes the j th entry in $\hat{\mathbf{x}}_i$ and $[\nabla \mathcal{L}(\hat{\mathbf{x}})]_{i(j)}$ denotes the j th entry in $[\nabla \mathcal{L}(\hat{\mathbf{x}})]_i$.

2 First-order d-stationary points

This section provides the local optimality and some properties of first-order d-stationary points of problem (1.1).

2.1 Local optimality of first-order d-stationary points

Let us introduce the concept of first-order d-stationary points [1, 7, 26, 27].

Definition 2.1 $\hat{\mathbf{x}} \in \mathbb{R}^n$ is called a first-order d-stationary point of problem (1.1) if the directional derivative satisfies

$$f'(\hat{\mathbf{x}}; \mathbf{x} - \hat{\mathbf{x}}) := \lim_{t \downarrow 0} \frac{f(\hat{\mathbf{x}} + t(\mathbf{x} - \hat{\mathbf{x}})) - f(\hat{\mathbf{x}})}{t} \geq 0, \quad \forall \mathbf{x} \in \mathbb{R}^n. \quad (2.1)$$

According to [1, 26], first-order d-stationary points are sharper than lifted stationary points, critical points, and C-stationary points for the local optimality. It is known that first-order d-stationary points have the following locally optimal properties.

Theorem 2.2 Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be locally Lipschitz continuous and directionally differentiable at $\hat{\mathbf{x}} \in \mathbb{R}^n$. The following two statements hold:

(i) If $\hat{\mathbf{x}}$ is a local minimizer of f , then $\hat{\mathbf{x}}$ is a first-order d -stationary point of f .

(ii) $\hat{\mathbf{x}}$ is a strict local minimizer that fulfils the first-order growth condition, i.e., there exists a neighborhood \mathcal{W} of $\hat{\mathbf{x}}$ and a positive number δ such that

$$f(\mathbf{x}) \geq f(\hat{\mathbf{x}}) + \delta \|\mathbf{x} - \hat{\mathbf{x}}\|, \quad \forall \mathbf{x} \in \mathcal{W}, \quad (2.2)$$

if and only if $\hat{\mathbf{x}}$ satisfies that

$$f'(\hat{\mathbf{x}}; \mathbf{x} - \hat{\mathbf{x}}) > 0, \quad \forall \mathbf{x} \in \mathbb{R}^n \setminus \{\hat{\mathbf{x}}\}. \quad (2.3)$$

If f is differentiable at \mathbf{x} , then $f'(\mathbf{x}; \mathbf{z}) = \langle \nabla f(\mathbf{x}), \mathbf{z} \rangle$. Inequality (2.3) does not hold at any differentiable point of f , but it may hold at some non-differentiable points of f . Many local minimizers of problem (1.1) are non-differentiable points of f , which makes conclusion (ii) of Theorem 2.2 very interesting. For example, let $f(t) = t^2 + \log(1 + |t|)$, then $f'(0; s) = |s| > 0$ ($s \neq 0$), and $f(t) \geq |t|$ for any $t \in \mathbb{R}$.

To have a clear presentation, we denote the ℓ_2 vector norm as a function

$$m(\mathbf{u}) := \|\mathbf{u}\| = \left(\sum_{j=1}^{d_i} u_j^2 \right)^{\frac{1}{2}}, \quad \forall \mathbf{u} \in \mathbb{R}^{d_i}, \quad i \in \{1, \dots, K\}. \quad (2.4)$$

Although the dimensions of the vectors may be different, we believe that it will not cause any confusion according to the context.

Since $m(\mathbf{u})$ is differentiable at all points except $\mathbf{u} = \mathbf{0}$, we have that for any $\mathbf{u}, \mathbf{w} \in \mathbb{R}^{d_i}$,

$$m'(\mathbf{u}; \mathbf{w}) = \lim_{t \downarrow 0} \frac{\|\mathbf{u} + t\mathbf{w}\| - \|\mathbf{u}\|}{t} = \begin{cases} \|\mathbf{w}\|, & \text{if } \mathbf{u} = \mathbf{0}, \\ \frac{\langle \mathbf{u}, \mathbf{w} \rangle}{\|\mathbf{u}\|}, & \text{if } \mathbf{u} \neq \mathbf{0}. \end{cases} \quad (2.5)$$

2.2 First-order d -stationary points of problem (1.1)

In this subsection, we use an explicit formula of directional derivative to provide sufficient and necessary conditions for first-order d -stationary points of problem (1.1).

Our analysis is based on a difference-of-convex (DC) form of the penalty function so that the directional derivative of the objective function in (1.1) can be explicitly expressed.

Assumption (A1): The penalty function $\varphi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ is a DC function given by

$$\varphi(t) \triangleq g(t) - h(t), \quad \text{with } h(t) \triangleq \max_{1 \leq \nu \leq \bar{\nu}} \{h_\nu(t)\} \text{ for some integer } \bar{\nu} \geq 1, \quad (2.6)$$

where g and h_ν ($1 \leq \nu \leq \bar{\nu}$) are convex and differentiable in $t \in (0, \infty)$ with $g'(0) := g'(0+)$ and $h'_\nu(0) := h'_\nu(0+)$ for $1 \leq \nu \leq \bar{\nu}$.

Consequently, our group sparse optimization model (1.1) is rewritten as

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) := \mathcal{L}(\mathbf{x}) + \sum_{i=1}^K [g(\|\mathbf{x}_i\|) - h(\|\mathbf{x}_i\|)]. \quad (2.7)$$

From the literatures (e.g., [1, 21]), we know that several folded concave penalty functions can be formulated as DC functions satisfying Assumption (A1), such as logarithm penalty, fraction penalty, CapL1, HTP, MCP and SCAD. In particular, as given in [1] we have the following expressions:

CapL1: $\varphi^{\text{CapL1}}(t) = g^{\text{CapL1}}(t) - h^{\text{CapL1}}(t)$ with

$$g^{\text{CapL1}}(t) = \frac{\lambda t}{\alpha}, \quad h^{\text{CapL1}}(t) = \max \left\{ 0, \frac{\lambda t}{\alpha} - \lambda \right\}, \quad (\alpha > 0, \lambda > 0);$$

MCP: $\varphi^{\text{MCP}}(t) = g^{\text{MCP}}(t) - h^{\text{MCP}}(t)$ with

$$g^{\text{MCP}}(t) = \lambda t, \quad h^{\text{MCP}}(t) = \begin{cases} \frac{t^2}{2\alpha}, & \text{if } 0 \leq t \leq \alpha\lambda, \\ \lambda t - \frac{\alpha\lambda^2}{2}, & \text{if } t > \alpha\lambda, \end{cases} \quad (\alpha > 1, \lambda > 0);$$

SCAD: $\varphi^{\text{SCAD}}(t) = g^{\text{SCAD}}(t) - h^{\text{SCAD}}(t)$ with

$$g^{\text{SCAD}}(t) = \lambda t, \quad h^{\text{SCAD}}(t) = \begin{cases} 0, & \text{if } 0 \leq t \leq \lambda, \\ \frac{(t-\lambda)^2}{2(\alpha-1)}, & \text{if } \lambda < t \leq \alpha\lambda, \\ \lambda t - \frac{(\alpha+1)\lambda^2}{2}, & \text{if } t > \alpha\lambda, \end{cases} \quad (\alpha > 2, \lambda > 0).$$

127 **Theorem 2.3** Under Assumption (A1), the directional derivative of the objective function
128 f in (1.1) has the following form

$$\begin{aligned} f'(\widehat{\mathbf{x}}; \mathbf{x} - \widehat{\mathbf{x}}) &= \langle \nabla \mathcal{L}(\widehat{\mathbf{x}}), \mathbf{x} - \widehat{\mathbf{x}} \rangle + \sum_{i=1}^K g'(\|\widehat{\mathbf{x}}_i\|) m'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) \\ &\quad - \sum_{i=1}^K \max_{\nu_i \in \mathcal{A}_i(\widehat{\mathbf{x}}_i)} h'_{\nu_i}(\|\widehat{\mathbf{x}}_i\|) m'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) \end{aligned} \quad (2.8)$$

129 with $\mathcal{A}_i(\widehat{\mathbf{x}}_i) = \{\nu_i \in \{1, \dots, \bar{\nu}\} : h_{\nu_i}(\|\widehat{\mathbf{x}}_i\|) = h(\|\widehat{\mathbf{x}}_i\|)\}$ and

$$m'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) = \begin{cases} \|\mathbf{x}_i\|, & \text{if } i \notin I(\widehat{\mathbf{x}}), \\ \frac{\langle \widehat{\mathbf{x}}_i, \mathbf{x}_i - \widehat{\mathbf{x}}_i \rangle}{\|\widehat{\mathbf{x}}_i\|}, & \text{if } i \in I(\widehat{\mathbf{x}}_i). \end{cases} \quad (2.9)$$

130 *Proof* Under Assumption (A1), problem (1.1) can be written as (2.7). Since $h : \mathbb{R}_+ \rightarrow \mathbb{R}_+$
131 and $m : \mathbb{R}^{d_i} \rightarrow \mathbb{R}_+$ are both convex, $h \circ m : \mathbb{R}^{d_i} \rightarrow \mathbb{R}_+$ is directionally differentiable.
132 According to the chain rule for directional derivatives and the differentiability of each h_{ν} ,
133 for $i = 1, \dots, K$, we have

$$(h \circ m)'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) = h'(\|\widehat{\mathbf{x}}_i\|; m'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i)) = \max_{\nu_i \in \mathcal{A}_i(\widehat{\mathbf{x}}_i)} h'_{\nu_i}(\|\widehat{\mathbf{x}}_i\|) m'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i).$$

134 Since \mathcal{L} and g are differentiable, we obtain the directional derivative at $\widehat{\mathbf{x}}$ for $\mathbf{x} - \widehat{\mathbf{x}}$ in (2.8).
135 \square

136 The following lemma shows that at any first-order d-stationary point of (1.1), the entries
137 of the gradient of the loss function \mathcal{L} for $i \in I(\widehat{\mathbf{x}})$ can be presented by the derivatives of g
138 and h_{ν_i} .

139 **Lemma 2.4** Suppose Assumption (A1) holds. Let $\widehat{\mathbf{x}} \in \mathbb{R}^n$ be a first-order d-stationary point
140 of problem (1.1). Then for $i \in I(\widehat{\mathbf{x}})$, we have

$$[\nabla \mathcal{L}(\widehat{\mathbf{x}})]_{i(j)} = 0, \quad \forall j \notin J_i(\widehat{\mathbf{x}}), \quad (2.10)$$

141 and

$$|[\nabla \mathcal{L}(\widehat{\mathbf{x}})]_{i(j)}| = \frac{|g'(\|\widehat{\mathbf{x}}_i\|) - h'_{\nu_i}(\|\widehat{\mathbf{x}}_i\|)| \cdot |\widehat{x}_{i(j)}|}{\|\widehat{\mathbf{x}}_i\|}, \quad \forall j \in J_i(\widehat{\mathbf{x}}), \quad \forall \nu_i \in \mathcal{A}_i(\widehat{\mathbf{x}}_i). \quad (2.11)$$

142 *Proof* From Theorem 2.3, we have

$$\begin{aligned}
& \langle \nabla \mathcal{L}(\widehat{\mathbf{x}}), \mathbf{x} - \widehat{\mathbf{x}} \rangle + \sum_{i=1}^K [g'(\|\widehat{\mathbf{x}}_i\|) - h'_{\nu_i}(\|\widehat{\mathbf{x}}_i\|)] m'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) \\
& \geq \langle \nabla \mathcal{L}(\widehat{\mathbf{x}}), \mathbf{x} - \widehat{\mathbf{x}} \rangle + \sum_{i=1}^K g'(\|\widehat{\mathbf{x}}_i\|) m'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) - \sum_{i=1}^K \max_{\nu_i \in \mathcal{A}_i(\widehat{\mathbf{x}}_i)} h'_{\nu_i}(\|\widehat{\mathbf{x}}_i\|) m'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) \\
& \geq 0, \quad \forall \mathbf{x} \in \mathbb{R}^n,
\end{aligned} \tag{2.12}$$

143 where $\nu_i \in \mathcal{A}_i(\widehat{\mathbf{x}}_i)$, $i = 1, \dots, K$, and $m'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i)$ is given by (2.9). It is obvious that
144 inequality (2.12) also holds for any $\mathbf{x} \in \mathcal{X}(\widehat{\mathbf{x}}) := \{\mathbf{x} \in \mathbb{R}^n : \mathbf{x}_i = \mathbf{0} \text{ whenever } i \notin I(\widehat{\mathbf{x}})\}$. This
145 combining with formula (2.9) yields that

$$\sum_{i \in I(\widehat{\mathbf{x}})} \left\langle [\nabla \mathcal{L}(\widehat{\mathbf{x}})]_i + \frac{[g'(\|\widehat{\mathbf{x}}_i\|) - h'_{\nu_i}(\|\widehat{\mathbf{x}}_i\|)]}{\|\widehat{\mathbf{x}}_i\|} \widehat{\mathbf{x}}_i, \mathbf{x}_i - \widehat{\mathbf{x}}_i \right\rangle \geq 0, \quad \forall \mathbf{x} \in \mathcal{X}(\widehat{\mathbf{x}}).$$

146 According to the arbitrariness of $\mathbf{x} \in \mathcal{X}(\widehat{\mathbf{x}})$, we obtain

$$[\nabla \mathcal{L}(\widehat{\mathbf{x}})]_i + \frac{g'(\|\widehat{\mathbf{x}}_i\|) - h'_{\nu_i}(\|\widehat{\mathbf{x}}_i\|)}{\|\widehat{\mathbf{x}}_i\|} \widehat{\mathbf{x}}_i = \mathbf{0}, \quad \forall i \in I(\widehat{\mathbf{x}}). \tag{2.13}$$

147 Therefore, we have

$$[\nabla \mathcal{L}(\widehat{\mathbf{x}})]_{i(j)} = 0, \quad \forall i \in I(\widehat{\mathbf{x}}), j \notin J_i(\widehat{\mathbf{x}}),$$

148 and

$$\frac{|g'(\|\widehat{\mathbf{x}}_i\|) - h'_{\nu_i}(\|\widehat{\mathbf{x}}_i\|)| \cdot |\widehat{x}_{i(j)}|}{\|\widehat{\mathbf{x}}_i\|} = |[\nabla \mathcal{L}(\widehat{\mathbf{x}})]_{i(j)}|, \quad \forall i \in I(\widehat{\mathbf{x}}), j \in J_i(\widehat{\mathbf{x}}).$$

149 The conclusion is obtained. \square

150 By applying Lemma 2.4 to CapL1, MCP and SCAD, we can get the following lower
151 bounds of the ℓ_2 vector norm of nonzero groups of first-order d-stationary points, whose
152 proof is omitted.

153 **Corollary 2.5** *Suppose there exists a nondecreasing function $C : \mathbb{R} \rightarrow \mathbb{R}_+$ such that $\|\nabla \mathcal{L}(\mathbf{x})\|$
154 $\leq C(\mathcal{L}(\mathbf{x}))$ for any $\mathbf{x} \in \mathbb{R}^n$. Let $\widehat{\mathbf{x}} \in \mathbb{R}^n$ be a first-order d-stationary point of problem (1.1),
155 and $\mathbf{x}^0 \in \mathbb{R}^n$ be a point such that $\mathcal{L}(\widehat{\mathbf{x}}) \leq \mathcal{L}(\mathbf{x}^0)$, then the following statements hold:*

- 156 (i) *For CapL1, if $\frac{\lambda}{\alpha} > C(\mathcal{L}(\mathbf{x}^0))$, then either $\|\widehat{\mathbf{x}}_i\| = 0$ or $\|\widehat{\mathbf{x}}_i\| \geq \alpha$, $i = 1, \dots, K$.*
157 (ii) *For MCP, if $\lambda > C(\mathcal{L}(\mathbf{x}^0))$, then either $\|\widehat{\mathbf{x}}_i\| = 0$ or $\|\widehat{\mathbf{x}}_i\| \geq \alpha\lambda - \alpha \cdot C(\mathcal{L}(\mathbf{x}^0)) > 0$,
158 $i = 1, \dots, K$.*
159 (iii) *For SCAD, if $\lambda > C(\mathcal{L}(\mathbf{x}^0))$, then either $\|\widehat{\mathbf{x}}_i\| = 0$ or $\|\widehat{\mathbf{x}}_i\| \geq \alpha\lambda - (\alpha - 1) \cdot C(\mathcal{L}(\mathbf{x}^0)) >$
160 λ , $i = 1, \dots, K$.*

161 **Remark 2.6** *The existence of the nondecreasing function $C : \mathbb{R} \rightarrow \mathbb{R}_+$ means that the norm
162 of the gradient $\nabla \mathcal{L}(\mathbf{x})$ can be bounded by the function value $\mathcal{L}(\mathbf{x})$ via $C(\cdot)$. This condition
163 can be easily satisfied, for example, for the square loss function $\mathcal{L}(\mathbf{x}) = \frac{1}{2m} \|A\mathbf{x} - b\|^2$, $C(t) =$
164 $\|A\|_2 \sqrt{\frac{2}{m} t}$ meets the requirements since*

$$\|\nabla \mathcal{L}(\mathbf{x})\| = \frac{1}{m} \|A^\top (A\mathbf{x} - b)\| \leq \frac{\|A\|_2}{m} \|A\mathbf{x} - b\| = \|A\|_2 \sqrt{\frac{2}{m} \mathcal{L}(\mathbf{x})}.$$

165 When φ is the difference of two differentiable convex functions in $(0, \infty)$, such as φ^{MCP}
 166 and φ^{SCAD} , we have the following corollary, which will be used in Theorem 4.5 to derive the
 167 consistency of the second-order stationary point.

168 **Corollary 2.7** Suppose Assumption (A1) holds with $\bar{\nu} = 1$, that is, $\varphi = g - h$ where g, h
 169 are both convex and differentiable in $(0, \infty)$. Let $\hat{\mathbf{x}} \in \mathbb{R}^n$ be a first-order d-stationary point
 170 of problem (1.1), then the following statements hold:

- 171 (i) $f'(\hat{\mathbf{x}}; \mathbf{x} - \hat{\mathbf{x}}) = \sum_{i \notin I(\hat{\mathbf{x}})} \left[\langle [\nabla \mathcal{L}(\hat{\mathbf{x}})]_i, \mathbf{x}_i \rangle + \varphi'(0) \|\mathbf{x}_i\| \right]$ for any $\mathbf{x} \in \mathbb{R}^n$.
 172 (ii) $\|[\nabla \mathcal{L}(\hat{\mathbf{x}})]_i\| \leq \varphi'(0)$ whenever $i \notin I(\hat{\mathbf{x}})$.
 173 (iii) $f'(\hat{\mathbf{x}}; \mathbf{x} - \hat{\mathbf{x}}) = 0$ implies $\mathbf{x}_i = \mathbf{0}$ whenever $i \notin I(\hat{\mathbf{x}})$ and $\|[\nabla \mathcal{L}(\hat{\mathbf{x}})]_i\| < \varphi'(0)$.

174 *Proof* (i) From Theorem 2.3, $f'(\hat{\mathbf{x}}; \mathbf{x} - \hat{\mathbf{x}})$ has the following form

$$\begin{aligned} f'(\hat{\mathbf{x}}; \mathbf{x} - \hat{\mathbf{x}}) &= \sum_{i \in I(\hat{\mathbf{x}})} \left\langle [\nabla \mathcal{L}(\hat{\mathbf{x}})]_i + \frac{[g'(\|\hat{\mathbf{x}}_i\|) - h'(\|\hat{\mathbf{x}}_i\|)]}{\|\hat{\mathbf{x}}_i\|} \hat{\mathbf{x}}_i, \mathbf{x}_i - \hat{\mathbf{x}}_i \right\rangle \\ &\quad + \sum_{i \notin I(\hat{\mathbf{x}})} \left[\langle [\nabla \mathcal{L}(\hat{\mathbf{x}})]_i, \mathbf{x}_i \rangle + (g'(0) - h'(0)) \|\mathbf{x}_i\| \right]. \end{aligned} \quad (2.14)$$

175 Since $\hat{\mathbf{x}}$ is a first-order d-stationary point of problem (1.1), equation (2.13) holds with $h_{\nu_i} = h$.
 176 Hence (2.14) can be simplified as

$$f'(\hat{\mathbf{x}}; \mathbf{x} - \hat{\mathbf{x}}) = \sum_{i \notin I(\hat{\mathbf{x}})} \left[\langle [\nabla \mathcal{L}(\hat{\mathbf{x}})]_i, \mathbf{x}_i \rangle + \varphi'(0) \|\mathbf{x}_i\| \right],$$

177 where $\varphi'(0) = g'(0) - h'(0) > 0$.

178 (ii) Since $\hat{\mathbf{x}}$ is a first-order d-stationary point of problem (1.1), $f'(\hat{\mathbf{x}}; \mathbf{x} - \hat{\mathbf{x}}) \geq 0$ for all
 179 $\mathbf{x} \in \mathbb{R}^n$, that is,

$$f'(\hat{\mathbf{x}}; \mathbf{x} - \hat{\mathbf{x}}) = \sum_{i \notin I(\hat{\mathbf{x}})} \left[\langle [\nabla \mathcal{L}(\hat{\mathbf{x}})]_i, \mathbf{x}_i \rangle + \varphi'(0) \|\mathbf{x}_i\| \right] \geq 0, \quad \forall \mathbf{x} \in \mathbb{R}^n. \quad (2.15)$$

180 For each fixed $i \notin I(\hat{\mathbf{x}})$, if we take $\check{\mathbf{x}}_i = -[\nabla \mathcal{L}(\hat{\mathbf{x}})]_i$ and the other entries of $\check{\mathbf{x}}$ are all zeros,
 181 then we get

$$f'(\hat{\mathbf{x}}; \check{\mathbf{x}} - \hat{\mathbf{x}}) = \|[\nabla \mathcal{L}(\hat{\mathbf{x}})]_i\| \cdot \left[\varphi'(0) - \|[\nabla \mathcal{L}(\hat{\mathbf{x}})]_i\| \right] \geq 0. \quad (2.16)$$

182 If $\|[\nabla \mathcal{L}(\hat{\mathbf{x}})]_i\| = 0$, then $\|[\nabla \mathcal{L}(\hat{\mathbf{x}})]_i\| = 0 < \varphi'(0)$. If $\|[\nabla \mathcal{L}(\hat{\mathbf{x}})]_i\| > 0$, then from (2.16), we
 183 obtain $\varphi'(0) \geq \|[\nabla \mathcal{L}(\hat{\mathbf{x}})]_i\|$.

184 (iii) It follows from (i), (ii) and Cauchy-Schwartz inequality that

$$\begin{aligned} f'(\hat{\mathbf{x}}; \mathbf{x} - \hat{\mathbf{x}}) &= \sum_{i \notin I(\hat{\mathbf{x}})} \left[\langle [\nabla \mathcal{L}(\hat{\mathbf{x}})]_i, \mathbf{x}_i \rangle + \varphi'(0) \|\mathbf{x}_i\| \right] \\ &\geq \sum_{i \notin I(\hat{\mathbf{x}})} \left[\varphi'(0) - \|[\nabla \mathcal{L}(\hat{\mathbf{x}})]_i\| \right] \|\mathbf{x}_i\| \geq 0. \end{aligned}$$

185 Hence, if $f'(\hat{\mathbf{x}}; \mathbf{x} - \hat{\mathbf{x}}) = 0$, it must hold that $\|\mathbf{x}_i\| = 0$ whenever $i \notin I(\hat{\mathbf{x}})$ and $\|[\nabla \mathcal{L}(\hat{\mathbf{x}})]_i\| <$
 186 $\varphi'(0)$. \square

3 Second-order d-stationary points

In this section, we provide second-order optimality conditions for problem (1.1) using second-order directional derivatives.

3.1 Local optimality of second-order d-stationary points

Second-order directional derivatives for nonsmooth functions have been studied by many authors (e.g., see [2, 7, 27, 31]) with different definitions for one direction or two directions. In this paper, we use the definition of the second-order directional derivative for one direction in [7, 27] to define the second-order d-stationary point of problem (1.1). We show that second-order d-stationary points of problem (1.1) are local minimizers fulfilling the second-order growth condition under some mild conditions.

Definition 3.1 [7, 27] *Let $\theta : \mathbb{R}^n \rightarrow \mathbb{R}$ be a locally Lipschitz continuous and directionally differentiable function, and $\hat{\mathbf{x}}, \mathbf{z} \in \mathbb{R}^n$. If the limit*

$$\lim_{\mathbf{y} \rightarrow \mathbf{z}, t \downarrow 0} \frac{\theta(\hat{\mathbf{x}} + t\mathbf{y}) - \theta(\hat{\mathbf{x}}) - t\theta'(\hat{\mathbf{x}}; \mathbf{y})}{\frac{1}{2}t^2} \quad (3.1)$$

exists, it is called the second-order directional derivative of θ at $\hat{\mathbf{x}}$ for \mathbf{z} , denoted by $\theta^{(2)}(\hat{\mathbf{x}}; \mathbf{z})$. If for every $\mathbf{z} \in \mathbb{R}^n$, $\theta^{(2)}(\hat{\mathbf{x}}; \mathbf{z})$ exists, θ is called twice directionally differentiable at $\hat{\mathbf{x}}$.

Indeed, to say that limit (3.1) exists and equals $\theta^{(2)}(\hat{\mathbf{x}}; \mathbf{z})$ is to say that whenever \mathbf{x}^ν converges to $\hat{\mathbf{x}}$ from the direction of \mathbf{z} , in the sense that $[\mathbf{x}^\nu - \hat{\mathbf{x}}]/t^\nu \rightarrow \mathbf{z}$ for some choice of $t^\nu \downarrow 0$, one has

$$\frac{\theta(\mathbf{x}^\nu) - \theta(\hat{\mathbf{x}}) - \theta'(\hat{\mathbf{x}}; \mathbf{x}^\nu - \hat{\mathbf{x}})}{\frac{1}{2}(t^\nu)^2} \rightarrow \theta^{(2)}(\hat{\mathbf{x}}; \mathbf{z}).$$

Clearly, if limit (3.1) exists, then

$$\theta^{(2)}(\hat{\mathbf{x}}; \mathbf{z}) = \lim_{t \downarrow 0} \frac{\theta(\hat{\mathbf{x}} + t\mathbf{z}) - \theta(\hat{\mathbf{x}}) - t\theta'(\hat{\mathbf{x}}; \mathbf{z})}{\frac{1}{2}t^2}.$$

It is obvious that if θ is twice directionally differentiable at $\hat{\mathbf{x}}$, then for any $\mathbf{z} \in \mathbb{R}^n$ there exists $\delta > 0$ such that

$$\theta(\hat{\mathbf{x}} + t\mathbf{y}) = \theta(\hat{\mathbf{x}}) + t\theta'(\hat{\mathbf{x}}; \mathbf{y}) + \frac{1}{2}t^2\theta^{(2)}(\hat{\mathbf{x}}; \mathbf{z}) + o(t^2), \quad \forall t \in (0, \delta) \text{ and } \forall \mathbf{y} \in \mathcal{N}(\mathbf{z}, \delta),$$

and particularly

$$\theta(\hat{\mathbf{x}} + t\mathbf{z}) = \theta(\hat{\mathbf{x}}) + t\theta'(\hat{\mathbf{x}}; \mathbf{z}) + \frac{1}{2}t^2\theta^{(2)}(\hat{\mathbf{x}}; \mathbf{z}) + o(t^2), \quad \forall t \in (0, \delta).$$

Moreover, if θ is twice differentiable at $\hat{\mathbf{x}}$, then

$$\theta^{(2)}(\hat{\mathbf{x}}; \mathbf{z}) = \langle \nabla^2 \theta(\hat{\mathbf{x}}) \mathbf{z}, \mathbf{z} \rangle, \quad \forall \mathbf{z} \in \mathbb{R}^n.$$

From [7, 27], we also know that if θ is convex and twice directionally differentiable at $\hat{\mathbf{x}}$, then

$$\theta^{(2)}(\hat{\mathbf{x}}; \mathbf{z}) \geq 0, \quad \forall \mathbf{z} \in \mathbb{R}^n.$$

For a vector-valued function $\Phi : \mathbb{R}^n \rightarrow \mathbb{R}^m$ with component functions Φ_i for $i = 1, \dots, m$, $\Phi^{(2)}(\mathbf{x}; \mathbf{z})$ is defined to be the m -vector with components $\Phi_i^{(2)}(\mathbf{x}; \mathbf{z})$ for $i = 1, \dots, m$.

210 **Lemma 3.2** Let $\varrho : \mathbb{R}^m \rightarrow \mathbb{R}$ be locally Lipschitz continuous at $\Phi(\mathbf{x}) \in \mathbb{R}^m$, and $\Phi : \mathbb{R}^n \rightarrow$
 211 \mathbb{R}^m be locally Lipschitz continuous at $\mathbf{x} \in \mathbb{R}^n$, then the composite function $\theta = \varrho \circ \Phi : \mathbb{R}^n \rightarrow \mathbb{R}$
 212 is twice directionally differentiable at \mathbf{x} under either one of the following three conditions:

213 (a) ϱ is semismoothly differentiable at $\Phi(\mathbf{x})$ (i.e., ϱ is differentiable near $\Phi(\mathbf{x})$ and $\nabla \varrho$ is
 214 semismooth at $\Phi(\mathbf{x})$), and Φ is twice directionally differentiable at \mathbf{x} .

215 (b) ϱ is twice directionally differentiable at $\Phi(\mathbf{x})$ and Φ is piecewise affine near \mathbf{x} .

216 (c) ϱ is piecewise affine near $\Phi(\mathbf{x})$ and Φ is twice directionally differentiable at \mathbf{x} .

217 Moreover, we have, for all $\mathbf{z} \in \mathbb{R}^n$,

$$\theta^{(2)}(\mathbf{x}; \mathbf{z}) = \Phi'(\mathbf{x}; \mathbf{z})^\top (\nabla \varrho)'(\Phi(\mathbf{x}); \Phi'(\mathbf{x}; \mathbf{z})) + \nabla \varrho(\Phi(\mathbf{x}))^\top \Phi^{(2)}(\mathbf{x}; \mathbf{z}), \text{ if (a) holds; } \quad (3.2)$$

218

$$\theta^{(2)}(\mathbf{x}; \mathbf{z}) = \varrho^{(2)}(\Phi(\mathbf{x}); \Phi'(\mathbf{x}; \mathbf{z})), \text{ if (b) holds; } \quad (3.3)$$

219 and

$$\theta^{(2)}(\mathbf{x}; \mathbf{z}) = \varrho'(\Phi(\mathbf{x}); \Phi^{(2)}(\mathbf{x}; \mathbf{z})), \text{ if (c) holds. } \quad (3.4)$$

220 *Proof* Conclusions (3.2) and (3.3) have been proved in [7, Prop. 3.2]. It is easy to prove
 221 conclusion (3.4) under condition (c) by noting that $\varrho'(\mathbf{u}; \mathbf{v})$ exists and $\varrho^{(2)}(\mathbf{u}; \mathbf{v}) = 0$ at any
 222 point \mathbf{u} for any direction \mathbf{v} when ϱ is piecewise affine. \square

223 **Definition 3.3** [7] Let $\theta : \mathbb{R}^n \rightarrow \mathbb{R}$ be twice directionally differentiable at $\hat{\mathbf{x}} \in \mathbb{R}^n$. $\hat{\mathbf{x}}$ is
 224 called a second-order d-stationary point of θ if $\hat{\mathbf{x}}$ is a first-order d-stationary point of θ , and
 225 for any $\mathbf{z} \in \mathbb{R}^n$,

$$\theta'(\hat{\mathbf{x}}; \mathbf{z}) = 0 \text{ implies } \theta^{(2)}(\hat{\mathbf{x}}; \mathbf{z}) \geq 0.$$

226 According to [7, Theorem 1] and [27, Theorem 13.24], if θ is twice directionally dif-
 227 ferentiable, then second-order d-stationary points of θ have the following locally optimal
 228 properties.

229 **Proposition 3.4** Let $\theta : \mathbb{R}^n \rightarrow \mathbb{R}$ be twice directionally differentiable at $\hat{\mathbf{x}} \in \mathbb{R}^n$. The
 230 following two statements hold:

231 (i) If $\hat{\mathbf{x}} \in \mathbb{R}^n$ is a local minimizer of θ , then $\hat{\mathbf{x}}$ is a second-order d-stationary point of θ .

232 (ii) $\hat{\mathbf{x}} \in \mathbb{R}^n$ is a **strong** local minimizer of θ , i.e., there exist a neighborhood \mathcal{W} of $\hat{\mathbf{x}}$ and
 233 a scalar $\delta > 0$ such that

$$\theta(\mathbf{x}) \geq \theta(\hat{\mathbf{x}}) + \delta \|\mathbf{x} - \hat{\mathbf{x}}\|^2, \quad \forall \mathbf{x} \in \mathcal{W},$$

234 if and only if $\hat{\mathbf{x}}$ is a first-order d-stationary point of θ and satisfies that for any $\hat{\mathbf{x}} \neq \mathbf{x} \in \mathbb{R}^n$,

$$\theta'(\hat{\mathbf{x}}; \mathbf{x} - \hat{\mathbf{x}}) = 0 \text{ implies } \theta^{(2)}(\hat{\mathbf{x}}; \mathbf{x} - \hat{\mathbf{x}}) > 0.$$

235 In the following parts, we will use the second-order directional derivative of ℓ_2 vector
 236 norm function. Recall that $m(\mathbf{u}) = \|\mathbf{u}\|$, and that

$$m'(\mathbf{u}; \mathbf{v}) = \lim_{t \downarrow 0} \frac{\|\mathbf{u} + t\mathbf{v}\| - \|\mathbf{u}\|}{t} = \begin{cases} \|\mathbf{v}\|, & \text{if } \mathbf{u} = \mathbf{0}, \\ \frac{\langle \mathbf{u}, \mathbf{v} \rangle}{\|\mathbf{u}\|}, & \text{if } \mathbf{u} \neq \mathbf{0}, \end{cases} \quad \forall \mathbf{u}, \mathbf{v} \in \mathbb{R}^{d_i}.$$

237 It is easy to know that $m(\cdot)$ is twice differentiable at all points except $\mathbf{u} = \mathbf{0}$, and that

$$\begin{aligned} m^{(2)}(\mathbf{u}; \mathbf{w}) &= \lim_{\mathbf{v} \rightarrow \mathbf{w}, t \downarrow 0} \frac{\|\mathbf{u} + t\mathbf{v}\| - \|\mathbf{u}\| - tm'(\mathbf{u}; \mathbf{v})}{\frac{1}{2}t^2} \\ &= \begin{cases} 0, & \text{if } \mathbf{u} = \mathbf{0}, \\ \frac{(\|\mathbf{u}\|\|\mathbf{w}\|)^2 - |\langle \mathbf{u}, \mathbf{w} \rangle|^2}{\|\mathbf{u}\|^3}, & \text{if } \mathbf{u} \neq \mathbf{0}, \end{cases} \quad \forall \mathbf{u}, \mathbf{w} \in \mathbb{R}^{d_i}. \end{aligned} \quad (3.5)$$

238 **3.2 Second-order sufficient and necessary conditions for problem (1.1)**

239 To study second-order d-stationary points of problem (1.1), we need the following assump-
240 tion.

241 **Assumption (A2)** The penalty function $\varphi : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ is a DC function given by

$$\varphi(t) \triangleq g(t) - h(t), \quad (3.6)$$

242 where g is affine in $t \in [0, \infty)$ with $g'(0) := g'(0+)$, and h is convex and semismoothly
243 differentiable in $t \in (0, \infty)$ with $h'(0) := h'(0+)$.

244 We can easily check that several folded concave penalty functions satisfy Assumption
245 (A2), such as logarithm penalty, fraction penalty, HTP, MCP and SCAD.

246 In general the second-order directional derivative of a function is not easy to compute.
247 The following lemma provides an explicit formula for computing the second-order directional
248 derivative of the objective function of problem (1.1).

249 **Lemma 3.5** *Under Assumption (A2), the second-order directional derivative of the objective*
250 *function f in (1.1) has the following form*

$$\begin{aligned} f^{(2)}(\widehat{\mathbf{x}}; \mathbf{x} - \widehat{\mathbf{x}}) = & \langle \nabla^2 \mathcal{L}(\widehat{\mathbf{x}})(\mathbf{x} - \widehat{\mathbf{x}}), \mathbf{x} - \widehat{\mathbf{x}} \rangle + \sum_{i=1}^K \left[g'(\|\widehat{\mathbf{x}}_i\|) - h'(\|\widehat{\mathbf{x}}_i\|) \right] m^{(2)}(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) \\ & - \sum_{i=1}^K m'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) H'(\|\widehat{\mathbf{x}}_i\|; m'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i)), \end{aligned} \quad (3.7)$$

251 where $m'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i)$ is given by (2.9),

$$m^{(2)}(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) = \begin{cases} 0, & \text{if } i \notin I(\widehat{\mathbf{x}}), \\ \frac{(\|\mathbf{x}_i - \widehat{\mathbf{x}}_i\| \|\widehat{\mathbf{x}}_i\|)^2 - |\langle \widehat{\mathbf{x}}_i, \mathbf{x}_i - \widehat{\mathbf{x}}_i \rangle|^2}{\|\widehat{\mathbf{x}}_i\|^3}, & \text{if } i \in I(\widehat{\mathbf{x}}), \end{cases} \quad (3.8)$$

252 and $H(t) := h'(t)$ for any $t \in [0, +\infty)$.

253 *Proof* Since \mathcal{L} is twice continuously differentiable, $\mathcal{L}^{(2)}(\widehat{\mathbf{x}}; \mathbf{x} - \widehat{\mathbf{x}}) = \langle \nabla^2 \mathcal{L}(\widehat{\mathbf{x}})(\mathbf{x} - \widehat{\mathbf{x}}), \mathbf{x} - \widehat{\mathbf{x}} \rangle$.
254 Since g is affine in $[0, \infty)$ with $g'(0) = g'(0+)$, $(g \circ m)'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) = g'(\|\widehat{\mathbf{x}}_i\|) m'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i)$,
255 $g^{(2)}(\|\widehat{\mathbf{x}}_i\|; m'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i)) = 0$. By Lemma 3.2,

$$\begin{aligned} (g \circ m)^{(2)}(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) &= g'(\|\widehat{\mathbf{x}}_i\|) m^{(2)}(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) + g^{(2)}(\|\widehat{\mathbf{x}}_i\|; m'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i)) m'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) \\ &= g'(\|\widehat{\mathbf{x}}_i\|) m^{(2)}(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) \end{aligned}$$

256 for $i = 1, \dots, K$.

257 Since h is semismoothly differentiable in $(0, \infty)$ with $h'(0) = h'(0+)$, h is twice direction-
258 ally differentiable and $(h \circ m)'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) = h'(\|\widehat{\mathbf{x}}_i\|) m'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i)$. By Lemma 3.2,

$$\begin{aligned} (h \circ m)^{(2)}(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) &= h'(\|\widehat{\mathbf{x}}_i\|) m^{(2)}(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) + h^{(2)}(\|\widehat{\mathbf{x}}_i\|; m'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i)) m'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) \\ &= h'(\|\widehat{\mathbf{x}}_i\|) m^{(2)}(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) + H'(\|\widehat{\mathbf{x}}_i\|; m'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i)) m'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) \end{aligned}$$

259 for $i = 1, \dots, K$.

260 Then we have

$$\begin{aligned} f^{(2)}(\widehat{\mathbf{x}}; \mathbf{x} - \widehat{\mathbf{x}}) &= \langle \nabla^2 \mathcal{L}(\widehat{\mathbf{x}})(\mathbf{x} - \widehat{\mathbf{x}}), \mathbf{x} - \widehat{\mathbf{x}} \rangle + \sum_{i=1}^K \left[g'(\|\widehat{\mathbf{x}}_i\|) - h'(\|\widehat{\mathbf{x}}_i\|) \right] m^{(2)}(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) \\ &\quad - \sum_{i=1}^K m'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) H'(\|\widehat{\mathbf{x}}_i\|; m'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i)), \end{aligned}$$

261 where $m'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i)$ and $m^{(2)}(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i)$ are given by (2.5) and (3.5) respectively. \square

262 From Definition 3.3, Proposition 3.4 and Lemma 3.5, we obtain the following theorem.

263 **Theorem 3.6** *Suppose Assumption (A2) holds and $\widehat{\mathbf{x}} \in \mathbb{R}^n$ is a first-order d-stationary point*
 264 *of problem (1.1), then the following two statements hold with $f'(\widehat{\mathbf{x}}; \mathbf{x} - \widehat{\mathbf{x}})$ and $f^{(2)}(\widehat{\mathbf{x}}; \mathbf{x} - \widehat{\mathbf{x}})$*
 265 *given by (2.8) and (3.7) respectively.*

266 (i) $\widehat{\mathbf{x}}$ is a second-order d-stationary point of problem (1.1) if and only if for any $\mathbf{x} \in \mathbb{R}^n$,
 267 $f'(\widehat{\mathbf{x}}; \mathbf{x} - \widehat{\mathbf{x}}) = 0$ implies $f^{(2)}(\widehat{\mathbf{x}}; \mathbf{x} - \widehat{\mathbf{x}}) \geq 0$.

268 (ii) $\widehat{\mathbf{x}}$ is a strong local minimizer of problem (1.1) if and only if for any $\widehat{\mathbf{x}} \neq \mathbf{x} \in \mathbb{R}^n$,
 269 $f'(\widehat{\mathbf{x}}; \mathbf{x} - \widehat{\mathbf{x}}) = 0$ implies $f^{(2)}(\widehat{\mathbf{x}}; \mathbf{x} - \widehat{\mathbf{x}}) > 0$.

270 The following theorem shows that the second-order directional derivative at a second-
 271 order d-stationary point can be simplified and is nonnegative on a special set.

272 **Theorem 3.7** *Under Assumption (A2), let $\widehat{\mathbf{x}} \in \mathbb{R}^n$ be a second-order d-stationary point of*
 273 *problem (1.1), and*

$$\mathcal{X}(\widehat{\mathbf{x}}) = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{x}_i = \mathbf{0} \text{ whenever } i \notin I(\widehat{\mathbf{x}})\}, \quad (3.9)$$

274 then for any $\mathbf{x} \in \mathcal{X}(\widehat{\mathbf{x}})$,

$$\langle \nabla^2 \mathcal{L}(\widehat{\mathbf{x}})(\mathbf{x} - \widehat{\mathbf{x}}), \mathbf{x} - \widehat{\mathbf{x}} \rangle + \sum_{i \in I(\widehat{\mathbf{x}})} \left[(g \circ m)^{(2)}(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) - (h \circ m)^{(2)}(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) \right] \geq 0,$$

275 where for $i \in I(\widehat{\mathbf{x}})$,

$$\begin{aligned} (g \circ m)^{(2)}(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) &= g'(\|\widehat{\mathbf{x}}_i\|)m^{(2)}(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i), \\ (h \circ m)^{(2)}(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) &= m'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i)H'(\|\widehat{\mathbf{x}}_i\|; m'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i)) \\ &\quad + h'(\|\widehat{\mathbf{x}}_i\|)m^{(2)}(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i), \\ m'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) &= \frac{\langle \widehat{\mathbf{x}}_i, \mathbf{x}_i - \widehat{\mathbf{x}}_i \rangle}{\|\widehat{\mathbf{x}}_i\|}, \\ m^{(2)}(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) &= \frac{(\|\mathbf{x}_i - \widehat{\mathbf{x}}_i\| \|\widehat{\mathbf{x}}_i\|)^2 - |\langle \widehat{\mathbf{x}}_i, \mathbf{x}_i - \widehat{\mathbf{x}}_i \rangle|^2}{\|\widehat{\mathbf{x}}_i\|^3}, \\ H(t) &= h'(t) \text{ for any } t \in (0, \infty). \end{aligned}$$

276 *Proof* Since $\widehat{\mathbf{x}}$ is a second-order d-stationary point of problem (1.1), it is also a first-order
 277 d-stationary point of problem (1.1), which means

$$\langle \nabla \mathcal{L}(\widehat{\mathbf{x}}), \mathbf{x} - \widehat{\mathbf{x}} \rangle + \sum_{i=1}^K (g'(\|\widehat{\mathbf{x}}_i\|) - h'(\|\widehat{\mathbf{x}}_i\|))m'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) \geq 0, \quad \forall \mathbf{x} \in \mathbb{R}^n.$$

278 By the same argument in the proof of (2.13), we have

$$[\nabla \mathcal{L}(\widehat{\mathbf{x}})]_i + \frac{g'(\|\widehat{\mathbf{x}}_i\|) - h'(\|\widehat{\mathbf{x}}_i\|)}{\|\widehat{\mathbf{x}}_i\|} \widehat{\mathbf{x}}_i = \mathbf{0}, \quad \forall i \in I(\widehat{\mathbf{x}}).$$

279 Therefore, we get

$$\sum_{i \in I(\widehat{\mathbf{x}})} \left\langle [\nabla \mathcal{L}(\widehat{\mathbf{x}})]_i + \frac{g'(\|\widehat{\mathbf{x}}_i\|) - h'(\|\widehat{\mathbf{x}}_i\|)}{\|\widehat{\mathbf{x}}_i\|} \widehat{\mathbf{x}}_i, \mathbf{x}_i - \widehat{\mathbf{x}}_i \right\rangle = 0, \quad \forall \mathbf{x} \in \mathbb{R}^n. \quad (3.10)$$

280 For any $\mathbf{x} \in \mathcal{X}(\widehat{\mathbf{x}})$, by (2.9), (3.10) and direct computation, we obtain

$$\langle \nabla \mathcal{L}(\widehat{\mathbf{x}}), \mathbf{x} - \widehat{\mathbf{x}} \rangle + \sum_{i=1}^K (g'(\|\widehat{\mathbf{x}}_i\|) - h'(\|\widehat{\mathbf{x}}_i\|))m'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) = 0,$$

281 that is, $f'(\widehat{\mathbf{x}}; \mathbf{x} - \widehat{\mathbf{x}}) = 0$, which together with that $\widehat{\mathbf{x}}$ is a second-order d-stationary point of
 282 problem (1.1) yields that $f^{(2)}(\widehat{\mathbf{x}}; \mathbf{x} - \widehat{\mathbf{x}}) \geq 0$. From (2.9), (3.8) and $\mathbf{x} \in \mathcal{X}(\widehat{\mathbf{x}})$, we have that
 283 for $i \notin I(\widehat{\mathbf{x}})$, $m'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) = m^{(2)}(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) = 0$, and that for $i \in I(\widehat{\mathbf{x}})$,

$$\begin{aligned} m'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) &= \frac{\langle \widehat{\mathbf{x}}_i, \mathbf{x}_i - \widehat{\mathbf{x}}_i \rangle}{\|\widehat{\mathbf{x}}_i\|}, \\ m^{(2)}(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) &= \frac{(\|\mathbf{x}_i - \widehat{\mathbf{x}}_i\| \|\widehat{\mathbf{x}}_i\|)^2 - |\langle \widehat{\mathbf{x}}_i, \mathbf{x}_i - \widehat{\mathbf{x}}_i \rangle|^2}{\|\widehat{\mathbf{x}}_i\|^3}, \\ (g \circ m)^{(2)}(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) &= g'(\|\widehat{\mathbf{x}}_i\|) m^{(2)}(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i), \\ (h \circ m)^{(2)}(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) &= m'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) H'(\|\widehat{\mathbf{x}}_i\|; m'(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i)) \\ &\quad + h'(\|\widehat{\mathbf{x}}_i\|) m^{(2)}(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i), \end{aligned}$$

284 where $H(t) = h'(t)$ for any $t \in (0, \infty)$. Hence we get

$$\begin{aligned} &\langle \nabla^2 \mathcal{L}(\widehat{\mathbf{x}})(\mathbf{x} - \widehat{\mathbf{x}}), \mathbf{x} - \widehat{\mathbf{x}} \rangle + \sum_{i \in I(\widehat{\mathbf{x}})} \left[(g \circ m)^{(2)}(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) - (h \circ m)^{(2)}(\widehat{\mathbf{x}}_i; \mathbf{x}_i - \widehat{\mathbf{x}}_i) \right] \\ &= f^{(2)}(\widehat{\mathbf{x}}; \mathbf{x} - \widehat{\mathbf{x}}) \geq 0. \end{aligned}$$

285 The proof is finished. □

286 3.3 Lower bound theory of second-order d-stationary points

287 In this subsection, we analyze the lower bound of the ℓ_2 vector norm of nonzero groups of
 288 second-order d-stationary points of problem (1.1). We will see that the second-order lower
 289 bounds are tighter than the corresponding first-order lower bounds. At first, we give a useful
 290 lemma which provides an upper bound for the second-order directional derivative of the
 291 penalty function h at any second-order d-stationary point.

292 **Lemma 3.8** *Under Assumption (A2), let $\widehat{\mathbf{x}} \in \mathbb{R}^n$ be a second-order d-stationary point of*
 293 *problem (1.1), then*

$$\langle \nabla_i^2 \mathcal{L}(\widehat{\mathbf{x}}) \widehat{\mathbf{x}}_i, \widehat{\mathbf{x}}_i \rangle \geq \|\widehat{\mathbf{x}}_i\|^2 \cdot \max\{H'(\|\widehat{\mathbf{x}}_i\|; 1), -H'(\|\widehat{\mathbf{x}}_i\|; -1)\}, \quad \forall i \in I(\widehat{\mathbf{x}}),$$

294 where $\nabla_i^2 \mathcal{L}(\mathbf{x})$ denotes the principal submatrix of $\nabla^2 \mathcal{L}(\mathbf{x})$ corresponding to the group \mathbf{x}_i .

Proof For each fixed $i \in I(\widehat{\mathbf{x}})$, let $\mathbf{x}^1, \mathbf{x}^2 \in \mathbb{R}^n$ be taken as

$$\mathbf{x}_{i'}^1 = \begin{cases} 2\widehat{\mathbf{x}}_i, & \text{if } i' = i, \\ \widehat{\mathbf{x}}_{i'}, & \text{if } i' \neq i, \end{cases} \quad \mathbf{x}_{i'}^2 = \begin{cases} \mathbf{0}, & \text{if } i' = i, \\ \widehat{\mathbf{x}}_{i'}, & \text{if } i' \neq i. \end{cases}$$

295 Then it is easy to check that $\mathbf{x}^1, \mathbf{x}^2 \in \mathcal{X}(\widehat{\mathbf{x}})$ which has been defined by (3.9). By Theorem
 296 3.7, we have

$$\begin{aligned} &\langle \nabla^2 \mathcal{L}(\widehat{\mathbf{x}})(\mathbf{x}^\eta - \widehat{\mathbf{x}}), \mathbf{x}^\eta - \widehat{\mathbf{x}} \rangle + \sum_{i' \in I(\widehat{\mathbf{x}})} \left[(g \circ m)^{(2)}(\widehat{\mathbf{x}}_{i'}; \mathbf{x}_{i'}^\eta - \widehat{\mathbf{x}}_{i'}) - (h \circ m)^{(2)}(\widehat{\mathbf{x}}_{i'}; \mathbf{x}_{i'}^\eta - \widehat{\mathbf{x}}_{i'}) \right] \\ &\geq 0, \quad \eta = 1, 2, \end{aligned} \tag{3.11}$$

297 where, according to the definitions of \mathbf{x}^1 and \mathbf{x}^2 as well as formulas (2.9) and (3.8),

$$\begin{aligned} \langle \nabla^2 \mathcal{L}(\widehat{\mathbf{x}})(\mathbf{x}^\eta - \widehat{\mathbf{x}}), \mathbf{x}^\eta - \widehat{\mathbf{x}} \rangle &= \langle \nabla_i^2 \mathcal{L}(\widehat{\mathbf{x}}) \widehat{\mathbf{x}}_i, \widehat{\mathbf{x}}_i \rangle, \quad \eta = 1, 2, \\ m'(\widehat{\mathbf{x}}_{i'}; \mathbf{x}_{i'}^1 - \widehat{\mathbf{x}}_{i'}) &= \begin{cases} \|\widehat{\mathbf{x}}_i\|, & \text{if } i' = i, \\ 0, & \text{if } i' \neq i, \end{cases} \quad m'(\widehat{\mathbf{x}}_{i'}; \mathbf{x}_{i'}^2 - \widehat{\mathbf{x}}_{i'}) = \begin{cases} -\|\widehat{\mathbf{x}}_i\|, & \text{if } i' = i, \\ 0, & \text{if } i' \neq i, \end{cases} \\ m^{(2)}(\widehat{\mathbf{x}}_{i'}; \mathbf{x}_{i'}^\eta - \widehat{\mathbf{x}}_{i'}) &= 0, \quad \forall i' = 1, \dots, K, \quad \eta = 1, 2, \\ (g \circ m)^{(2)}(\widehat{\mathbf{x}}_{i'}; \mathbf{x}_{i'}^\eta - \widehat{\mathbf{x}}_{i'}) &= g'(\|\widehat{\mathbf{x}}_{i'}\|) m^{(2)}(\widehat{\mathbf{x}}_{i'}; \mathbf{x}_{i'}^\eta - \widehat{\mathbf{x}}_{i'}) = 0, \quad \eta = 1, 2, \\ (h \circ m)^{(2)}(\widehat{\mathbf{x}}_{i'}; \mathbf{x}_{i'}^\eta - \widehat{\mathbf{x}}_{i'}) &= m'(\widehat{\mathbf{x}}_{i'}; \mathbf{x}_{i'}^\eta - \widehat{\mathbf{x}}_{i'}) H'(\|\widehat{\mathbf{x}}_{i'}\|; m'(\widehat{\mathbf{x}}_{i'}; \mathbf{x}_{i'}^\eta - \widehat{\mathbf{x}}_{i'})) \\ &\quad + h'(\|\widehat{\mathbf{x}}_{i'}\|) m^{(2)}(\widehat{\mathbf{x}}_{i'}; \mathbf{x}_{i'}^\eta - \widehat{\mathbf{x}}_{i'}) \\ &= m'(\widehat{\mathbf{x}}_{i'}; \mathbf{x}_{i'}^\eta - \widehat{\mathbf{x}}_{i'}) H'(\|\widehat{\mathbf{x}}_{i'}\|; m'(\widehat{\mathbf{x}}_{i'}; \mathbf{x}_{i'}^\eta - \widehat{\mathbf{x}}_{i'})), \quad \eta = 1, 2, \\ H(t) &= h'(t) \text{ for any } t \in (0, \infty). \end{aligned}$$

298 Therefore, by taking the above terms into inequality (3.11), we get

$$\langle \nabla_i^2 \mathcal{L}(\widehat{\mathbf{x}}) \widehat{\mathbf{x}}_i, \widehat{\mathbf{x}}_i \rangle \geq \|\widehat{\mathbf{x}}_i\| \cdot \max\{H'(\|\widehat{\mathbf{x}}_i\|; \|\widehat{\mathbf{x}}_i\|), -H'(\|\widehat{\mathbf{x}}_i\|; -\|\widehat{\mathbf{x}}_i\|)\}, \quad \forall i \in I(\widehat{\mathbf{x}}).$$

299 By the positive homogeneity of $H'(\|\widehat{\mathbf{x}}_i\|; \cdot)$ and $\|\widehat{\mathbf{x}}_i\| > 0$, we derive the desired result. \square

300 **Theorem 3.9** Suppose Assumption (A2) holds and there exists $M > 0$ such that $\|\nabla^2 \mathcal{L}(\mathbf{x})\|_2 \leq$
301 M for all $\mathbf{x} \in \mathbb{R}^n$. Let $\widehat{\mathbf{x}} \in \mathbb{R}^n$ be a second-order d-stationary point of problem (1.1), then
302 the following statements hold:

303 (i) For MCP, if $M < \frac{1}{\alpha}$, then either $\|\widehat{\mathbf{x}}_i\| = 0$ or $\|\widehat{\mathbf{x}}_i\| > \alpha\lambda$, $i = 1, \dots, K$.

304 (ii) For SCAD, if $M < \frac{1}{\alpha-1}$, then either $\|\widehat{\mathbf{x}}_i\| < \lambda$ or $\|\widehat{\mathbf{x}}_i\| > \alpha\lambda$, $i = 1, \dots, K$.

305 (iii) For SCAD, suppose, in addition, there exists a nondecreasing function $C : \mathbb{R} \rightarrow \mathbb{R}_+$
306 such that $\|\nabla \mathcal{L}(\mathbf{x})\| \leq C(\mathcal{L}(\mathbf{x}))$ for all $\mathbf{x} \in \mathbb{R}^n$. If there exists $\mathbf{x}^0 \in \mathbb{R}^n$ satisfying $\mathcal{L}(\mathbf{x}^0) \geq$
307 $\mathcal{L}(\widehat{\mathbf{x}})$, $\varphi'(0) > C(\mathcal{L}(\mathbf{x}^0))$, and $\frac{1}{\alpha-1} > M$, then either $\|\widehat{\mathbf{x}}_i\| = 0$ or $\|\widehat{\mathbf{x}}_i\| > \alpha\lambda$, $i = 1, \dots, K$.

308 *Proof* Since $\|\nabla^2 \mathcal{L}(\mathbf{x})\|_2 \leq M$ for all $\mathbf{x} \in \mathbb{R}^n$, we have

$$\langle \nabla^2 \mathcal{L}(\mathbf{x}) \mathbf{z}, \mathbf{z} \rangle \leq M \|\mathbf{z}\|^2, \quad \forall \mathbf{x}, \mathbf{z} \in \mathbb{R}^n. \quad (3.12)$$

(i) For MCP: recall that $H^{\text{MCP}}(t) = (h^{\text{MCP}})'(t) = \begin{cases} \frac{t}{\alpha}, & \text{if } 0 \leq t \leq \alpha\lambda, \\ \lambda, & \text{if } t > \alpha\lambda, \end{cases}$ we have

$$(H^{\text{MCP}})'(t; 1) = \begin{cases} \frac{1}{\alpha}, & \text{if } 0 \leq t < \alpha\lambda, \\ 0, & \text{if } t \geq \alpha\lambda, \end{cases} \quad (H^{\text{MCP}})'(t; -1) = \begin{cases} -\frac{1}{\alpha}, & \text{if } 0 < t \leq \alpha\lambda, \\ 0, & \text{if } t > \alpha\lambda. \end{cases}$$

Assume, on the contrary, that $0 < \|\widehat{\mathbf{x}}_i\| \leq \alpha\lambda$, then

$$(H^{\text{MCP}})'(\|\widehat{\mathbf{x}}_i\|; 1) \leq -(H^{\text{MCP}})'(\|\widehat{\mathbf{x}}_i\|; -1) = \frac{1}{\alpha}.$$

309 From Lemma 3.8 and (3.12), we have

$$M \geq -(H^{\text{MCP}})'(\|\widehat{\mathbf{x}}_i\|; -1) = \frac{1}{\alpha},$$

310 which contradicts the condition $M < \frac{1}{\alpha}$. Therefore, we have $\|\widehat{\mathbf{x}}_i\| > \alpha\lambda$ for any $i \in I(\widehat{\mathbf{x}})$.

311 (ii) For SCAD: recall that $H^{\text{SCAD}}(t) = (h^{\text{SCAD}})'(t) = \begin{cases} 0, & \text{if } 0 \leq t \leq \lambda, \\ \frac{t-\lambda}{\alpha-1}, & \text{if } \lambda < t \leq \alpha\lambda, \\ \lambda, & \text{if } t > \alpha\lambda, \end{cases}$ we have

$$(H^{\text{SCAD}})'(t; 1) = \begin{cases} 0, & \text{if } t \in [0, \lambda) \cup [\alpha\lambda, +\infty), \\ \frac{1}{\alpha-1}, & \text{if } t \in [\lambda, \alpha\lambda), \end{cases}$$

$$(H^{\text{SCAD}})'(t; -1) = \begin{cases} 0, & \text{if } t \in (0, \lambda] \cup (\alpha\lambda, +\infty), \\ -\frac{1}{\alpha-1}, & \text{if } t \in (\lambda, \alpha\lambda]. \end{cases}$$

Assume, on the contrary, that $\lambda \leq \|\widehat{\mathbf{x}}_i\| \leq \alpha\lambda$, then

$$\max\{(H^{\text{SCAD}})'(\|\widehat{\mathbf{x}}_i\|; 1), -(H^{\text{SCAD}})'(\|\widehat{\mathbf{x}}_i\|; -1)\} = \frac{1}{\alpha-1}.$$

312 From Lemma 3.8 and (3.12), we have

$$M \geq \frac{1}{\alpha-1},$$

313 which contradicts the condition $M < \frac{1}{\alpha-1}$. Therefore, we have either $\|\widehat{\mathbf{x}}_i\| < \lambda$ or $\|\widehat{\mathbf{x}}_i\| > \alpha\lambda$.

314 (iii) Since $\widehat{\mathbf{x}}$ is a second-order d-stationary point of problem (1.1), it is also a first-order
315 d-stationary point of problem (1.1). Combining (ii) with Corollary 2.5 (iii), we derive the
316 desired result. \square

317 **Remark 3.10** *The condition in Theorem 3.9 means that the operator $\nabla^2 \mathcal{L}(\mathbf{x})$ has an uni-*
318 *form bound M on \mathbb{R}^n . We can easily check that $\mathcal{L}(\mathbf{x}) = \frac{1}{2m} \|A\mathbf{x} - b\|^2$ satisfies this condition*
319 *since $\|\nabla^2 \mathcal{L}(\mathbf{x})\|_2 = \frac{\|A^\top A\|_2}{m} = \frac{\|A\|_2^2}{m}$.*

320 4 Smoothing functions and consistency of stationary points

321 As we have seen, first-order and second-order d-stationary points have good locally optimal
322 properties. How to compute such points is an interesting and challenging problem. Smooth
323 approximations are widely used in optimization and scientific computing, e.g., see [8, 9, 10]. In
324 this section, we construct a twice continuously differentiable smoothing function of the objec-
325 tive function f of problem (1.1), and show that the first-order and second-order d-stationary
326 points of problem (1.1) can be obtained via the first-order and second-order stationary points
327 of the smoothing problem. We should notice that in problem (1.1), the term $\varphi(\|\widehat{\mathbf{x}}_i\|)$ is a
328 composite of two nonsmooth functions φ and $\|\cdot\|$. Using the special structure of these two
329 functions, our smoothing function can be easily constructed.

330 For $\mu \in (0, \infty)$ and $m(\mathbf{u}) = \|\mathbf{u}\|$, let

$$\tilde{m}_\mu(\mathbf{u}) = \sqrt{\|\mathbf{u}\|^2 + \mu}, \quad \forall \mathbf{u} \in \mathbb{R}^{d_i}, \quad (4.1)$$

331 then $\tilde{m}_\mu(\mathbf{u})$ is always positive and twice continuously differentiable with

$$\nabla \tilde{m}_\mu(\mathbf{u}) = \frac{\mathbf{u}}{\sqrt{\|\mathbf{u}\|^2 + \mu}}, \quad \nabla^2 \tilde{m}_\mu(\mathbf{u}) = \frac{(\|\mathbf{u}\|^2 + \mu)\mathbf{I} - \mathbf{u}\mathbf{u}^\top}{(\|\mathbf{u}\|^2 + \mu)^{3/2}}, \quad (4.2)$$

332 and

$$0 < \tilde{m}_\mu(\mathbf{u}) - m(\mathbf{u}) = \sqrt{\|\mathbf{u}\|^2 + \mu} - \|\mathbf{u}\| \leq \mu^{\frac{1}{2}}, \quad (4.3)$$

333 where \mathbf{I} denotes the identity matrix. One can also check that $\tilde{m}_\mu(\mathbf{u})$ satisfies the following
 334 three properties:

- 335 (i) $\lim_{\mathbf{v} \rightarrow \mathbf{u}, \mu \downarrow 0} \tilde{m}_\mu(\mathbf{v}) = m(\mathbf{u})$ for all $\mathbf{u} \in \mathbb{R}^{d_i}$;
 336 (ii) (Consistency or weak consistency of directional derivatives)

$$\lim_{\mathbf{v} \rightarrow \mathbf{u}, \mu \downarrow 0} \langle \nabla \tilde{m}_\mu(\mathbf{v}), \mathbf{w} \rangle = \langle \nabla m(\mathbf{u}), \mathbf{w} \rangle = m'(\mathbf{u}; \mathbf{w}), \quad \forall \mathbf{0} \neq \mathbf{u} \in \mathbb{R}^{d_i}, \forall \mathbf{w} \in \mathbb{R}^{d_i}, \quad (4.4)$$

$$\begin{aligned} \limsup_{\mathbf{v} \rightarrow \mathbf{0}, \mu \downarrow 0} \langle \nabla \tilde{m}_\mu(\mathbf{v}), \mathbf{w} \rangle &= \limsup_{\mathbf{v} \rightarrow \mathbf{0}, \mu \downarrow 0} \frac{\langle \mathbf{v}, \mathbf{w} \rangle}{\sqrt{\|\mathbf{v}\|^2 + \mu}} = \limsup_{t \downarrow 0, \mu \downarrow 0} \frac{t \|\mathbf{w}\|^2}{\sqrt{t^2 \|\mathbf{w}\|^2 + \mu}} \\ &= \limsup_{t \downarrow 0, \mu \downarrow 0} \frac{\|\mathbf{w}\|^2}{\sqrt{\|\mathbf{w}\|^2 + \frac{\mu}{t^2}}} = \|\mathbf{w}\| = m'(\mathbf{0}; \mathbf{w}), \quad \forall \mathbf{w} \in \mathbb{R}^{d_i}; \end{aligned} \quad (4.5)$$

- 337 (iii) (Consistency or weak consistency of second-order directional derivatives)

$$\begin{aligned} \lim_{\mathbf{v} \rightarrow \mathbf{u}, \mu \downarrow 0} \langle \nabla^2 \tilde{m}_\mu(\mathbf{v}) \mathbf{w}, \mathbf{w} \rangle &= \langle \nabla^2 m(\mathbf{u}) \mathbf{w}, \mathbf{w} \rangle \\ &= m^{(2)}(\mathbf{u}; \mathbf{w}), \quad \forall \mathbf{0} \neq \mathbf{u} \in \mathbb{R}^{d_i}, \forall \mathbf{w} \in \mathbb{R}^{d_i}, \end{aligned} \quad (4.6)$$

$$\begin{aligned} \liminf_{\mathbf{v} \rightarrow \mathbf{0}, \mu \downarrow 0} \langle \nabla^2 \tilde{m}_\mu(\mathbf{v}) \mathbf{w}, \mathbf{w} \rangle &= \liminf_{\mathbf{v} \rightarrow \mathbf{0}, \mu \downarrow 0} \frac{\|\mathbf{v}\|^2 \|\mathbf{w}\|^2 - (\mathbf{v}^\top \mathbf{w})^2 + \mu \|\mathbf{w}\|^2}{(\|\mathbf{v}\|^2 + \mu)^{3/2}} \\ &= \liminf_{\mathbf{v} \rightarrow \mathbf{0}, \mu \downarrow 0} \frac{\|\mathbf{w}\|^2}{\left(\frac{\|\mathbf{v}\|^2}{\mu^{2/3}} + \mu^{1/3}\right)^{3/2}} \\ &= 0 = m^{(2)}(\mathbf{0}; \mathbf{w}), \quad \forall \mathbf{w} \in \mathbb{R}^{d_i}. \end{aligned} \quad (4.7)$$

338 Under Assumption (A2), h is semismoothly differentiable in $(0, \infty)$. If h is not twice con-
 339 tinuously differentiable in $(0, \infty)$, for each $\mu > 0$, let \tilde{h}_μ be a twice continuously differentiable
 340 function in $(0, \infty)$ such that

$$\lim_{s \rightarrow t, \mu \downarrow 0} \tilde{h}_\mu(s) = h(t), \quad \lim_{s \rightarrow t, \mu \downarrow 0} \tilde{h}'_\mu(s) = h'(t), \quad \lim_{s \downarrow 0, \mu \downarrow 0} \tilde{h}'_\mu(s) = h'(0+), \quad (4.8)$$

$$\liminf_{s \rightarrow t, \mu \downarrow 0} \tilde{h}''_\mu(s) = \min\{H'(t; 1), -H'(t; -1)\}, \quad \text{and} \quad (4.9)$$

$$\limsup_{s \rightarrow t, \mu \downarrow 0} \tilde{h}''_\mu(s) = \max\{H'(t; 1), -H'(t; -1)\}. \quad (4.10)$$

341 Note that if h is twice continuously differentiable at $t > 0$, then $H'(t; 1) = -H'(t; -1) =$
 342 $h''(t)$.

For example, in MCP,

$$h^{\text{MCP}}(t) = \begin{cases} \frac{t^2}{2\alpha}, & \text{if } 0 \leq t \leq \alpha\lambda, \\ \lambda t - \frac{\alpha\lambda^2}{2}, & \text{if } t > \alpha\lambda, \end{cases} = \lambda t - \lambda \int_0^t \left(1 - \frac{\tau}{\alpha\lambda}\right)_+ d\tau \quad (\alpha > 1, \lambda > 0).$$

343 Let

$$\tilde{h}_\mu^{\text{MCP}}(t) = \lambda t - \frac{\lambda}{2} \int_0^t \left[\left(\left(1 - \frac{\tau}{\alpha\lambda}\right)^2 + \mu \right)^{1/2} + \left(1 - \frac{\tau}{\alpha\lambda}\right) \right] d\tau, \quad (4.11)$$

344 then one can check that for each $\mu > 0$, $\tilde{h}_\mu^{\text{MCP}}$ is twice continuously differentiable in $t \in (0, \infty)$
 345 with

$$\begin{aligned} (\tilde{h}_\mu^{\text{MCP}})'(t) &= \lambda - \frac{\lambda}{2} \left[\left(\left(1 - \frac{t}{\alpha\lambda}\right)^2 + \mu \right)^{1/2} + \left(1 - \frac{t}{\alpha\lambda}\right) \right], \\ (\tilde{h}_\mu^{\text{MCP}})''(t) &= \frac{1}{2\alpha} \left[\frac{1 - \frac{t}{\alpha\lambda}}{\sqrt{\left(1 - \frac{t}{\alpha\lambda}\right)^2 + \mu}} + 1 \right], \end{aligned}$$

346 and satisfies the following three properties:

- 347 (i) $\lim_{s \rightarrow t, \mu \downarrow 0} \tilde{h}_\mu^{\text{MCP}}(s) = h^{\text{MCP}}(t)$ for all $t \in [0, \infty)$;
 348 (ii) $\lim_{s \rightarrow t, \mu \downarrow 0} (\tilde{h}_\mu^{\text{MCP}})'(s) = (h^{\text{MCP}})'(t)$ for all $t \in (0, \infty)$, and $\lim_{s \downarrow 0, \mu \downarrow 0} (\tilde{h}_\mu^{\text{MCP}})'(s) = (h^{\text{MCP}})'(0+)$;
 349 (iii) For any $t \in (0, \alpha\lambda) \cup (\alpha\lambda, \infty)$,

$$\begin{aligned} \lim_{s \rightarrow t, \mu \downarrow 0} (\tilde{h}_\mu^{\text{MCP}})''(s) &= \lim_{s \rightarrow t, \mu \downarrow 0} \frac{1}{2\alpha} \left[\frac{1 - \frac{s}{\alpha\lambda}}{\sqrt{\left(1 - \frac{s}{\alpha\lambda}\right)^2 + \mu}} + 1 \right] \\ &= \lim_{s \rightarrow t, \mu \downarrow 0} \frac{1}{2\alpha} \left[\frac{\text{sign}\left(1 - \frac{s}{\alpha\lambda}\right)}{\sqrt{1 + \frac{\mu}{\left(1 - \frac{s}{\alpha\lambda}\right)^2}}} + 1 \right] = \begin{cases} \frac{1}{\alpha}, & \text{if } 0 < t < \alpha\lambda, \\ 0, & \text{if } t > \alpha\lambda, \end{cases} \\ &= (H^{\text{MCP}})'(t; 1) = -(H^{\text{MCP}})'(t; -1) = (h^{\text{MCP}})''(t); \end{aligned}$$

350 for $t = \alpha\lambda$,

$$\begin{aligned} \liminf_{s \rightarrow t, \mu \downarrow 0} (\tilde{h}_\mu^{\text{MCP}})''(s) &= \liminf_{s \rightarrow t, \mu \downarrow 0} \frac{1}{2\alpha} \left[\frac{\text{sign}\left(1 - \frac{s}{\alpha\lambda}\right)}{\sqrt{1 + \frac{\mu}{\left(1 - \frac{s}{\alpha\lambda}\right)^2}}} + 1 \right] = 0 = (H^{\text{MCP}})'(t; 1), \text{ and} \\ \limsup_{s \rightarrow t, \mu \downarrow 0} (\tilde{h}_\mu^{\text{MCP}})''(s) &= \limsup_{s \rightarrow t, \mu \downarrow 0} \frac{1}{2\alpha} \left[\frac{\text{sign}\left(1 - \frac{s}{\alpha\lambda}\right)}{\sqrt{1 + \frac{\mu}{\left(1 - \frac{s}{\alpha\lambda}\right)^2}}} + 1 \right] = \frac{1}{\alpha} = -(H^{\text{MCP}})'(t; -1). \end{aligned}$$

351 Now, under Assumption (A2), we have a twice continuously differentiable approximation
 352 $\tilde{f}_\mu(\mathbf{x})$ of the objective function $f(\mathbf{x})$ in problem (1.1),

$$\tilde{f}_\mu(\mathbf{x}) = \mathcal{L}(\mathbf{x}) + \sum_{i=1}^K \left[g \circ \tilde{m}_\mu(\mathbf{x}_i) - \tilde{h}_\mu \circ \tilde{m}_\mu(\mathbf{x}_i) \right],$$

353 with $\lim_{\mathbf{z} \rightarrow \mathbf{x}, \mu \downarrow 0} \tilde{f}_\mu(\mathbf{z}) = f(\mathbf{x})$ for any $\mathbf{x} \in \mathbb{R}^n$. It should be noted that although $g(\|\cdot\|) - \tilde{h}_\mu(\|\cdot\|)$
 354 is not differentiable at $\mathbf{x}_i = \mathbf{0}$, $g \circ \tilde{m}_\mu(\cdot) - \tilde{h}_\mu \circ \tilde{m}_\mu(\cdot)$ is twice continuously differentiable at any
 355 point $\mathbf{x}_i \in \mathbb{R}^{d_i}$ since $\tilde{m}_\mu(\mathbf{x}_i)$ is always strictly positive for any $\mu > 0$. Consequently, $\tilde{f}_\mu(\cdot)$ is
 356 twice continuously differentiable at any point $\mathbf{x} \in \mathbb{R}^n$. Thus we obtain a twice continuously
 357 differentiable optimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^n} \tilde{f}_\mu(\mathbf{x}). \quad (4.12)$$

358 By the standard definitions for twice differentiable optimization problems, $\hat{\mathbf{x}}^\mu$ is called a first-
 359 order stationary point of problem (4.12) if $\nabla \tilde{f}_\mu(\hat{\mathbf{x}}^\mu) = \mathbf{0}$; and $\hat{\mathbf{x}}^\mu$ is called a second-order
 360 stationary point of problem (4.12) if $\nabla \tilde{f}_\mu(\hat{\mathbf{x}}^\mu) = \mathbf{0}$ and

$$\langle \nabla^2 \tilde{f}_\mu(\hat{\mathbf{x}}^\mu) \mathbf{z}, \mathbf{z} \rangle \geq 0, \quad \forall \mathbf{z} \in \mathbb{R}^n. \quad (4.13)$$

Let $\{\widehat{\mathbf{x}}^{\mu_k}\}$ denote a sequence of first-order or second-order stationary points of problem (4.12) with $\mu_k > 0$, $k = 1, 2, \dots$, and $\mu_k \rightarrow 0$ as $k \rightarrow \infty$. We will investigate the accumulation points of $\{\widehat{\mathbf{x}}^{\mu_k}\}$.

Theorem 4.1 (*Consistency of first-order stationary points*) Suppose Assumption (A2) holds. Let $\{\widehat{\mathbf{x}}^{\mu_k}\}$ be a sequence of first-order stationary points of problem (4.12) with $\mu = \mu_k$. Then any accumulation point of $\{\widehat{\mathbf{x}}^{\mu_k}\}$ is a first-order d-stationary point of problem (1.1).

Proof Let $\widehat{\mathbf{x}}$ be an accumulation point of $\{\widehat{\mathbf{x}}^{\mu_k}\}$. Without loss of generality, we may assume that $\{\widehat{\mathbf{x}}^{\mu_k}\}$ converges to $\widehat{\mathbf{x}}$.

Since $\widehat{\mathbf{x}}^{\mu_k}$ is a first-order stationary point of problem (4.12) with $\mu = \mu_k$, then

$$\nabla \widetilde{f}_{\mu_k}(\widehat{\mathbf{x}}^{\mu_k}) = \nabla \mathcal{L}(\widehat{\mathbf{x}}^{\mu_k}) + \begin{pmatrix} \left[g' \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_1^{\mu_k}) - \widetilde{h}'_{\mu_k} \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_1^{\mu_k}) \right] \nabla \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_1^{\mu_k}) \\ \vdots \\ \left[g' \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_K^{\mu_k}) - \widetilde{h}'_{\mu_k} \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_K^{\mu_k}) \right] \nabla \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_K^{\mu_k}) \end{pmatrix} = \mathbf{0}.$$

Therefore, for any $\mathbf{z} \in \mathbb{R}^n$ we have

$$\begin{aligned} 0 &= \langle \nabla \widetilde{f}_{\mu_k}(\widehat{\mathbf{x}}^{\mu_k}), \mathbf{z} \rangle \\ &= \langle \nabla \mathcal{L}(\widehat{\mathbf{x}}^{\mu_k}), \mathbf{z} \rangle + \sum_{i=1}^K \left[g' \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) - \widetilde{h}'_{\mu_k} \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) \right] \langle \nabla \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}), \mathbf{z}_i \rangle. \end{aligned} \quad (4.14)$$

Let $k \rightarrow \infty$, then we get $\mu_k \rightarrow 0$ and $\widehat{\mathbf{x}}^{\mu_k} \rightarrow \widehat{\mathbf{x}}$, consequently, $\widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) \rightarrow m(\widehat{\mathbf{x}}_i)$, $g' \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) \rightarrow g' \circ m(\widehat{\mathbf{x}}_i)$ and $\widetilde{h}'_{\mu_k} \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) \rightarrow h' \circ m(\widehat{\mathbf{x}}_i)$. Moreover, from (4.4) and (4.5), we have

$$\lim_{k \rightarrow \infty} \langle \nabla \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}), \mathbf{z}_i \rangle = m'(\widehat{\mathbf{x}}_i; \mathbf{z}_i) \quad \text{if } \widehat{\mathbf{x}}_i \neq \mathbf{0},$$

and

$$\limsup_{k \rightarrow \infty} \langle \nabla \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}), \mathbf{z}_i \rangle = m'(\widehat{\mathbf{x}}_i; \mathbf{z}_i) \quad \text{if } \widehat{\mathbf{x}}_i = \mathbf{0}.$$

By the condition $\varphi'(0) := \varphi'(0+) = g'(0+) - h'(0+) > 0$, we know that $g' \circ m(\widehat{\mathbf{x}}_i) - h' \circ m(\widehat{\mathbf{x}}_i) > 0$ if $\widehat{\mathbf{x}}_i = \mathbf{0}$. Hence when k is sufficiently large, $g' \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) - \widetilde{h}'_{\mu_k} \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) > 0$ for the index i such that $\widehat{\mathbf{x}}_i = \mathbf{0}$. From (4.14), we derive that for any $\mathbf{z} \in \mathbb{R}^n$,

$$\begin{aligned} 0 &= \lim_{k \rightarrow \infty} \langle \nabla \widetilde{f}_{\mu_k}(\widehat{\mathbf{x}}^{\mu_k}), \mathbf{z} \rangle \\ &= \lim_{k \rightarrow \infty} \langle \nabla \mathcal{L}(\widehat{\mathbf{x}}^{\mu_k}), \mathbf{z} \rangle + \lim_{k \rightarrow \infty} \sum_{i=1}^K \left[g' \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) - \widetilde{h}'_{\mu_k} \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) \right] \langle \nabla \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}), \mathbf{z}_i \rangle \\ &= \langle \nabla \mathcal{L}(\widehat{\mathbf{x}}), \mathbf{z} \rangle + \lim_{k \rightarrow \infty} \sum_{i: \widehat{\mathbf{x}}_i \neq \mathbf{0}} \left[g' \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) - \widetilde{h}'_{\mu_k} \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) \right] \langle \nabla \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}), \mathbf{z}_i \rangle \\ &\quad + \lim_{k \rightarrow \infty} \sum_{i: \widehat{\mathbf{x}}_i = \mathbf{0}} \left[g' \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) - \widetilde{h}'_{\mu_k} \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) \right] \langle \nabla \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}), \mathbf{z}_i \rangle \\ &\leq \langle \nabla \mathcal{L}(\widehat{\mathbf{x}}), \mathbf{z} \rangle + \sum_{i: \widehat{\mathbf{x}}_i \neq \mathbf{0}} \lim_{k \rightarrow \infty} \left[g' \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) - \widetilde{h}'_{\mu_k} \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) \right] \cdot \lim_{k \rightarrow \infty} \langle \nabla \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}), \mathbf{z}_i \rangle \\ &\quad + \sum_{i: \widehat{\mathbf{x}}_i = \mathbf{0}} \lim_{k \rightarrow \infty} \left[g' \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) - \widetilde{h}'_{\mu_k} \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) \right] \cdot \limsup_{k \rightarrow \infty} \langle \nabla \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}), \mathbf{z}_i \rangle \end{aligned}$$

378

$$\begin{aligned}
&= \langle \nabla \mathcal{L}(\widehat{\mathbf{x}}), \mathbf{z} \rangle + \sum_{i: \widehat{\mathbf{x}}_i \neq \mathbf{0}} \left[g' \circ m(\widehat{\mathbf{x}}_i) - h' \circ m(\widehat{\mathbf{x}}_i) \right] m'(\widehat{\mathbf{x}}_i; \mathbf{z}_i) \\
&\quad + \sum_{i: \widehat{\mathbf{x}}_i = \mathbf{0}} \left[g' \circ m(\widehat{\mathbf{x}}_i) - h' \circ m(\widehat{\mathbf{x}}_i) \right] m'(\widehat{\mathbf{x}}_i; \mathbf{z}_i) \\
&= \langle \nabla \mathcal{L}(\widehat{\mathbf{x}}), \mathbf{z} \rangle + \sum_{i=1}^K \left[g' \circ m(\widehat{\mathbf{x}}_i) - h' \circ m(\widehat{\mathbf{x}}_i) \right] m'(\widehat{\mathbf{x}}_i; \mathbf{z}_i) \\
&= f'(\widehat{\mathbf{x}}; \mathbf{z}),
\end{aligned}$$

379 which shows that $\widehat{\mathbf{x}}$ is a first-order d-stationary point of problem (4.1). \square

380 Before discussing the consistency of second-order stationary points, we first study the
381 property of second-order stationary points of the smoothing problem (4.12).

382 **Lemma 4.2** Under Assumption (A2), let $\widehat{\mathbf{x}}^{\mu_k} \in \mathbb{R}^n$ be a second-order stationary point of
383 problem (4.12) with $\mu = \mu_k$, then the following two statements hold for $i = 1, \dots, K$:

$$\begin{aligned}
384 \quad (i) \quad & \|\nabla \mathcal{L}(\widehat{\mathbf{x}}^{\mu_k})\|_i = \left| g' \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) - \widetilde{h}'_{\mu_k} \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) \right| \frac{\|\widehat{\mathbf{x}}_i^{\mu_k}\|}{\sqrt{\|\widehat{\mathbf{x}}_i^{\mu_k}\|^2 + \mu_k}}. \\
385 \quad (ii) \quad & \widetilde{h}''_{\mu_k} \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) \frac{\|\widehat{\mathbf{x}}_i^{\mu_k}\|^4}{\|\widehat{\mathbf{x}}_i^{\mu_k}\|^2 + \mu_k} \leq \langle \nabla_i^2 \mathcal{L}(\widehat{\mathbf{x}}^{\mu_k}) \widehat{\mathbf{x}}_i^{\mu_k}, \widehat{\mathbf{x}}_i^{\mu_k} \rangle \\
386 \quad & \quad \quad \quad + \left[g' \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) - \widetilde{h}'_{\mu_k} \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) \right] \frac{\mu_k \|\widehat{\mathbf{x}}_i^{\mu_k}\|^2}{(\|\widehat{\mathbf{x}}_i^{\mu_k}\|^2 + \mu_k)^{\frac{3}{2}}}.
\end{aligned}$$

387 *Proof* (i) Since $\widehat{\mathbf{x}}^{\mu_k}$ is a second-order stationary point of problem (4.12), we have

$$\nabla \widetilde{f}_{\mu_k}(\widehat{\mathbf{x}}^{\mu_k}) = \nabla \mathcal{L}(\widehat{\mathbf{x}}^{\mu_k}) + \begin{pmatrix} \left[g' \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_1^{\mu_k}) - \widetilde{h}'_{\mu_k} \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_1^{\mu_k}) \right] \nabla \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_1^{\mu_k}) \\ \vdots \\ \left[g' \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_K^{\mu_k}) - \widetilde{h}'_{\mu_k} \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_K^{\mu_k}) \right] \nabla \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_K^{\mu_k}) \end{pmatrix} = \mathbf{0},$$

388 where, according to (4.2), $\nabla \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) = \frac{\widehat{\mathbf{x}}_i^{\mu_k}}{\sqrt{\|\widehat{\mathbf{x}}_i^{\mu_k}\|^2 + \mu_k}}$ for $i = 1, \dots, K$. Therefore, we get

$$\|\nabla \mathcal{L}(\widehat{\mathbf{x}}^{\mu_k})\|_i = \left| g' \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) - \widetilde{h}'_{\mu_k} \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) \right| \frac{\|\widehat{\mathbf{x}}_i^{\mu_k}\|}{\sqrt{\|\widehat{\mathbf{x}}_i^{\mu_k}\|^2 + \mu_k}}.$$

389 (ii) Since $\widehat{\mathbf{x}}^{\mu_k}$ is a second-order stationary point of problem (4.12), we know that $\nabla^2 \widetilde{f}_{\mu_k}(\widehat{\mathbf{x}}^{\mu_k})$
390 is positive semi-definite, and then $\langle \nabla^2 \widetilde{f}_{\mu_k}(\widehat{\mathbf{x}}^{\mu_k}) \mathbf{z}, \mathbf{z} \rangle \geq 0$ for any $\mathbf{z} \in \mathbb{R}^n$. For each fixed
391 $i = 1, \dots, K$, let $\bar{\mathbf{z}}_i = \widehat{\mathbf{x}}_i^{\mu_k}$ and other entries of $\bar{\mathbf{z}}$ are all zeros, then we get

$$\begin{aligned}
0 &\leq \langle \nabla^2 \widetilde{f}_{\mu_k}(\widehat{\mathbf{x}}^{\mu_k}) \bar{\mathbf{z}}, \bar{\mathbf{z}} \rangle \\
&= \langle \nabla_i^2 \mathcal{L}(\widehat{\mathbf{x}}^{\mu_k}) \widehat{\mathbf{x}}_i^{\mu_k}, \widehat{\mathbf{x}}_i^{\mu_k} \rangle + \left[g' \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) - \widetilde{h}'_{\mu_k} \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) \right] \cdot \langle \nabla^2 \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) \widehat{\mathbf{x}}_i^{\mu_k}, \widehat{\mathbf{x}}_i^{\mu_k} \rangle \\
&\quad - \widetilde{h}''_{\mu_k} \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) \left[\langle \nabla \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}), \widehat{\mathbf{x}}_i^{\mu_k} \rangle \right]^2,
\end{aligned} \tag{4.15}$$

where, according to (4.2) and $\widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) = \sqrt{\|\widehat{\mathbf{x}}_i^{\mu_k}\|^2 + \mu_k}$,

$$\left[\langle \nabla \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}), \widehat{\mathbf{x}}_i^{\mu_k} \rangle \right]^2 = \frac{\|\widehat{\mathbf{x}}_i^{\mu_k}\|^4}{\|\widehat{\mathbf{x}}_i^{\mu_k}\|^2 + \mu_k}, \quad \langle \nabla^2 \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) \widehat{\mathbf{x}}_i^{\mu_k}, \widehat{\mathbf{x}}_i^{\mu_k} \rangle = \frac{\mu_k \|\widehat{\mathbf{x}}_i^{\mu_k}\|^2}{(\|\widehat{\mathbf{x}}_i^{\mu_k}\|^2 + \mu_k)^{\frac{3}{2}}}.$$

392 Thus, from (4.15), we obtain

$$\begin{aligned} & \tilde{h}_{\mu_k}'' \circ \tilde{m}_{\mu_k}(\hat{\mathbf{x}}_i^{\mu_k}) \frac{\|\hat{\mathbf{x}}_i^{\mu_k}\|^4}{\|\hat{\mathbf{x}}_i^{\mu_k}\|^2 + \mu_k} \\ & \leq \langle \nabla_i^2 \mathcal{L}(\hat{\mathbf{x}}^{\mu_k}) \hat{\mathbf{x}}_i^{\mu_k}, \hat{\mathbf{x}}_i^{\mu_k} \rangle + \left[g' \circ \tilde{m}_{\mu_k}(\hat{\mathbf{x}}_i^{\mu_k}) - \tilde{h}_{\mu_k}' \circ \tilde{m}_{\mu_k}(\hat{\mathbf{x}}_i^{\mu_k}) \right] \frac{\mu_k \|\hat{\mathbf{x}}_i^{\mu_k}\|^2}{(\|\hat{\mathbf{x}}_i^{\mu_k}\|^2 + \mu_k)^{\frac{3}{2}}}. \end{aligned}$$

393 The proof is completed. \square

394 Now we begin to discuss the consistency of second-order stationary points. If h is twice
395 differentiable in $(0, \infty)$, such as logarithm penalty and fraction penalty, there is no need to
396 smooth h , but h^{MCP} and h^{SCAD} are not twice differentiable in $(0, \infty)$. In the following part,
397 we focus on that h is not twice differentiable in $(0, \infty)$.

398 **Assumption (A3)** Under Assumption (A2),

$$D(h) := \{t \in (0, \infty) : h \text{ is not twice differentiable at } t\} \quad (4.16)$$

399 has finite many points. In this case, we denote $l_h := \min\{t : t \in D(h)\}$, $L_h := \max\{t : t \in$
400 $D(h)\}$.

401 We can easily check that several penalty functions satisfy Assumption (A3), such as MCP
402 ($l_h = L_h = \alpha\lambda$) and SCAD ($l_h = \lambda, L_h = \alpha\lambda$). We also observe that the values of l_h and
403 L_h are highly consistent with the corresponding lower bounds obtained in Corollary 2.5 and
404 Theorem 3.9. Since g is affine in Assumption (A2), we know $\varphi = g - h$ is also not twice
405 differentiable at t for $t \in D(h)$.

406 **Lemma 4.3** Suppose Assumption (A3) holds and the following four conditions hold.

407 (a) There exists a nondecreasing function $C : \mathbb{R} \rightarrow \mathbb{R}_+$ such that $\|\nabla \mathcal{L}(\mathbf{x})\| \leq C(\mathcal{L}(\mathbf{x}))$
408 for any $\mathbf{x} \in \mathbb{R}^n$.

409 (b) There exists $\mathbf{x}^0 \in \mathbb{R}^n$ satisfying $\varphi'(0) > C(\mathcal{L}(\mathbf{x}^0))$.

410 (c) There exists $M > 0$ such that $\|\nabla^2 \mathcal{L}(\mathbf{x})\|_2 \leq M$ for all $\mathbf{x} \in \mathbb{R}^n$.

411 (d) If $l_h = L_h$, it holds that $\inf_{t \in (0, L_h]} \max\{H'(t; 1), -H'(t; -1)\} > M$; if $l_h < L_h$, it holds

412 that $\inf_{t \in (0, l_h]} \varphi'(t) \geq \varphi'(0)$ and that $\inf_{t \in (l_h, L_h]} \max\{H'(t; 1), -H'(t; -1)\} > M$, where $H(t) =$
413 $h'(t)$.

414 Let $\{\hat{\mathbf{x}}^{\mu_k}\}$ be a sequence of second-order stationary points of problem (4.12) with $\mu = \mu_k$
415 satisfying $\mathcal{L}(\hat{\mathbf{x}}^{\mu_k}) \leq \mathcal{L}(\mathbf{x}^0)$, and $\hat{\mathbf{x}}$ be any accumulation point of $\{\hat{\mathbf{x}}^{\mu_k}\}$, then the following
416 three statements hold:

417 (i) $\|\nabla \mathcal{L}(\hat{\mathbf{x}})\| < \varphi'(0)$.

418 (ii) $\min_{i: \hat{\mathbf{x}}_i \neq 0} \|\hat{\mathbf{x}}_i\| > L_h$.

419 (iii) For any subsequence $\{\hat{\mathbf{x}}^{\mu_k}\}_{k \in \mathcal{K}}$ converging to $\hat{\mathbf{x}}$, we have

$$\Gamma^{\mu_k} := \{i \in \{1, \dots, K\} : \|\hat{\mathbf{x}}_i^{\mu_k}\| \leq \frac{L_h}{2}\} = \{i \in \{1, \dots, K\} : \|\hat{\mathbf{x}}_i\| = 0\} := \Gamma$$

420 for all sufficiently large $k \in \mathcal{K}$,

421 *Proof* Without loss of generality, we may assume that $\{\hat{\mathbf{x}}^{\mu_k}\}$ converges to $\hat{\mathbf{x}}$.

422 (i) By Condition (a) and $\mathcal{L}(\hat{\mathbf{x}}^{\mu_k}) \leq \mathcal{L}(\mathbf{x}^0)$, we have

$$\|\nabla \mathcal{L}(\hat{\mathbf{x}}^{\mu_k})\| \leq C(\mathcal{L}(\hat{\mathbf{x}}^{\mu_k})) \leq C(\mathcal{L}(\mathbf{x}^0)).$$

423 Then it follows from the continuity of $\nabla \mathcal{L}(\cdot)$, $\hat{\mathbf{x}}^{\mu_k} \rightarrow \hat{\mathbf{x}}$ and Condition (b) that

$$\|\nabla \mathcal{L}(\hat{\mathbf{x}})\| \leq C(\mathcal{L}(\mathbf{x}^0)) < \varphi'(0).$$

The first conclusion is proved.

(ii) We consider an arbitrary nonzero group of $\widehat{\mathbf{x}}$, say $\widehat{\mathbf{x}}_i \neq \mathbf{0}$. Since $\mu_k \rightarrow 0$ and $\widehat{\mathbf{x}}_i^{\mu_k} \rightarrow \widehat{\mathbf{x}}_i$, it follows from (4.3) and (4.8) that

$$\begin{aligned} \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) &\rightarrow m(\widehat{\mathbf{x}}_i) = \|\widehat{\mathbf{x}}_i\| \neq 0, \quad \frac{\|\widehat{\mathbf{x}}_i^{\mu_k}\|}{\sqrt{\|\widehat{\mathbf{x}}_i^{\mu_k}\|^2 + \mu_k}} \rightarrow 1, \quad \text{and} \\ \left| g' \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) - \widetilde{h}'_{\mu_k} \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) \right| &\rightarrow |g' \circ m(\widehat{\mathbf{x}}_i) - h' \circ m(\widehat{\mathbf{x}}_i)| = \varphi'(\|\widehat{\mathbf{x}}_i\|) \geq 0. \end{aligned} \quad (4.17)$$

As a consequence of Lemma 4.2 (i), (4.17) and $\|[\nabla \mathcal{L}(\widehat{\mathbf{x}}^{\mu_k})]_i\| \rightarrow \|[\nabla \mathcal{L}(\widehat{\mathbf{x}})]_i\|$, we get

$$\|[\nabla \mathcal{L}(\widehat{\mathbf{x}})]_i\| = \varphi'(\|\widehat{\mathbf{x}}_i\|). \quad (4.18)$$

From Lemma 4.2 (ii) and Condition (c), we derive

$$\begin{aligned} &\widetilde{h}''_{\mu_k} \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) \frac{\|\widehat{\mathbf{x}}_i^{\mu_k}\|^4}{\|\widehat{\mathbf{x}}_i^{\mu_k}\|^2 + \mu_k} \\ &\leq \langle \nabla_i^2 \mathcal{L}(\widehat{\mathbf{x}}^{\mu_k}) \widehat{\mathbf{x}}_i^{\mu_k}, \widehat{\mathbf{x}}_i^{\mu_k} \rangle + \left[g' \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) - \widetilde{h}'_{\mu_k} \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) \right] \frac{\mu_k \|\widehat{\mathbf{x}}_i^{\mu_k}\|^2}{(\|\widehat{\mathbf{x}}_i^{\mu_k}\|^2 + \mu_k)^{\frac{3}{2}}} \\ &\leq M \|\widehat{\mathbf{x}}_i^{\mu_k}\|^2 + \left[g' \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) - \widetilde{h}'_{\mu_k} \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) \right] \frac{\mu_k \|\widehat{\mathbf{x}}_i^{\mu_k}\|^2}{(\|\widehat{\mathbf{x}}_i^{\mu_k}\|^2 + \mu_k)^{\frac{3}{2}}}. \end{aligned}$$

Since $\|\widehat{\mathbf{x}}_i^{\mu_k}\| \rightarrow \|\widehat{\mathbf{x}}_i\| > 0$, when k is sufficiently large the above inequality can be simplified as

$$\widetilde{h}''_{\mu_k} \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) \frac{\|\widehat{\mathbf{x}}_i^{\mu_k}\|^2}{\|\widehat{\mathbf{x}}_i^{\mu_k}\|^2 + \mu_k} \leq M + \frac{\left[g' \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) - \widetilde{h}'_{\mu_k} \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) \right] \mu_k}{(\|\widehat{\mathbf{x}}_i^{\mu_k}\|^2 + \mu_k)^{\frac{3}{2}}}.$$

Let $k \rightarrow 0$ in the above inequality. By (4.10) and (4.17), we obtain

$$\max\{H'(\|\widehat{\mathbf{x}}_i\|; 1), -H'(\|\widehat{\mathbf{x}}_i\|; -1)\} = \limsup_{k \rightarrow \infty} \widetilde{h}''_{\mu_k} \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) \leq M. \quad (4.19)$$

To verify the second conclusion, let us consider two cases.

Case 1: $l_h = L_h$. In this case, assume, on the contrary, that $\|\widehat{\mathbf{x}}_i\| \leq L_h$. Then by the first part of Condition (d), we obtain

$$\max\{H'(\|\widehat{\mathbf{x}}_i\|; 1), -H'(\|\widehat{\mathbf{x}}_i\|; -1)\} \geq \inf_{t \in (0, L_h]} \max\{H'(t; 1), -H'(t; -1)\} > M,$$

which is in contradiction with (4.19). Hence, we must have $\|\widehat{\mathbf{x}}_i\| > L_h$.

Case 2: $l_h < L_h$. In this case, assume at first that $\|\widehat{\mathbf{x}}_i\| \leq l_h$. Then by the second part of Condition (d), we obtain

$$\varphi'(\|\widehat{\mathbf{x}}_i\|) \geq \inf_{t \in (0, l_h]} \varphi'(t) \geq \varphi'(0).$$

But equality (4.18) and Conclusion (i) yield that

$$\varphi'(\|\widehat{\mathbf{x}}_i\|) = \|[\nabla \mathcal{L}(\widehat{\mathbf{x}})]_i\| \leq \|\nabla \mathcal{L}(\widehat{\mathbf{x}})\| < \varphi'(0),$$

which is a contradiction. Hence, we must have $\|\widehat{\mathbf{x}}_i\| > l_h$.

Secondly, assume that $l_h < \|\widehat{\mathbf{x}}_i\| \leq L_h$. Then by the second part of Condition (d), we obtain

$$\max\{H'(\|\widehat{\mathbf{x}}_i\|; 1), -H'(\|\widehat{\mathbf{x}}_i\|; -1)\} \geq \inf_{t \in (l_h, L_h]} \max\{H'(t; 1), -H'(t; -1)\} > M,$$

442 which is in contradiction with inequality (4.19). Hence, we must have $\|\widehat{\mathbf{x}}_i\| > L_h$.

443 Taken together, we have shown that $\|\widehat{\mathbf{x}}_i\| > L_h$ whenever $\widehat{\mathbf{x}}_i \neq \mathbf{0}$, which means $\min_{i:\widehat{\mathbf{x}}_i \neq \mathbf{0}} \|\widehat{\mathbf{x}}_i\| >$
444 L_h .

445 (iii) Let a subsequence $\{\widehat{\mathbf{x}}^{\mu_k}\}_{k \in \mathcal{K}} \rightarrow \widehat{\mathbf{x}}$, then $\{\widehat{\mathbf{x}}_i^{\mu_k}\}_{k \in \mathcal{K}} \rightarrow \widehat{\mathbf{x}}_i$ for each $i = 1, \dots, K$.
446 Suppose $i \in \Gamma$, then $\widehat{\mathbf{x}}_i = \mathbf{0}$. Since $\{\|\widehat{\mathbf{x}}_i^{\mu_k}\|\}_{k \in \mathcal{K}} \rightarrow \|\widehat{\mathbf{x}}_i\| = 0$, we have $\|\widehat{\mathbf{x}}_i^{\mu_k}\| < \frac{L_h}{2}$ for all
447 sufficiently large $k \in \mathcal{K}$. That is, $i \in \Gamma^{\mu_k}$ for any sufficiently large $k \in \mathcal{K}$, which shows
448 $\Gamma \subset \Gamma^{\mu_k}$. Now we suppose $i \in \Gamma^{\mu_k}$, then $\|\widehat{\mathbf{x}}_i^{\mu_k}\| \leq \frac{L_h}{2}$. If $i \notin \Gamma$, then $\widehat{\mathbf{x}}_i \neq \mathbf{0}$, therefore
449 $\|\widehat{\mathbf{x}}_i\| > L_h$ according to (ii). It follows from $\{\widehat{\mathbf{x}}_i^{\mu_k}\}_{k \in \mathcal{K}} \rightarrow \widehat{\mathbf{x}}_i$ that $\|\widehat{\mathbf{x}}_i^{\mu_k}\| > \frac{L_h}{2}$ for any
450 sufficiently large $k \in \mathcal{K}$. This contradiction shows $i \in \Gamma$, thus $\Gamma^{\mu_k} \subset \Gamma$ for all sufficiently
451 large $k \in \mathcal{K}$. Therefore, $\Gamma^{\mu_k} = \Gamma$ for all sufficiently large $k \in \mathcal{K}$. \square

452 **Remark 4.4** Condition (d) in Lemma 4.3 is very important to ensure the lower bound
453 given by Conclusion (ii) when h is differentiable but not twice differentiable in $(0, \infty)$. We
454 can see that MCP and SCAD meet this condition. In fact, for MCP, $l_h = L_h = \alpha\lambda$ ($\alpha >$
455 1), then $\inf_{t \in (0, L_h)} \max\{H'(t; 1), -H'(t; -1)\} = \frac{1}{\alpha} > M$ whenever α is taken such that $1 <$
456 $\alpha < \frac{1}{M}$; and for SCAD, $l_h = \lambda < \alpha\lambda = L_h$ ($\alpha > 2$), then $\inf_{t \in (0, l_h)} \varphi'(t) = \lambda = \varphi'(0)$ and
457 $\inf_{t \in (l_h, L_h)} \max\{H'(t; 1), -H'(t; -1)\} = \frac{1}{\alpha-1} > M$ whenever α is taken such that $2 < \alpha < \frac{1}{M} + 1$.

458 **Theorem 4.5** (Consistency of second-order stationary points) Under the conditions of Lem-
459 ma 4.3, let $\{\widehat{\mathbf{x}}^{\mu_k}\}$ be a sequence of second-order stationary points of problem (4.12) with
460 $\mu = \mu_k$ satisfying $\mathcal{L}(\mathbf{x}^{\mu_k}) \leq \mathcal{L}(\mathbf{x}^0)$, then any accumulation point of $\{\widehat{\mathbf{x}}^{\mu_k}\}$ is a second-order
461 d-stationary point of problem (1.1).

Proof Without loss of generality, we may assume that $\{\widehat{\mathbf{x}}^{\mu_k}\}$ converges to $\widehat{\mathbf{x}}$. Since $\widehat{\mathbf{x}}^{\mu_k}$ is a
second-order stationary point of problem (4.12) with $\mu = \mu_k$, we have

$$\nabla \widetilde{f}_{\mu_k}(\widehat{\mathbf{x}}^{\mu_k}) = \mathbf{0} \quad \text{and} \quad \langle \nabla^2 \widetilde{f}_{\mu_k}(\widehat{\mathbf{x}}^{\mu_k}) \mathbf{z}, \mathbf{z} \rangle \geq 0, \quad \forall \mathbf{z} \in \mathbb{R}^n.$$

462 According to Theorem 4.1, $\widehat{\mathbf{x}}$ is a first-order d-stationary point of problem (1.1), that is,
463 $f'(\widehat{\mathbf{x}}; \mathbf{z}) \geq 0$ for any $\mathbf{z} \in \mathbb{R}^n$.

464 In the following arguments, we only consider such $\mathbf{z} \in \mathbb{R}^n$ that makes $f'(\widehat{\mathbf{x}}; \mathbf{z}) = 0$.
465 According to Lemma 4.3 (i), it holds that $\max_{i:\widehat{\mathbf{x}}_i = \mathbf{0}} \|\nabla \mathcal{L}(\widehat{\mathbf{x}})\|_i \leq \|\nabla \mathcal{L}(\widehat{\mathbf{x}})\| < \varphi'(0)$. By virtue
466 of this inequality and Corollary 2.7 (iii), it yields from $f'(\widehat{\mathbf{x}}; \mathbf{z}) = 0$ that $\mathbf{z}_i = \mathbf{0}$ whenever
467 $\widehat{\mathbf{x}}_i = \mathbf{0}$.

468 By using $\mathbf{z}_i = \mathbf{0}$ whenever $\widehat{\mathbf{x}}_i = \mathbf{0}$, we have

$$\begin{aligned} 0 &\leq \langle \nabla^2 \widetilde{f}_{\mu_k}(\widehat{\mathbf{x}}^{\mu_k}) \mathbf{z}, \mathbf{z} \rangle \\ &= \langle \nabla^2 \mathcal{L}(\widehat{\mathbf{x}}^{\mu_k}) \mathbf{z}, \mathbf{z} \rangle \\ &\quad + \sum_{i=1}^K \left[g' \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) - \widetilde{h}'_{\mu_k} \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) \right] \cdot \langle \nabla^2 \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) \mathbf{z}_i, \mathbf{z}_i \rangle \\ &\quad - \sum_{i=1}^K \widetilde{h}''_{\mu_k} \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) \left[\langle \nabla \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}), \mathbf{z}_i \rangle \right]^2 \\ &= \langle \nabla^2 \mathcal{L}(\widehat{\mathbf{x}}^{\mu_k}) \mathbf{z}, \mathbf{z} \rangle + \sum_{i:\widehat{\mathbf{x}}_i \neq \mathbf{0}} \left[g' \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) - \widetilde{h}'_{\mu_k} \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) \right] \langle \nabla^2 \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) \mathbf{z}_i, \mathbf{z}_i \rangle \\ &\quad - \sum_{i:\widehat{\mathbf{x}}_i \neq \mathbf{0}} \widetilde{h}''_{\mu_k} \circ \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) \cdot \left[\langle \nabla \widetilde{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}), \mathbf{z}_i \rangle \right]^2. \end{aligned} \tag{4.20}$$

469 According to Lemma 4.3 (ii), we have $\min_{i:\widehat{\mathbf{x}}_i \neq \mathbf{0}} \|\widehat{\mathbf{x}}_i\| > L_h$. Under Assumption (A3), this in-
 470 equality means that h is twice continuously differentiable at each $\|\widehat{\mathbf{x}}_i\|$ whenever $\widehat{\mathbf{x}}_i \neq \mathbf{0}$.

471 Since for each i , $\lim_{k \rightarrow \infty} \widehat{\mathbf{x}}_i^{\mu_k} = \widehat{\mathbf{x}}_i$, $\lim_{k \rightarrow \infty} \widehat{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) = m(\widehat{\mathbf{x}}_i) = \|\widehat{\mathbf{x}}_i\|$,

$$\lim_{k \rightarrow \infty} \left[g' \circ \widehat{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) - \widetilde{h}'_{\mu_k} \circ \widehat{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) \right] = g' \circ m(\widehat{\mathbf{x}}_i) - h' \circ m(\widehat{\mathbf{x}}_i),$$

and for $\widehat{\mathbf{x}}_i \neq \mathbf{0}$,

$$\lim_{k \rightarrow \infty} \langle \nabla \widehat{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}), \mathbf{z}_i \rangle = m'(\widehat{\mathbf{x}}_i; \mathbf{z}_i),$$

$$\lim_{k \rightarrow \infty} \langle \nabla^2 \widehat{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) \mathbf{z}_i, \mathbf{z}_i \rangle = m^{(2)}(\widehat{\mathbf{x}}_i; \mathbf{z}_i),$$

$$\lim_{k \rightarrow \infty} \widetilde{h}''_{\mu_k} \circ \widehat{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) = H'(\|\widehat{\mathbf{x}}_i\|; 1) = -H'(\|\widehat{\mathbf{x}}_i\|; -1),$$

472 then from (4.20), we obtain

$$\begin{aligned} 0 &\leq \lim_{k \rightarrow \infty} \langle \nabla^2 \widetilde{f}_{\mu_k}(\widehat{\mathbf{x}}^{\mu_k}) \mathbf{z}, \mathbf{z} \rangle \\ &= \lim_{k \rightarrow \infty} \langle \nabla^2 \mathcal{L}(\widehat{\mathbf{x}}^{\mu_k}) \mathbf{z}, \mathbf{z} \rangle \\ &\quad + \sum_{i:\widehat{\mathbf{x}}_i \neq \mathbf{0}} \lim_{k \rightarrow \infty} \left[g' \circ \widehat{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) - \widetilde{h}'_{\mu_k} \circ \widehat{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) \right] \lim_{k \rightarrow \infty} \langle \nabla^2 \widehat{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) \mathbf{z}_i, \mathbf{z}_i \rangle \\ &\quad - \sum_{i:\widehat{\mathbf{x}}_i \neq \mathbf{0}} \lim_{k \rightarrow \infty} \widetilde{h}''_{\mu_k} \circ \widehat{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}) \cdot \left[\lim_{k \rightarrow \infty} \langle \nabla \widehat{m}_{\mu_k}(\widehat{\mathbf{x}}_i^{\mu_k}), \mathbf{z}_i \rangle \right]^2 \\ &= \langle \nabla^2 \mathcal{L}(\widehat{\mathbf{x}}) \mathbf{z}, \mathbf{z} \rangle + \sum_{i:\widehat{\mathbf{x}}_i \neq \mathbf{0}} \left[g' \circ m(\widehat{\mathbf{x}}_i) - h' \circ m(\widehat{\mathbf{x}}_i) \right] m^{(2)}(\widehat{\mathbf{x}}_i; \mathbf{z}_i) \\ &\quad - \sum_{i:\widehat{\mathbf{x}}_i \neq \mathbf{0}} H'(\widehat{\mathbf{x}}_i; 1) [m'(\widehat{\mathbf{x}}_i, \mathbf{z}_i)]^2 \\ &= \langle \nabla^2 \mathcal{L}(\widehat{\mathbf{x}}) \mathbf{z}, \mathbf{z} \rangle + \sum_{i:\widehat{\mathbf{x}}_i \neq \mathbf{0}} \left[g' \circ m(\widehat{\mathbf{x}}_i) - h' \circ m(\widehat{\mathbf{x}}_i) \right] m^{(2)}(\widehat{\mathbf{x}}_i; \mathbf{z}_i) \\ &\quad - \sum_{i:\widehat{\mathbf{x}}_i \neq \mathbf{0}} H'(\widehat{\mathbf{x}}_i; m'(\widehat{\mathbf{x}}_i, \mathbf{z}_i)) m'(\widehat{\mathbf{x}}_i, \mathbf{z}_i) \\ &= \langle \nabla^2 \mathcal{L}(\widehat{\mathbf{x}}) \mathbf{z}, \mathbf{z} \rangle + \sum_{i=1}^K \left[g' \circ m(\widehat{\mathbf{x}}_i) - h' \circ m(\widehat{\mathbf{x}}_i) \right] m^{(2)}(\widehat{\mathbf{x}}_i; \mathbf{z}_i) - \sum_{i=1}^K H'(\widehat{\mathbf{x}}_i; m'(\widehat{\mathbf{x}}_i, \mathbf{z}_i)) m'(\widehat{\mathbf{x}}_i, \mathbf{z}_i) \\ &= f^{(2)}(\widehat{\mathbf{x}}; \mathbf{z}), \end{aligned}$$

where the third equality is due to

$$H'(\widehat{\mathbf{x}}_i; 1) [m'(\widehat{\mathbf{x}}_i, \mathbf{z}_i)]^2 = -H'(\widehat{\mathbf{x}}_i; -1) [m'(\widehat{\mathbf{x}}_i, \mathbf{z}_i)]^2 = H'(\widehat{\mathbf{x}}_i; m'(\widehat{\mathbf{x}}_i, \mathbf{z}_i)) m'(\widehat{\mathbf{x}}_i, \mathbf{z}_i)$$

when $\|\widehat{\mathbf{x}}_i\| > L_h$, and the fourth equality is due to

$$\mathbf{z}_i = \mathbf{0}, \quad m'(\widehat{\mathbf{x}}_i, \mathbf{z}_i) = m^{(2)}(\widehat{\mathbf{x}}_i; \mathbf{z}_i) = 0$$

473 when $\widehat{\mathbf{x}}_i = \mathbf{0}$.

474 As a summary, we have shown that $\widehat{\mathbf{x}}$ is a first-order d-stationary point of problem (1.1)
 475 and that for any $\mathbf{z} \in \mathbb{R}^n$, $f'(\widehat{\mathbf{x}}; \mathbf{z}) = 0$ implies $f^{(2)}(\widehat{\mathbf{x}}; \mathbf{z}) \geq 0$. Therefore, $\widehat{\mathbf{x}}$ is a second-order
 476 d-stationary point of problem (1.1). \square

477 Now, we use an example of problem (1.1) to illustrate how to compute a second-order
478 directional stationary point by the smoothing method.

479 **Example 4.1.** Consider the following problem

$$\min_{x_1, x_2 \in \mathbb{R}} f(x_1, x_2) := \frac{1}{2}(x_1 + x_2 - 1)^2 + \varphi^{\text{MCP}}(|x_1|) + \varphi^{\text{MCP}}(|x_2|), \quad (4.21)$$

where the parameters in φ^{MCP} satisfy $\alpha > 1$ and $\lambda > 0$. In Tables 1,2,3, we present the sets of the first-order d-stationary points, second-order d-stationary points, local minimizers, and global minimizers of (4.21) with different parameters. From the tables, we can see the relation between these sets for problem (4.21):

first-order d-stationary \Leftarrow second-order d-stationary \Leftrightarrow local minimizer \Leftarrow global minimizer

480 For example, when $0 < \alpha\lambda \leq \frac{1}{2}$, $\lambda < 1$, let $\bar{x} := (\bar{x}_1, \bar{x}_2)^\top = (1 + \alpha\lambda, -\alpha\lambda)^\top$, then
481 $\bar{x}_1 + \bar{x}_2 = 1$, $|\bar{x}_1| \geq \alpha\lambda$, and $|\bar{x}_2| \geq \alpha\lambda$. It is easy to check that for any $d := (d_1, d_2)^\top \in \mathbb{R}^2$,
482 $f'(\bar{x}; d) = 0$ and

$$\begin{aligned} f^{(2)}(\bar{x}; d) &= (d_1, d_2) \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} (d_1, d_2)^\top + \begin{cases} -\frac{d_2^2}{\alpha}, & d_2 > 0, \\ 0, & d_2 \leq 0, \end{cases} \\ &= \begin{cases} (d_1 + d_2)^2 - \frac{d_2^2}{\alpha}, & d_2 > 0, \\ (d_1 + d_2)^2, & d_2 \leq 0. \end{cases} \end{aligned}$$

483 Since it cannot ensure $f^{(2)}(\bar{x}; d) \geq 0$ for any d , \bar{x} is a first-order d-stationary point but not
484 a second-order d-stationary point of (4.21).

Table 1 First-order d-stationary points of (4.21)

parameters	first-order d-stationary points
$0 < \alpha\lambda \leq \frac{1}{2}$ $\lambda < 1$	$(1, 0)^\top, (0, 1)^\top, \{(x_1, x_2)^\top : x_1 + x_2 = 1, x_1 \geq \alpha\lambda, x_2 \geq \alpha\lambda\}$
$\frac{1}{2} < \alpha\lambda \leq 1$ $\lambda < 1$	$(1, 0)^\top, (0, 1)^\top, (\frac{\alpha(1-\lambda)}{2\alpha-1}, \frac{\alpha(1-\lambda)}{2\alpha-1})^\top, \{(x_1, x_2)^\top : x_1 + x_2 = 1, x_1 \geq \alpha\lambda, x_2 \geq \alpha\lambda\}$
$\alpha\lambda > 1$ $\lambda < 1$	$(\frac{\alpha(1-\lambda)}{\alpha-1}, 0)^\top, (0, \frac{\alpha(1-\lambda)}{\alpha-1})^\top, (\frac{\alpha(1-\lambda)}{2\alpha-1}, \frac{\alpha(1-\lambda)}{2\alpha-1})^\top, \{(x_1, x_2)^\top : x_1 + x_2 = 1, x_1 \geq \alpha\lambda, x_2 \geq \alpha\lambda\}$
$\alpha\lambda > 1$ $\lambda \geq 1$	$(0, 0)^\top, \{(x_1, x_2)^\top : x_1 + x_2 = 1, x_1 \geq \alpha\lambda, x_2 \geq \alpha\lambda\}$

Table 2 Second-order d-stationary points / local minimizers of (4.21)

parameters	second-order d-stationary points / local minimizers
$0 < \alpha\lambda \leq \frac{1}{2}$ $\lambda < 1$	$(1, 0)^\top, (0, 1)^\top, \{(x_1, x_2)^\top : x_1 + x_2 = 1, x_1 > \alpha\lambda, x_2 > \alpha\lambda\}$
$\frac{1}{2} < \alpha\lambda \leq 1$ $\lambda < 1$	$(1, 0)^\top, (0, 1)^\top, \{(x_1, x_2)^\top : x_1 + x_2 = 1, x_1 > \alpha\lambda, x_2 > \alpha\lambda\}$
$\alpha\lambda > 1$ $\lambda < 1$	$(\frac{\alpha(1-\lambda)}{\alpha-1}, 0)^\top, (0, \frac{\alpha(1-\lambda)}{\alpha-1})^\top, \{(x_1, x_2)^\top : x_1 + x_2 = 1, x_1 > \alpha\lambda, x_2 > \alpha\lambda\}$
$\alpha\lambda > 1$ $\lambda \geq 1$	$(0, 0)^\top, \{(x_1, x_2)^\top : x_1 + x_2 = 1, x_1 > \alpha\lambda, x_2 > \alpha\lambda\}$

485 To test the smoothing method and the consistency theory of stationary points, we use the
486 smoothing trust region Newton (STRN) method proposed in [9] with an initial point $(1, 1)^\top$
487 to solve problem (4.21) where the smoothing function of h^{MCP} is taken $\tilde{h}_\mu^{\text{MCP}}$ as (4.11).
488 The numerical results are listed in Table 4, where f^* means the global minimum of (4.21),
489 and \bar{x} is the output solution of the STRN method. Table 4 shows that \bar{x} is a second-order
490 d-stationary point and a global minimizer of problem (4.21).

Table 3 Global minimizers of (4.21)

parameters	global minimizers
$0 < \alpha\lambda \leq \frac{1}{2}$ $\lambda < 1$	$(1, 0)^\top, (0, 1)^\top$
$\frac{1}{2} < \alpha\lambda \leq 1$ $\lambda < 1$	$(1, 0)^\top, (0, 1)^\top$
$\alpha\lambda > 1$ $\lambda < 1$	$(\frac{\alpha(1-\lambda)}{\alpha-1}, 0)^\top, (0, \frac{\alpha(1-\lambda)}{\alpha-1})^\top$
$\alpha\lambda > 1$ $\lambda \geq 1$	$(0, 0)^\top$

Table 4 Numerical results of the STRN method for (4.21) with different values of α and λ

α	λ	global minimizers	f^*	output solution \bar{x}	$f(\bar{x})$
$\alpha = 2$	$\lambda = 0.25$	$(1, 0)^\top, (0, 1)^\top$	0.0625	$(1, 0)^\top$	0.0625
$\alpha = 1.5$	$\lambda = 0.5$	$(1, 0)^\top, (0, 1)^\top$	0.1875	$(1, 0)^\top$	0.1875
$\alpha = 3$	$\lambda = 0.5$	$(0.75, 0)^\top, (0, 0.75)^\top$	0.3125	$(0.75, 0)^\top$	0.3125
$\alpha = 2$	$\lambda = 1$	$(0, 0)^\top$	0.5	$(0, 0)^\top$	0.5

491 5 Concluding remarks

492 This paper shows that the first-order and second-order d-stationary points of folded concave
 493 penalized group sparse optimization problem (1.1) are local minimizers fulfilling the first-
 494 order and second-order growth conditions respectively under some mild conditions. Moreover,
 495 we construct a twice continuously differentiable smoothing approximation for the nonsmooth
 496 objective function, and show that any accumulation point of the sequence of second-order
 497 stationary points of the smoothing problem is a second-order d-stationary point of problem
 498 (1.1). The result provides a theoretic basis for computing first-order and second-order d-
 499 stationary points of the problem by using the gradient and Hessian of smoothing functions.
 500 Our results can be used for developing second-order algorithms for folded concave penalized
 501 group sparse optimization problems, and verifying the optimality of numerical solutions
 502 obtained by any algorithms. A simple example shows the validity of our theory and numerical
 503 method.

504 **Acknowledgements.** The authors would like to thank two referees for their helpful
 505 comments.

506 References

- 507 1. Ahn, M., Pang, J.-S., Xin, J.: Difference-of-convex learning: directional stationarity, optimality, and sparsity. *SIAM*
 508 *J. Optim.* **27**, 1637-1665 (2017)
- 509 2. Ben-Tal, A., Zowe, J.: Necessary and sufficient optimality conditions for a class of nonsmooth minimization prob-
 510 lems. *Math. Program.* **24**, 70-91 (1982)
- 511 3. Bian, W., Chen, X.: Optimality and complexity for constrained optimization problems with nonconvex regulariza-
 512 tion. *Math. Oper. Res.* **42**, 1063-1084 (2017)
- 513 4. Bian, W., Chen, X., Ye, Y.: Complexity analysis of interior point algorithms for non-Lipschitz and nonconvex
 514 minimization. *Math. Program.* **149**, 301-327 (2015)
- 515 5. Breheny, P., Huang, J.: Group descent algorithms for nonconvex penalized linear and logistic regression models
 516 with grouped predictors. *Stat. Comput.* **25**, 173-187 (2015)
- 517 6. Bruckstein, A.M., Donoho, D.L., Elad, M.: From sparse solutions of systems of equations to sparse modeling of
 518 signals and images. *SIAM Rev.* **51**, 34-81 (2009)
- 519 7. Chang, T.H., Hong, M., Pang, J.-S.: Local minimizers and second-order conditions in composite piecewise program-
 520 ming via directional derivatives. arxiv.org/abs/1709.05758 (2017)
- 521 8. Chen, X.: Smoothing methods for nonsmooth, nonconvex minimization. *Math. Program.* **134**, 71-99 (2012)
- 522 9. Chen, X., Niu L., Yuan Y.: Optimality conditions and a smoothing trust region Newton method for non-Lipschitz
 523 optimization. *SIAM J. Optim.* **23**, 1528-1552 (2013)

- 524 10. Chen, X., Xu, F., Ye, Y.: Lower bound theory of nonzero entries in solutions of $\ell_2 - \ell_p$ minimization. *SIAM J. Sci.*
525 *Comput.* **32**, 2832-2852 (2010)
- 526 11. Eren Ahsen, M., Vidyasagar, M.: Error bounds for compressed sensing algorithms with group sparsity: A unified
527 approach. *Appl. Comput. Harmon. Anal.* **43**, 212-232 (2017)
- 528 12. Fan, J., Li, R.: Variable selection via nonconcave penalized likelihood and its oracle properties, *J. Amer. Stat.*
529 *Assoc.* **96**, 1348-1360 (2001)
- 530 13. Fan, J., Li, R.: Statistical challenges with high dimensionality: feature selection in knowledge discovery. *Proceedings*
531 *of the International Congress of Mathematicians* **3**, 595-622, Madrid, Spain (2006)
- 532 14. Fan, J., Xue, L., Zou, H.: Strong oracle optimality of folded concave penalized estimation. *Ann. Stat.* **42**, 819-849
533 (2014)
- 534 15. Hu, Y., Li, C., Meng, K., Qin, J., Yang, X.: Group sparse optimization via $\ell_{p,q}$ regularization. *J. Mach. Learn. Res.*
535 **18**, 1-52 (2017)
- 536 16. Huang, J., Ma, S., Xue, H., Zhang, C.H.: A group bridge approach for variable selection. *Biometrika* **96**, 339-355
537 (2009)
- 538 17. Huang, J., Zhang, T.: The benefit of group sparsity. *Ann. Stat.* **38**, 1978-2004 (2010)
- 539 18. Jacob, L., Obozinski, G., Vert, J.P.: Group lasso with overlap and graph lasso. *Proceedings of the 26th International*
540 *Conference on Machine Learning*, 433-440, Montreal, Canada (2009)
- 541 19. Jiao, Y., Jin, B., Lu, X.: Group sparse recovery via the $\ell_0(\ell_2)$ penalty: theory and algorithm. *IEEE Trans. Signal*
542 *Process.* **65**, 998-1012 (2017)
- 543 20. Knight, K., Fu, W. J.: Asymptotics for lasso-type estimators. *Ann. Stat.* **28**, 1356-1378 (2000)
- 544 21. Le Thi, H.A., Pham Dinh, T., Le, H.M., Vo, X.T.: DC approximation approaches for sparse optimization. *Eur. J.*
545 *Oper. Res.* **244**, 26-46 (2015)
- 546 22. Lee, S., Oh, M., Kim, Y.: Sparse optimization for nonconvex group penalized estimation. *J. Stat. Comput. Simul.*
547 **86**, 597-610 (2016)
- 548 23. Liu, H., Yao, T., Li, R., Ye, Y.: Folded concave penalized sparse linear regression: sparsity, statistical performance,
549 and algorithmic theory for local solutions. *Math. Program.* **166**, 207-240 (2017)
- 550 24. Meier, L., van de Geer, S., Bühlmann, P.: The group Lasso for logistic regression. *J. R. Stat. Soc. B* **70**, 53-71
551 (2008)
- 552 25. Nikolova, M., Ng, M.K., Zhang, S., Ching, W.: Efficient reconstruction of piecewise constant images using nonsmooth
553 nonconvex minimization. *SIAM J. Imag. Sci.* **1**, 2-25 (2008)
- 554 26. Pang, J.-S., Razaviyayn, M., Alvarado, A.: Computing B-stationary points of nonsmooth DC programs. *Math.*
555 *Oper. Res.* **42**, 95-118 (2017)
- 556 27. Rochafellar, R.T., Wets, R.J.-B.: *Variational Analysis* (3rd). Springer-Verlag, Berlin, 2009
- 557 28. Tibshirani, R.: Regression shrinkage and selection via the LASSO. *J. R. Stat. Soc. B* **58**, 267-288 (1996)
- 558 29. Wei, F., Zhu, H.: Group coordinate descent algorithms for nonconvex penalized regression. *Comput. Stat. Data*
559 *Anal.* **56**, 316-326 (2012)
- 560 30. Yang, E., Lozano, A.C.: Sparse + group-sparse dirty models: statistical guarantees without unreasonable conditions
561 and a case for non-convexity. *Proceedings of the 34th International Conference on Machine Learning*, PMLR **70**,
562 3911-3920, Sydney, Australia (2017)
- 563 31. Yang, X.Q.: On second-order directional derivatives. *Nonlin. Anal. TMA* **26**, 55-66, 1996.
- 564 32. Yang, Y., Zou, H.: A fast unified algorithm for solving group-lasso penalize learning problems. *Stat. Comput.* **25**,
565 1129-1141 (2015)
- 566 33. Y. Yuan, Conditions for convergence of trust region algorithms for nonsmooth optimization. *Math. Program.* **31**,
567 220-228 (1985)
- 568 34. Yuan, M., Lin, Y.: Model selection and estimation in regression with grouped variables. *J. R. Stat. Soc. B* **68**, 49-67
569 (2006)
- 570 35. Zhang, C.H.: Nearly unbiased variable selection under minimax concave penalty. *Ann. Stat.* **38**, 894-942 (2010)
- 571 36. Zhang, T.: Analysis of multi-stage convex relaxation for sparse regularization. *J. Mach. Learn. Res.* **11**, 1081-1107
572 (2010)
- 573 37. Zhao, P., Rocha, G., Yu, B.: The composite absolute penalties family for grouped and hierarchical variable selection.
574 *Ann. Stat.* **37**, 3468-3497 (2009)