

1 **A SMOOTHING PROXIMAL GRADIENT ALGORITHM FOR NONSMOOTH**
2 **CONVEX REGRESSION WITH CARDINALITY PENALTY***

3 WEI BIAN[†] AND XIAOJUN CHEN[‡]

4 **Abstract.** In this paper, we focus on the constrained sparse regression problem, where the loss function is convex
5 but nonsmooth, and the penalty term is defined by the cardinality function. Firstly, we give an exact continuous
6 relaxation problem in the sense that both problems have the same optimal solution set. Moreover, we show that
7 a vector is a local minimizer with the lower bound property of the original problem if and only if it is a lifted
8 stationary point of the relaxation problem. Secondly, we propose a *smoothing proximal gradient* (SPG) algorithm
9 for finding a lifted stationary point of the continuous relaxation model. Our algorithm is a novel combination of
10 the classical proximal gradient algorithm and the smoothing method. We prove that the proposed SPG algorithm
11 globally converges to a lifted stationary point of the relaxation problem, has the local convergence rate of $o(k^{-\tau})$
12 with $\tau \in (0, \frac{1}{2})$ on the objective function value, and identifies the zero entries of the lifted stationary point in finite
13 iterations. Finally, we use three examples to illustrate the validity of the continuous relaxation model and good
14 numerical performance of the SPG algorithm.

15 **Key words.** nonsmooth convex regression; cardinality penalty; proximal gradient method; smoothing method;
16 global sequence convergence.

17 **AMS subject classifications.** 90C46, 49K35, 90C30, 65K05

18 **1. Introduction.** For a vector $x \in \mathbb{R}^n$, denote its support set by $\mathcal{A}(x) = \{i \in \{1, \dots, n\} : x_i \neq 0\}$,
19 its cardinality by $|\mathcal{A}(x)|$, and its ℓ_0 -norm by $\|x\|_0 = |\mathcal{A}(x)|$. We call $x \in \mathbb{R}^n$ is sparse if
20 $|\mathcal{A}(x)| \ll n$. Sparse optimization problems emerge in many scientific and engineering problems, such
21 as regression [52], imaging decomposition [51], visual coding [44], source separation [10], compressed
22 sensing [12, 22], variable selection [39], etc. Sparse optimization is also the core problem of high-
23 dimensional statistical learning [11, 24]. These problems aim to find the sparse solutions of a
24 system of linear or nonlinear equations. The optimization model with the ℓ_0 -norm penalty can
25 improve estimation accuracy by effectively identifying the important predictors, and also enhance
26 its interpretability. However, it is known that the ℓ_0 penalized optimization problems are NP-hard.

27 Under some conditions on the sensing matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ (such as the RIP and incoherence
28 conditions), Donoho [22], and Candès, Romberg, Tao [12] proved that solving the ℓ_1 minimization
29 can find a sparsest solution satisfying the system of linear equations $\mathbf{A}x = b$ with $b \in \mathbb{R}^m$. However,
30 in 2001, Fan and Li [23] pointed out that using the ℓ_1 penalty often results in a biased estimator,
31 and introduced a *smoothly clipped absolute deviation* (SCAD) penalty. Besides SCAD, there are
32 many variant of continuous nonconvex penalties, such as the hard thresholding penalty [56], log-sum
33 penalty [13], bridge ℓ_p ($0 < p < 1$) penalty [17, 25], capped- ℓ_1 penalty [45, 47, 55] and minimax
34 concave penalty (MCP) [54]. These continuous but nonconvex penalties would bring better sparse
35 solutions than the ℓ_1 penalty in many cases [6, 15, 28, 31]. The estimators obtained by the SCAD,
36 MCP and capped- ℓ_1 penalty functions satisfy the three important properties: unbiasedness, con-
37 tinuity in data and sparsity [23]. Meantime, there are many algorithms for solving these continuous

* **Funding:** This work was funded in part by the NSF foundation (11871178,61773136) of China and Hong Kong Research Grant Council grant (153000/17P).

[†]School of Mathematics, Harbin Institute of Technology, Harbin, China; Institute of Advanced Study in Mathematics, Harbin Institute of Technology, Harbin 150001, China (bianweilvse520@163.com).

[‡]Department of Applied Mathematics, The Hong Kong Polytechnic University, Hong Kong, China (maxjchen@polyu.edu.hk).

38 nonconvex optimization problems, such as the iterative reweighted algorithm [13, 43, 36], interior
 39 point method [7], trust region method [18], cubic method [14], DC (difference of convex) function
 40 algorithm [1, 37], iterative thresholding algorithm [8], primal dual active set method [27], etc.

41 Despite the existing literature on the nonconvex but continuous penalties for replacing the ℓ_0 -
 42 norm, some important questions still remain. First of all, the relationships between the cardinality
 43 penalty problem and its continuous relaxations are not very clear for most cases regarding the mini-
 44 mizers. Apart from the theoretical results for the convex ℓ_1 relaxation under restrictive hypotheses,
 45 only a few special cases have been analyzed for the consistency. With a suitable condition on the
 46 sensing matrix \mathbf{A} , the equivalence between ℓ_0 and ℓ_p ($0 < p \leq 1$) problems with constraint $\mathbf{A}x = b$
 47 was proved in [25] and then this result was extended to the problem with equality and inequality
 48 constraints in [26]. In [19], the authors gave a class of smooth nonconvex penalties to approximate
 49 the ℓ_0 penalty in terms of the consistency of global minimizers. In the DC programming frame-
 50 work, an approximation of the ℓ_0 penalty with the consistency of global minimizers was studied in
 51 [37]. Recently, Soubies, Blanc-Féraud and Aubert proposed a continuous exact ℓ_0 (CEL0) penalty
 52 for the ℓ_2 - ℓ_0 problem [51], where the global minimizers of both problems can be the same, and in
 53 [50], they verified that the capped- ℓ_1 and SCAD penalties could only guarantee the consistency of
 54 global minimizers to the ℓ_2 - ℓ_0 problem, while the MCP, truncated- ℓ_p with $0 < p < 1$ and CEL0
 55 penalties could not only own the consistence of global minimizers, but also ensure that its local
 56 minimizers are in the set of local minimizers of the ℓ_2 - ℓ_0 problem. Next, due to the nonconvexity
 57 of the penalties, finding global minimizers of these nonconvex problems is often NP-hard. Most
 58 existing work for these continuous nonconvex penalized problems focuses on the stationary points
 59 in different sense [1, 7, 8, 14, 18, 31, 35, 36, 46]. Moreover, due to their nonconvexity, only the
 60 subsequence convergence to a stationary point can be proved for the proposed algorithms. The K-L
 61 (Kurdyka-Łojasiewicz) condition is a popular tool to obtain the algorithmic sequence convergence.
 62 In [2], the sequence convergence to a critical point of a class of nonconvex semi-algebra problems is
 63 established, where the K-L condition plays the key role. Most recently, the authors in [46] stated
 64 that it would be interesting whether the sequence convergence can be established to the DC problem
 65 by a given algorithm without the K-L condition on the objective function.

66 Denote x^* the true estimator, which is the true solution of the considered (linear or nonlinear)
 67 regression problem. Then, the oracle estimator is defined by

$$68 \quad (1.1) \quad x^{\text{oracle}} \in \arg \min_{x_{\mathcal{A}(x^*)^c} = 0} f(x),$$

69 where $\mathcal{A}(x^*)^c$ means the complementary set of $\mathcal{A}(x^*)$ and $f : \mathbb{R}^n \rightarrow [0, \infty)$ is the loss function to
 70 evaluate the regression. The oracle estimator can be used as a theoretic benchmark for comparison
 71 of computed solutions. We call that the penalized model has the oracle property if it owns a local
 72 solution having the same asymptotic distribution as the oracle estimator. The penalized problem
 73 with the SCAD, MCP or capped- ℓ_1 penalty owns the oracle property simultaneously [23, 54, 55].
 74 A folded concave penalized problem often has multiple local solutions and the oracle property is
 75 established only for one of local solutions [24]. Hence, deriving some appealing properties, such
 76 as the optimality, sparsity or statistical properties, of the relevant stationary points is interesting.
 77 Ahn, Pang and Xin [1] established some optimality and sparsity properties of the d-stationary
 78 points (its definition will be reminded in Section 2) of the continuous relaxation problems. Fan,
 79 Xue and Zou [24] proved that as long as there is a reasonable initial estimator, an oracle estimator
 80 can be obtained via the one-step local linear approximation algorithm.

81 In the recent years, algorithmic research on the sparse regression problems with cardinality
 82 penalty has received much attention [4, 3, 29, 31, 32]. However, to the best of our knowledge, all

83 the existing results are built up for the problem with a continuously differentiable loss function.
 84 The primal dual active set methods are proposed in [29, 31, 32] for the ℓ_2 - ℓ_0 problems. Under
 85 some regularity conditions, such as the strict complementarity condition [31] or RIP condition
 86 on the sensing matrix [29, 32], some variants of the primal dual active set methods were proved
 87 to be convergent in finite iterations. The loss functions considered in [4, 3, 40] are continuously
 88 differentiable and with Lipschitz continuous gradients.

89 **Our focuses and contributions.** In this paper, we consider the following penalized sparse
 90 regression problem with cardinality penalty, that is,

$$91 \quad (1.2) \quad \min_{x \in \mathcal{X}} \mathcal{F}_{\ell_0}(x) := f(x) + \lambda \|x\|_0,$$

92 where $\mathcal{X} = \{x \in \mathbb{R}^n : l \leq x \leq u\}$, $f : \mathbb{R}^n \rightarrow [0, \infty)$ is convex (not necessarily smooth), λ is a
 93 positive parameter, and $l, u \in \{\mathbb{R}, \pm\infty\}^n$ with $l \leq 0 \leq u$ and $l < u$.

94 One application of problem (1.2) comes from the linear regression problem. It is well-known
 95 that the least squares estimate with the ℓ_2 - ℓ_0 model is not robust for many cases [23]. We need to
 96 consider the problem with the outlier-resistant loss function, such as the ℓ_1 loss function given by

$$97 \quad (1.3) \quad f(x) = \frac{1}{m} \|Ax - b\|_1,$$

98 or Huber's functions [30], which are convex, but not smooth. Another important application of
 99 problem (1.2) comes from the censored regression problem with the nonsmooth convex loss function

$$100 \quad (1.4) \quad f(x) = \frac{1}{pm} \sum_{i=1}^m |\max\{A_i x - c_i, 0\} - b_i|^p,$$

101 where $p \in [1, 2]$, $A_i^T \in \mathbb{R}^n$ and $c_i, b_i \in \mathbb{R}$, $i = 1, \dots, m$. There are some other nonsmooth convex
 102 loss functions, for example the negative log-quasi-likelihood function [23] or the check loss function
 103 in penalized quantile regression [24, 33]. To the best of our knowledge, only little work has been
 104 dedicated to the penalized sparse regression problem (1.2) with a general convex loss function.

105 For a given parameter $\nu > 0$, let $\Phi(x) = \sum_{i=1}^n \phi(x_i)$ be a continuous relaxation of the ℓ_0 penalty
 106 with the capped- ℓ_1 function ϕ given by

$$107 \quad (1.5) \quad \phi(t) = \min\{1, |t|/\nu\}.$$

108 We consider the following Lipschitz continuous optimization problem for solving (1.2):

$$109 \quad (1.6) \quad \min_{x \in \mathcal{X}} \mathcal{F}(x) := f(x) + \lambda \Phi(x).$$

110 Differently from the previous work [1, 4, 3, 7, 8, 14, 18, 29, 31, 32, 35, 36, 46], this paper considers
 111 the original cardinality penalty problem with a continuous convex loss function and uses an exact
 112 continuous relaxation problem to solve it. In particular, we focus on problem (1.2) with a continuous
 113 convex loss function, which is nonsmooth or whose gradient is not Lipschitz continuous. The main
 114 contributions of this paper include the following two aspects. First, we prove that the continuous
 115 relaxation problem (1.6) with certain $\nu > 0$ has two advantages: global minimizers of (1.2) and
 116 (1.6) are same; any lifted stationary point of (1.6) (its definition will be reminded in Section 2) is
 117 a local minimizer of (1.2) with a desired lower bound property. Second, we propose a smoothing
 118 proximal gradient (SPG) algorithm with global sequence convergence to a lifted stationary point of

119 (1.6) without using the K-L condition. Moreover, the SPG algorithm owns a local convergence rate
 120 on the objective function value of (1.6) and the finite iterative identification for the zero entries of
 121 the limit point.

122 **Notations.** We denote $\mathbb{N} = \{0, 1, \dots\}$ and $\mathbb{D}^n = \{d \in \mathbb{R}^n : d_i \in \{1, 2, 3\}, i = 1, \dots, n\}$. For
 123 $x \in \mathbb{R}^n$ and $\delta > 0$, let $\|x\| := \|x\|_2$ and $\mathbb{B}_\delta(x)$ means the open ball centered at x with radius δ . For
 124 a nonempty, closed and convex set $\mathcal{X} \subseteq \mathbb{R}^n$, $N_{\mathcal{X}}(x)$ means the normal cone to \mathcal{X} at $x \in \mathcal{X}$. Let
 125 $\mathbf{1}_n \in \mathbb{R}^n$ be the all-ones vector and $e_i \in \mathbb{R}^n$ be the i th column of the n dimensional identity matrix.
 126 For a locally Lipschitz continuous function $\psi : \mathbb{R}^n \rightarrow \mathbb{R}$, we denote $\partial\psi(x)$ the Clarke subgradient
 127 [20] of ψ at $x \in \mathbb{R}^n$.

128 **2. An exact continuous relaxation for (1.2).** In this section, we analyze the relationships
 129 between (1.2) and (1.6), where the capped- ℓ_1 penalty can let problem (1.6) own the oracle property
 130 and then can be seen as one of the best continuous relaxations to the ℓ_0 -norm penalty [45].

131 ASSUMPTION 1. f is Lipschitz continuous on \mathcal{X} with Lipschitz constant L_f .

132 ASSUMPTION 2. Positive parameter ν in (1.5) satisfies $\nu < \bar{\nu} := \lambda/L_f$.

If there is no special explanation, we suppose Assumption 1 and Assumption 2 hold throughout the
 paper, and assume that L_f is large enough such that $L_f \geq \frac{\lambda}{\Gamma}$, where

$$\Gamma := \min\{|l_i|, u_j : l_i \neq 0, u_j \neq 0, i = 1, \dots, n, j = 1, \dots, n\}.$$

133 When f is defined by the ℓ_1 loss function or the loss function in (1.4) with $p = 1$, we can let
 134 $L_f = \max\{\|A\|_\infty, \frac{\lambda}{\Gamma}\}$.

2.1. Lifted stationary points of (1.6). Though ϕ is piecewise linear, problem (1.6) is still
 a nonconvex optimization problem. It has been proved in [6] that finding a global minimizer of
 (1.6) is NP-hard in general. Note that ϕ in (1.5) can be reformulated as a DC function, i.e.

$$\phi(t) = \frac{1}{\nu}|t| - \max\{\theta_1(t), \theta_2(t), \theta_3(t)\}$$

135 with $\theta_1(t) = 0$, $\theta_2(t) = t/\nu - 1$ and $\theta_3(t) = -t/\nu - 1$. For $t \in \mathbb{R}$, denote

$$136 \quad (2.1) \quad \mathcal{D}(t) = \{i \in \{1, 2, 3\} : \theta_i(t) = \max\{\theta_1(t), \theta_2(t), \theta_3(t)\}\}.$$

137 DEFINITION 2.1. [46] We say that $x \in \mathcal{X}$ is a lifted stationary point of (1.6) if there exist
 138 $d_i \in \mathcal{D}(x_i)$ for $i = 1, \dots, n$, such that

$$139 \quad (2.2) \quad \lambda \sum_{i=1}^n \theta'_{d_i}(x_i) e_i \in \partial f(x) + \frac{\lambda}{\nu} \partial \left(\sum_{i=1}^n |x_i| \right) + N_{\mathcal{X}}(x).$$

140 If (2.2) holds for all $d_i \in \mathcal{D}(x_i)$, $\forall i = 1, \dots, n$, then we call x a d-stationary point [46]. Due to the
 141 piecewise linearity of $\max\{\theta_1(t), \theta_2(t), \theta_3(t)\}$, x is a d-stationary point of (1.6) if and only if it is a
 142 local minimizer. Recall that \bar{x} is a limiting stationary point [48] of (1.6), if

$$143 \quad (2.3) \quad 0 \in \bar{\partial}(f + \lambda\Phi)(\bar{x}) + N_{\mathcal{X}}(\bar{x}),$$

where “ $\bar{\partial}$ ” indicates the limiting subgradient. And \bar{x} is a Clarke stationary point of (1.6), $0 \in$
 $\partial(f + \lambda\Phi)(\bar{x}) + N_{\mathcal{X}}(\bar{x})$. We call $\bar{x} \in \mathcal{X}$ a critical point of (1.6) if it satisfies $0 \in \partial f(\bar{x}) + \lambda\partial\Phi(\bar{x}) + N_{\mathcal{X}}(\bar{x})$.
 It holds that

$$\mathcal{S}_d \subseteq \mathcal{S}_{lim} \subseteq \mathcal{S}_{lif} \subseteq \mathcal{S}_{cl} \subseteq \mathcal{S}_{cr},$$

144 but their inverse may not hold, where \mathcal{S}_d , \mathcal{S}_{lim} , \mathcal{S}_{lif} , \mathcal{S}_{cl} and \mathcal{S}_{cr} denote the d-stationary point set,
 145 limiting stationary point set, lifted stationary point set, Clarke stationary point set and critical
 146 point set of (1.6), respectively.

147 A natural question arises why we focus on the lifted stationary points rather than the others.
 148 First, the lifted stationary points satisfy a sharper optimal necessary condition than the Clarke
 149 and critical stationary points. Second, the d-stationary and limiting stationary points of (1.6) are
 150 difficult to be computed. Though Pang, Razaviyayn and Alvarado [46] developed a novel algorithm
 151 for computing a d-stationary point of the DC optimization problems, the algorithm in [46] cannot
 152 be directly used to solve problem (1.6).

153 **2.2. Characterizations of lifted stationary points of (1.6).** With the computable con-
 154 dition on ν defined in Assumption 2, we first verify that the element in $\Pi_{i=1}^n \mathcal{D}(x_i)$ for a lifted
 155 stationary point satisfying (2.2) is unique and well-defined.

156 **PROPOSITION 2.2.** *If \bar{x} is a lifted stationary point of (1.6), then the vector $d^{\bar{x}} = (d_1^{\bar{x}}, \dots, d_n^{\bar{x}})^T \in$
 157 $\prod_{i=1}^n \mathcal{D}(\bar{x}_i)$ satisfying (2.2) is unique. In particular, for $i = 1, \dots, n$,*

$$158 \quad (2.4) \quad d_i^{\bar{x}} = \begin{cases} 1 & \text{if } |\bar{x}_i| < \nu, \\ 2 & \text{if } \bar{x}_i \geq \nu, \\ 3 & \text{if } \bar{x}_i \leq -\nu. \end{cases}$$

159 *Proof.* If $|\bar{x}_i| \neq \nu$, then the statement in this proposition holds naturally. Hence, we only
 160 need to consider the case $|\bar{x}_i| = \nu$. When $\bar{x}_i = \nu$, since $\mathcal{D}(\bar{x}_i) = \{1, 2\}$, arguing by contradiction,
 161 we assume (2.2) holds with $d_i^{\bar{x}} = 1$. By $\nu < \bar{\nu}$, we have $\bar{x}_i \in (l_i, u_i)$, and by (2.2), there exists
 162 $\xi(\bar{x}) \in \partial f(\bar{x})$ such that $0 = \xi_i(\bar{x}) + \lambda/\nu$, which implies $\lambda/\nu = |\xi_i(\bar{x})| \leq L_f$. This leads to a
 163 contradiction to $\nu < \lambda/L_f$. Then, (2.4) holds for $\bar{x}_i = \nu$. Similar analysis can be given for the case
 164 that $\bar{x}_i = -\nu$, which completes the proof. \square

165 For a given $d = (d_1, \dots, d_n)^T \in \mathbb{D}^n$, we define

$$166 \quad (2.5) \quad \Phi^d(x) := \sum_{i=1}^n |x_i|/\nu - \sum_{i=1}^n \theta_{d_i}(x_i),$$

167 which is convex with respect to x . It can be verified that $\Phi(x) = \min_{d \in \mathbb{D}^n} \Phi^d(x)$, $\forall x \in \mathcal{X}$. In
 168 particular, for a fixed $\bar{x} \in \mathcal{X}$, $\Phi(\bar{x}) = \Phi^{d^{\bar{x}}}(\bar{x})$ with $d^{\bar{x}}$ defined in (2.4).

169 **Remark 2.1.** *Proposition 2.2 implies that \bar{x} is a local minimizer of (1.6) if and only if \bar{x} is
 170 a lifted stationary point of (1.6) and $|\bar{x}_i| \neq \nu$, $\forall i = 1, \dots, n$. Moreover, due to the convexity of
 171 $f(x) + \lambda\Phi^d(x)$ and the linearity of $\sum_{i=1}^n \theta_{d_i}(x_i)$ for a fixed $d \in \mathbb{D}^n$, the assertion in Proposition 2.2
 172 implies the following equivalent results:*

$$173 \quad \bar{x} \text{ is a lifted stationary point of (1.6)} \Leftrightarrow (2.2) \text{ holds at } \bar{x} \in \mathcal{X} \text{ with } d = d^{\bar{x}} \text{ defined in (2.4)}$$

$$174 \quad (2.6) \quad \Leftrightarrow \bar{x} \in \arg \min_{x \in \mathcal{X}} f(x) + \lambda\Phi^{d^{\bar{x}}}(x)$$

$$175 \quad (2.7) \quad \Leftrightarrow \bar{x} \in \arg \min_{x \in \mathcal{X}, d^x = d^{\bar{x}}} f(x) + \lambda\Phi(x),$$

176 where the last equivalence uses $\Phi^{d^{\bar{x}}}(\bar{x}) = \Phi(\bar{x})$ and $\Phi^{d^{\bar{x}}}(x) \geq \Phi(x)$, $\forall x \in \mathbb{R}^n$.

177 We then show a lower bound property of the lifted stationary points of (1.6).

178 LEMMA 2.3. *If $\bar{x} \in \mathcal{X}$ is a lifted stationary point of (1.6), then it holds that*

$$179 \quad (2.8) \quad \bar{x}_i \in (-\nu, \nu) \quad \Rightarrow \quad \bar{x}_i = 0, \quad \forall i = 1, \dots, n.$$

180 *Proof.* Suppose \bar{x} is a lifted stationary point of (1.6). Assume that $\bar{x}_i \in (-\nu, \nu) \setminus \{0\}$ for some
 181 $i \in \{1, \dots, n\}$. Then, $d_i^{\bar{x}} = 1$ and $\bar{x}_i \in (l_i, u_i)$. By Definition 2.1, there exists $\xi(\bar{x}) \in \partial f(\bar{x})$ such
 182 that $\xi_i(\bar{x}) + (\lambda/\nu)\text{sign}(\bar{x}_i) = 0$. Then, $\lambda/\nu = |\xi_i(\bar{x})| \leq \|\xi(\bar{x})\| \leq L_f$, which leads to a contradiction
 183 to $\nu < \lambda/L_f$. Thus, for any $i \in \{1, \dots, n\}$, $\bar{x}_i \in (-\nu, \nu)$ implies $\bar{x}_i = 0$. \square

184 **Remark 2.2.** *On the one hand, if f is not continuously differentiable on $\mathcal{X}_\nu = \{x \in \mathcal{X} : |x_i| = \nu \text{ for some } i \in \{1, \dots, n\}\}$, a lifted stationary point of (1.6) is not necessary to be a Clarke
 185 stationary point [46]. On the other hand, if f is continuously differentiable on \mathcal{X}_ν , then \bar{x} is a
 186 lifted stationary point of (1.6) if and only if it is a limiting stationary point, but is not necessary to
 187 be a Clarke stationary point. A counterexample can be provided by setting $f(x) = (x_1 + x_2 - 1)^2$,
 188 $l = (0, 0)^T$, $u = (1, 1)^T$, $\lambda = 1$ and $\nu = 0.2$ in (1.6), where $\nu < \bar{\nu} = 0.25$. It follows from Lemma
 189 2.3 that $\mathcal{S}_{\text{cl}} = \mathcal{S}_{\text{lif}} \cup \{(0, 0.2)^T, (0.2, 0)^T\}$, where $\mathcal{S}_{\text{lif}} = \{x \in \mathbb{R}^2 : x_1 + x_2 = 1, x_1 \geq 0.2, x_2 \geq$
 190 $0.2\} \cup \{(0, 0)^T, (1, 0)^T, (0, 1)^T\}$.*

192 **2.3. Links between (1.2) and (1.6).** The goal of this subsection is to study the links
 193 between the ℓ_0 penalized minimization problem (1.2) and its continuous relaxation (1.6). In light
 194 of the lower bound characterization of the lifted stationary points of (1.6) given in Lemma 2.3, we
 195 show the links between (1.2) and (1.6) by the two following results, where the first result focuses
 196 on global minimizers, and the second is on local minimizers.

197 **THEOREM 2.4.** *$\bar{x} \in \mathcal{X}$ is a global minimizer of (1.2) if and only if it is a global minimizer of*
 198 *(1.6). Moreover, problems (1.2) and (1.6) have the same optimal value.*

Proof. First, let $\bar{x} \in \mathcal{X}$ be a global minimizer of (1.6), then \bar{x} is a lifted stationary point of
 (1.6). By (2.8), it gives $\Phi(\bar{x}) = \|\bar{x}\|_0$. Then,

$$f(\bar{x}) + \lambda\|\bar{x}\|_0 = f(\bar{x}) + \lambda\Phi(\bar{x}) \leq f(x) + \lambda\Phi(x) \leq f(x) + \lambda\|x\|_0, \quad \forall x \in \mathcal{X},$$

199 where the last inequality uses $\Phi(x) \leq \|x\|_0, \forall x \in \mathbb{R}^n$. Thus, \bar{x} is a global minimizer of (1.2).

Next, suppose $\bar{x} \in \mathcal{X}$ is a global minimizer of (1.2) but not a global minimizer of (1.6). Then
 there exists a global minimizer of (1.6) denoted by \hat{x} such that

$$f(\hat{x}) + \lambda\Phi(\hat{x}) < f(\bar{x}) + \lambda\Phi(\bar{x}).$$

200 From $\Phi(\hat{x}) = \|\hat{x}\|_0$ and $\Phi(\bar{x}) \leq \|\bar{x}\|_0$, we get $f(\hat{x}) + \lambda\|\hat{x}\|_0 < f(\bar{x}) + \lambda\|\bar{x}\|_0$, which leads to a
 201 contradiction. Thus, any global minimizer of (1.2) must be a global minimizer of (1.6). Hence,
 202 using Lemma 2.3, we ensure that problems (1.2) and (1.6) have the same optimal value. \square

203 Theorem 2.4 provides that problems (1.2) and (1.6) have the same global solution set. The
 204 following proposition and the subsequent example show that this is not always true for their local
 205 minimizers.

206 **PROPOSITION 2.5.** *If \bar{x} is a lifted stationary point of (1.6), then it is a local minimizer of (1.2)*
 207 *and the objective functions have the same value at \bar{x} , i.e. $\mathcal{F}_{\ell_0}(\bar{x}) = \mathcal{F}(\bar{x})$.*

Proof. Coming back to the definition of $\Phi^{d^{\bar{x}}}$ defined in (2.5) and from the lower bound property
 of \bar{x} in (2.8), for any $x \in \mathbb{R}^n$, we have

$$\Phi^{d^{\bar{x}}}(x) = \sum_{i=1}^n |x_i|/\nu - \sum_{i=1}^n \theta_{d_i^{\bar{x}}}(x_i) = \sum_{i:|\bar{x}_i| \geq \nu} 1 + \sum_{i:|\bar{x}_i| < \nu} |x_i|/\nu = \|\bar{x}\|_0 + \sum_{i:\bar{x}_i=0} |x_i|/\nu.$$

Then, there exists $\varrho > 0$ such that $\Phi^{d^{\bar{x}}}(x) \leq \|x\|_0, \forall x \in \mathbb{B}_\varrho(\bar{x})$. Combining this with $\Phi(x) \leq \|x\|_0$ and (2.6) gives

$$f(\bar{x}) + \lambda\|\bar{x}\|_0 \leq f(x) + \lambda\|x\|_0, \quad \forall x \in \mathcal{X} \cap \mathbb{B}_\varrho(\bar{x}).$$

208 Thus, \bar{x} is a local minimizer of (1.2). □

209 Proposition 2.5 states that any lifted stationary point of (1.6) is a local minimizer of (1.2),
 210 which implies that any local minimizer of (1.6) is certainly a local minimizer of (1.2). Due to the
 211 special structure of the cardinality norm, any minimizer of $\min_{x \in \mathcal{X}} f(x)$ is a local minimizer of
 212 (1.2). The following example shows that a lifted stationary point of (1.6) is a local minimizer of
 213 (1.2) with the lower bound property in (2.8) and is likely a global minimizer.

214 **Example 2.1.** Let problem (1.2) be in the form of

$$215 \quad (2.9) \quad \min_{0 \leq x_1, x_2 \leq 1} \mathcal{F}_{\ell_0}(x_1, x_2) := |x_1 + x_2 - 1| + \lambda\|x\|_0.$$

216 We can easily find that $\mathcal{LM} := \{x \in \mathbb{R}^2 : x_1 + x_2 = 1, 0 \leq x_1, x_2 \leq 1\} \cup \{(0, 0)^T\}$ is the set of local
 217 minimizers of (2.9). Moreover, $(0, 0)^T$ is the unique global minimizer when $\lambda > 1$, the global minimi-
 218 zers are $\{(0, 1)^T, (1, 0)^T\}$ when $\lambda < 1$, and the global minimizers are $\{(0, 1)^T, (1, 0)^T, (0, 0)^T\}$ when
 219 $\lambda = 1$. Here, $\bar{\nu}$ in Lemma 2.3 can be $\min\{\sqrt{2}\lambda/2, 1\}$. With $\nu < \min\{\sqrt{2}\lambda/2, 1\}$, the lifted stationary
 220 points of (1.6) for this example are $\{x \in \mathbb{R}^2 : x_1 + x_2 = 1, \nu \leq x_1, x_2 \leq 1\} \cup \{(0, 0)^T, (1, 0)^T, (0, 1)^T\}$,
 221 which is a proper subset of \mathcal{LM} . Specially, if $\sqrt{2}/2 < \lambda \leq 1$ and $1/2 < \nu < \min\{\sqrt{2}\lambda/2, 1\}$, the
 222 lifted stationary points of (1.6) are $\{(1, 0)^T, (0, 1)^T, (0, 0)^T\}$.

223 When f is convex, \bar{x} is a local minimizer of (1.2) if and only if $\bar{x} \in \mathcal{X}$ satisfies

$$224 \quad (2.10) \quad 0 \in [\partial f(\bar{x}) + N_{\mathcal{X}}(\bar{x})]_i, \quad \forall i \in \mathcal{A}(\bar{x}),$$

225 which is a criterion for the local minimizers of (1.2) [40]. From Lemma 2.3 and Theorem 2.4, we
 226 find that the lower bound property in (2.8) holds for any global minimizer of (1.2), but is not true
 227 for all of its local minimizers. This inspires us to define a class of strong local minimizers of (1.2)
 228 by combining the optimality condition in (2.10) and the lower bound property in (2.8).

DEFINITION 2.6. We call $\bar{x} \in \mathcal{X}$ a ν -strong local minimizer of (1.2), if there exist $\bar{\xi} \in \partial f(\bar{x})$
 and $\bar{\eta} \in N_{\mathcal{X}}(\bar{x})$ such that for any $i \in \mathcal{A}(\bar{x})$, it holds

$$\bar{\xi}_i + \bar{\eta}_i = 0 \quad \text{and} \quad |\bar{x}_i| \geq \nu.$$

229 By (2.10), any ν -strong local minimizer of (1.2) is a local minimizer of it. To close this section,
 230 we give a result on the relationship between the ν -strong local minimizers of (1.2) and the lifted
 231 stationary points of (1.6).

232 PROPOSITION 2.7. $\bar{x} \in \mathcal{X}$ is a ν -strong local minimizer of (1.2) if and only if it is a lifted
 233 stationary point of (1.6). Moreover, if $\bar{x} \in \mathcal{X}$ is a ν -strong local minimizer of (1.2), then it holds

$$234 \quad \mathcal{F}_{\ell_0}(\bar{x}) \leq \mathcal{F}_{\ell_0}(x), \quad \forall x \in \mathcal{X} \cap (\bar{x} - \nu e, \bar{x} + \nu e),$$

$$235 \quad (2.11) \quad f(\bar{x}) \leq f(x), \quad \forall x \in \{x \in \mathcal{X} : \mathcal{A}(x) \subseteq \mathcal{A}(\bar{x})\},$$

$$236 \quad (2.12) \quad \bar{x} \text{ is an oracle solution defined in (1.1) if } \mathcal{A}(\bar{x}) = \mathcal{A}(x^*).$$

Proof. From Lemma 2.3, we can easily verify the first statement. By (2.6), we see that if \bar{x} is
 a lifted stationary point of (1.6), then

$$\mathcal{F}_{\ell_0}(\bar{x}) = f(\bar{x}) + \lambda\|\bar{x}\|_0 = f(\bar{x}) + \lambda\Phi(\bar{x}) = f(\bar{x}) + \lambda\Phi^{d^{\bar{x}}}(\bar{x}) \leq f(x) + \lambda\Phi^{d^{\bar{x}}}(x), \quad \forall x \in \mathcal{X}.$$

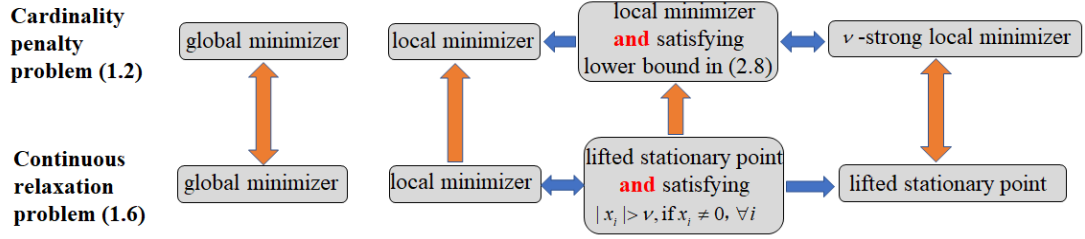


Fig. 2.1: Links between problems (1.2) and (1.6)

237 Due to Lemma 2.3, we then have $\mathcal{F}_{\ell_0}(\bar{x}) \leq \mathcal{F}_{\ell_0}(x)$, $\forall x \in \mathcal{X} \cap (\bar{x} - \nu \mathbf{1}_n, \bar{x} + \nu \mathbf{1}_n)$, which holds
 238 from $\Phi^{q^x}(x) \leq \|x\|_0$, $\forall x \in (\bar{x} - \nu \mathbf{1}_n, \bar{x} + \nu \mathbf{1}_n)$. Recalling (2.6) again, we obtain $f(\bar{x}) \leq f(x) +$
 239 $\lambda \sum_{i \notin \mathcal{A}(\bar{x})} |x_i|/\nu$, $\forall x \in \mathcal{X}$. If $\mathcal{A}(x) \subseteq \mathcal{A}(\bar{x})$, then $x_i = 0$ for $i \notin \mathcal{A}(\bar{x})$. Hence, (2.11) holds, which
 240 immediately implies (2.12). \square

241 **Remark 2.3.** In [50], the authors gave a unified view of exact continuous penalties for ℓ_2 - ℓ_0
 242 minimization, which derives necessary and sufficient conditions on ℓ_0 continuous relaxations such
 243 that each (local and global) minimizer of the underlying relaxation is also a minimizer of the ℓ_2 - ℓ_0
 244 problem. However, the property that any local minimizer of the relaxation problem with the capped-
 245 ℓ_1 penalty is a local minimizer of the ℓ_2 - ℓ_0 problem cannot be verified by the results in [50]. In this
 246 paper, we prove this property for the capped- ℓ_1 penalty by its lifted stationary points.

247 To end this section, we use Fig. 2.1 to give a brief description on the links between problems
 248 (1.2) and (1.6) when $\nu < \bar{\nu}$.

249 **3. Numerical Algorithm and its convergence analysis.** In this section, we focus on the
 250 numerical algorithm for finding a lifted stationary point of (1.6), which is a ν -strong local minimizer
 251 of (1.2). The first two subsections briefly introduce some useful preliminary results on smoothing
 252 methods and the proximal gradient algorithm, the third subsection presents a new proximal gradient
 253 algorithm combined with the smoothing method, and the last two subsections show the convergence
 254 of the proposed algorithm for solving (1.6).

255 **3.1. Smoothing approximation method.** A well-known method for solving nonsmooth
 256 optimization problems is to approximate the original problem by a sequence of smooth problems,
 257 which own rich theory and powerful numerical algorithms [42]. For the sake of completeness, we
 258 formally define a class of smoothing functions for f in (1.6).

259 **DEFINITION 3.1.** We call $\tilde{f} : \mathbb{R}^n \times [0, \bar{\mu}] \rightarrow \mathbb{R}$ with $\bar{\mu} > 0$ a smoothing function of the convex
 260 function f in (1.6), if $\tilde{f}(\cdot, \mu)$ is continuously differentiable in \mathbb{R}^n for any fixed $\mu > 0$ and satisfies
 261 the following conditions:

- 262 (i) $\lim_{z \rightarrow x, \mu \downarrow 0} \tilde{f}(z, \mu) = f(x)$, $\forall x \in \mathcal{X}$;
 263 (ii) (convexity) $\tilde{f}(x, \mu)$ is convex with respect to x in \mathcal{X} for any fixed $\mu > 0$;
 264 (iii) (gradient consistency) $\{\lim_{z \rightarrow x, \mu \downarrow 0} \nabla_z \tilde{f}(z, \mu)\} \subseteq \partial f(x)$, $\forall x \in \mathcal{X}$;
 (iv) (Lipschitz continuity with respect to μ) there exists a positive constant κ such that

$$|\tilde{f}(x, \mu_2) - \tilde{f}(x, \mu_1)| \leq \kappa |\mu_1 - \mu_2|, \quad \forall x \in \mathcal{X}, \mu_1, \mu_2 \in [0, \bar{\mu}];$$

- 265 (v) (Lipschitz continuity with respect to x) there exists a constant $L > 0$ such that for any

266 $\mu \in (0, \bar{\mu}]$, $\nabla_x \tilde{f}(\cdot, \mu)$ is Lipschitz continuous on \mathcal{X} with Lipschitz constant $L\mu^{-1}$.

267 Throughout this paper, we denote \tilde{f} a smoothing function of f in (1.6). When it is clear from
 268 the context, the derivative of $\tilde{f}(x, \mu)$ with respect to x is simply denoted as $\nabla \tilde{f}(x, \mu)$. Definition
 269 3.1-(iv) implies

$$270 \quad (3.1) \quad |\tilde{f}(x, \mu) - f(x)| \leq \kappa\mu, \quad \forall x \in \mathcal{X}, 0 < \mu \leq \bar{\mu}.$$

271
 272 **Example 3.1.** Many existing results in [16, 34, 49] give us some theoretical basis for con-
 273 structing smoothing functions satisfying the conditions in Definition 3.1. A smoothing function of
 274 the ℓ_1 loss function in (1.3) can be defined by

$$275 \quad (3.2) \quad \tilde{f}(x, \mu) = \frac{1}{m} \sum_{i=1}^m \tilde{\theta}(A_i x - b_i, \mu) \quad \text{with} \quad \tilde{\theta}(s, \mu) = \begin{cases} |s| & \text{if } |s| > \mu, \\ \frac{s^2}{2\mu} + \frac{\mu}{2} & \text{if } |s| \leq \mu. \end{cases}$$

276 For the loss function in (1.4) with $p = 1$, a smoothing function of it can be defined by

$$277 \quad (3.3) \quad \tilde{f}(x, \mu) = \frac{1}{m} \sum_{i=1}^m \tilde{\theta}(\tilde{\phi}(A_i x, \mu) - b_i, \mu) \quad \text{with} \quad \tilde{\phi}(s, \mu) = \begin{cases} \max\{s, 0\} & \text{if } |s| > \mu, \\ \frac{(s + \mu)^2}{4\mu} & \text{if } |s| \leq \mu. \end{cases}$$

We end this subsection by giving the following notations:

$$\tilde{\mathcal{F}}^d(x, \mu) \triangleq \tilde{f}(x, \mu) + \lambda \Phi^d(x) \quad \text{and} \quad \tilde{\mathcal{F}}(x, \mu) \triangleq \tilde{f}(x, \mu) + \lambda \Phi(x),$$

278 where \tilde{f} is a smoothing function of f , $\mu > 0$ and $d \in \mathbb{D}^n$. For any fixed $\mu > 0$ and $d \in \mathbb{D}^n$, both
 279 $\tilde{\mathcal{F}}^d(x, \mu)$ and $\tilde{\mathcal{F}}(x, \mu)$ are nonsmooth, $\tilde{\mathcal{F}}^d(x, \mu)$ is convex, but $\tilde{\mathcal{F}}(x, \mu)$ is nonconvex. Moreover,

$$280 \quad (3.4) \quad \tilde{\mathcal{F}}^d(x, \mu) \geq \tilde{\mathcal{F}}(x, \mu), \quad \forall d \in \mathbb{D}^n, x \in \mathcal{X}, \mu \in (0, \bar{\mu}].$$

281 **3.2. Proximal gradient method.** In this subsection, we consider the following constrained
 282 convex optimization problem with given smoothing parameter $\mu > 0$ and vector $d \in \mathbb{D}^n$

$$283 \quad (3.5) \quad \min_{x \in \mathcal{X}} \tilde{\mathcal{F}}^d(x, \mu).$$

284 It is good news that, for any given vectors $d \in \mathbb{D}^n$, $w \in \mathbb{R}^n$ and a positive number $\tau > 0$, the
 285 proximal operator of $\tau \Phi^d$ on \mathcal{X} has a closed form solution, i.e.

$$286 \quad (3.6) \quad \hat{x} = \arg \min_{x \in \mathcal{X}} \left\{ \tau \Phi^d(x) + \frac{1}{2} \|x - w\|^2 \right\}$$

287 can be calculated by $\hat{x}_i = \min\{\max\{l_i, y_i\}, u_i\}$ for $i = 1, \dots, n$, where

$$288 \quad (3.7) \quad y_i = \begin{cases} 0 & \text{if } |\bar{w}_i| \leq \tau/\nu, \\ \bar{w}_i - \tau/\nu & \text{if } \bar{w}_i > \tau/\nu, \\ \bar{w}_i + \tau/\nu & \text{if } \bar{w}_i < -\tau/\nu, \end{cases}$$

289 with $\bar{w}_i = w_i$ for $d_i = 1$, $\bar{w}_i = w_i + \tau/\nu$ for $d_i = 2$ and $\bar{w}_i = w_i - \tau/\nu$ for $d_i = 3$. Toward this end,
 290 we consider an approximation of $\tilde{\mathcal{F}}^d(\cdot, \mu)$ around a given point z , given by

$$291 \quad (3.8) \quad Q_{d,\gamma}(x, z, \mu) = \tilde{f}(z, \mu) + \langle x - z, \nabla \tilde{f}(z, \mu) \rangle + \frac{1}{2} \gamma \mu^{-1} \|x - z\|^2 + \lambda \Phi^d(x)$$

292 with a constant $\gamma > 0$. Since $\Phi^d(x)$ is convex with respect to x for any fixed $d \in \mathbb{D}^n$, function
 293 $Q_{d,\gamma}(x, z, \mu)$ is a strongly convex function with respect to x for any fixed d, γ, z and μ . Then,
 294 minimization problem $\min_{x \in \mathcal{X}} Q_{d,\gamma}(x, z, \mu)$ admits a unique minimizer, denoted by \hat{x} , which can
 295 be calculated by (3.7) with $\tau = \lambda \gamma^{-1} \mu$ and $w = z - \gamma^{-1} \mu \nabla \tilde{f}(z, \mu)$.

296 **3.3. Smoothing proximal gradient (SPG) algorithm.** In this subsection, we propose a
 297 new algorithm for finding a lifted stationary point of (1.6). Since the proposed algorithm combines
 298 the smoothing method and the proximal gradient algorithm, we call it *Smoothing Proximal Gradient*
 299 (SPG) algorithm.

300 For convenience of further reading, we begin this subsection by emphasizing the following
 301 assumptions needed in the convergence analysis of the SPG algorithm.

- 302 • (A1) Assumption 1 and Assumption 2 hold.
- 303 • (A2) \tilde{f} is a smoothing function of f defined in Definition 3.1.
- 304 • (A3) \mathcal{F} in (1.6) (or \mathcal{F}_{ℓ_0} in (1.2)) is level bounded on \mathcal{X}^1 .

305 As the feasible region \mathcal{X} is bounded, then assumption (A3) holds naturally. We give some more
 306 details on the parameters in these assumptions. Parameter L_f in Assumption 1 is used to define
 307 ν such that problems (1.2) and (1.6) have the consistency in Theorem 2.4 and Proposition 2.5.
 308 Parameter κ in Definition 3.1 is used in the SPG algorithm, which can be calculated exactly for
 309 most smoothing functions [16] and $\kappa = \frac{1}{2}$ for the smoothing functions in (3.2) and (3.3). The
 310 value of L in Definition 3.1 is not necessary and we will use a simple line search method to find an
 311 acceptable value at each iteration of the SPG algorithm. Upon the above assumptions, we present
 312 the SPG algorithm for solving (1.6). See Algorithm 3.1.

At each iteration, this algorithm takes the proximal gradient algorithm for solving (3.5) with
 fixed μ_k, γ_k and d^k , and uses a simple criterion for updating μ_k . The values of γ_k are chosen
 independently in Step 1 of each iteration. Step 3 updates the smoothing parameter μ_k by using
 (3.12), where $\tilde{\mathcal{F}}(x^{k+1}, \mu_k) + \kappa \mu_k$ can be seen as an energy function and its monotone non-increasing
 property will be proved in Lemma 3.3. If the energy function is decreased more than the given
 scale at the current iteration, then the current smoothing parameter is still acceptable for the next
 iteration, otherwise we reduce its value by the updating rule in (3.13) for the next iteration. Let

$$\mathcal{N}^s = \{k \in \mathbb{N} : \mu_{k+1} \neq \mu_k\},$$

313 and denote n_r^s the r th smallest number in \mathcal{N}^s . Then, we can obtain following updating method of
 314 $\{\mu_k\}$

$$315 \quad (3.14) \quad \mu^k = \mu^{n_r^s+1} = \frac{\mu_0}{(n_r^s + 1)^\sigma}, \quad \forall n_r^s + 1 \leq k \leq n_{r+1}^s,$$

316 which will be used in the proof of Lemma 3.2 and Lemma 3.5.

¹We call function \mathcal{F} is level bounded on \mathcal{X} , if for any $\Gamma > 0$, the level set $\{x \in \mathcal{X} : \mathcal{F}(x) \leq \Gamma\}$ is bounded.

Algorithm 3.1 Smoothing Proximal Gradient (SPG) algorithm

Input: Take initial iterates $x^{-1} = x^0 \in \mathcal{X}$ and $\mu_{-1} = \mu_0 \in (0, \bar{\mu}]$. Choose constants $\rho > 1$, $\sigma \in (\frac{1}{2}, 1)$, $\alpha > 0$ and $0 < \underline{\gamma} \leq \bar{\gamma}$. Set $k = 0$.

While a termination criterion is not met, **do**

Step 1. Choose $\gamma_k \in [\underline{\gamma}, \bar{\gamma}]$ and let $d^k \triangleq d^{x^k}$, where d^{x^k} is defined in (2.4).

Step 2. 2a) Compute

$$(3.9) \quad \hat{x}^{k+1} = \arg \min_{x \in \mathcal{X}} Q_{d^k, \gamma_k}(x, x^k, \mu_k).$$

2b) If \hat{x}^{k+1} satisfies

$$(3.10) \quad \tilde{\mathcal{F}}^{d^k}(\hat{x}^{k+1}, \mu_k) \leq Q_{d^k, \gamma_k}(\hat{x}^{k+1}, x^k, \mu_k),$$

set

$$(3.11) \quad x^{k+1} = \hat{x}^{k+1}$$

and go to **Step 3**. Otherwise, let $\gamma_k = \rho\gamma_k$ and return to 2a).

Step 3: If

$$(3.12) \quad \tilde{\mathcal{F}}(x^{k+1}, \mu_k) + \kappa\mu_k - \tilde{\mathcal{F}}(x^k, \mu_{k-1}) - \kappa\mu_{k-1} \leq -\alpha\mu_k^2,$$

set $\mu_{k+1} = \mu_k$, otherwise, set

$$(3.13) \quad \mu_{k+1} = \frac{\mu_0}{(k+1)^\sigma}.$$

Increment k by one and return to **Step 1**.

end while

317 **3.4. Basic convergence analysis of the SPG algorithm.** Denote $\{x^k\}$, $\{\gamma_k\}$ and $\{\mu_k\}$
 318 be the sequences generated by the SPG algorithm. In this subsection, we first establish some basic
 319 properties of the iterates $\{x^k\}$, $\{\gamma_k\}$ and $\{\mu_k\}$ in Lemma 3.2. Then, by the level boundedness
 320 assumption of \mathcal{F} (or \mathcal{F}_{ℓ_0}) on \mathcal{X} , the boundedness of $\{x^k\}$ is obtained in Lemma 3.3. At last, the
 321 subsequential convergence of $\{x^k : k \in \mathcal{N}^s\}$ to a lifted stationary point of (1.6) is established in
 322 Proposition 3.4.

323 **LEMMA 3.2.** *The proposed SPG algorithm is well-defined, and the sequences $\{x^k\}$, $\{\gamma_k\}$ and*
 324 *$\{\mu_k\}$ generated by it own the following properties:*

- 325 (i) $\{x^k\} \subseteq \mathcal{X}$ and $\{\gamma_k\} \subseteq [\underline{\gamma}, \max\{\bar{\gamma}, \rho L\}]$;
 326 (ii) *there are infinite elements in \mathcal{N}^s and $\lim_{k \rightarrow \infty} \mu_k = 0$.*

Proof. (i). Upon rearranging terms, (3.10) can be rewritten as

$$\tilde{f}(\hat{x}^{k+1}, \mu_k) \leq \tilde{f}(x^k, \mu_k) + \langle \nabla \tilde{f}(x^k, \mu_k), \hat{x}^{k+1} - x^k \rangle + \frac{1}{2} \gamma_k \mu_k^{-1} \|\hat{x}^{k+1} - x^k\|^2.$$

327 Invoking Definition 3.1-(v), (3.10) holds when $\gamma_k \geq L$. Thus the updating of γ_k in Step 2 is at
 328 most $\log_\eta(L/\underline{\gamma}) + 1$ times at each iteration. Hence, the SPG algorithm is well-defined and we have

329 that $\gamma_k \leq \max\{\bar{\gamma}, \rho L\}$, $\forall k \in \mathbb{N}$. From (3.11), it is easy to verify that $x^{k+1} \in \mathcal{X}$ by $x^k \in \mathcal{X}$ and
 330 $\hat{x}^{k+1} \in \mathcal{X}$.

(ii). Since $\{\mu_k\}$ is non-increasing, to prove (ii), we assume that $\lim_{k \rightarrow \infty} \mu_k = \hat{\mu} > 0$ by contradiction. Then, (3.13) happens finite times at most, which means that there exists $K \in \mathbb{N}$ such that $\mu_k = \hat{\mu}$, $\forall k \geq K$. Then,

$$\tilde{\mathcal{F}}(x^{k+1}, \mu_k) + \kappa \mu_k - \tilde{\mathcal{F}}(x^k, \mu_{k-1}) - \kappa \mu_{k-1} \leq -\alpha \hat{\mu}^2, \quad \forall k \geq K + 1.$$

331 We obtain from the above inequality that

$$332 \quad (3.15) \quad \lim_{k \rightarrow \infty} \tilde{\mathcal{F}}(x^{k+1}, \mu_k) + \kappa \mu_k = -\infty.$$

333 However, by $\{x^k\} \subseteq \mathcal{X}$ and (3.1), we see that

$$334 \quad (3.16) \quad \tilde{\mathcal{F}}(x^{k+1}, \mu_k) + \kappa \mu_k \geq \mathcal{F}(x^{k+1}) \geq \min_{x \in \mathcal{X}} \mathcal{F}(x) = \min_{x \in \mathcal{X}} \mathcal{F}_{\ell_0}(x), \quad \forall k \geq K,$$

335 where the last equality follows from Theorem 2.4. Thus, the contradiction between (3.15) and (3.16)
 336 implies (ii). \square

337 LEMMA 3.3. *For any $k \in \mathbb{N}$, we have*

$$338 \quad (3.17) \quad \tilde{\mathcal{F}}(x^{k+1}, \mu_k) - \tilde{\mathcal{F}}(x^k, \mu_k) \leq -\frac{1}{2} \gamma_k \mu_k^{-1} \|x^{k+1} - x^k\|^2,$$

339 which implies $\{\tilde{\mathcal{F}}(x^{k+1}, \mu_k) + \kappa \mu_k\}$ is non-increasing and $\lim_{k \rightarrow \infty} \tilde{\mathcal{F}}(x^{k+1}, \mu_k) = \lim_{k \rightarrow \infty} \mathcal{F}(x^k)$.

340 Moreover, there exists $R > 0$ such that $\|x^k\| \leq R$, $\forall k \in \mathbb{N}$.

Proof. Since $Q_{d^k, \gamma_k}(x, x^k, \mu_k)$ is strongly convex with modulus $\gamma_k \mu_k^{-1}$, using the definition of \hat{x}^{k+1} in (3.9) and $x^{k+1} = \hat{x}^{k+1}$ when (3.10) holds, we obtain

$$Q_{d^k, \gamma_k}(x^{k+1}, x^k, \mu_k) \leq Q_{d^k, \gamma_k}(x, x^k, \mu_k) - \frac{1}{2} \gamma_k \mu_k^{-1} \|x^{k+1} - x\|^2, \quad \forall x \in \mathcal{X}.$$

341 By the definition of function Q_{d^k, γ_k} given in (3.8), upon rearranging the terms, we have

$$342 \quad (3.18) \quad \begin{aligned} \lambda \Phi^{d^k}(x^{k+1}) &\leq \lambda \Phi^{d^k}(x) + \langle x - x^{k+1}, \nabla \tilde{f}(x^k, \mu_k) \rangle \\ &\quad + \frac{1}{2} \gamma_k \mu_k^{-1} \|x - x^k\|^2 - \frac{1}{2} \gamma_k \mu_k^{-1} \|x^{k+1} - x^k\|^2 - \frac{1}{2} \gamma_k \mu_k^{-1} \|x^{k+1} - x\|^2. \end{aligned}$$

343 Moreover, (3.10) can be written as

$$344 \quad (3.19) \quad \tilde{\mathcal{F}}^{d^k}(x^{k+1}, \mu_k) \leq \tilde{f}(x^k, \mu_k) + \langle x^{k+1} - x^k, \nabla \tilde{f}(x^k, \mu_k) \rangle + \frac{1}{2} \gamma_k \mu_k^{-1} \|x^{k+1} - x^k\|^2 + \lambda \Phi^{d^k}(x^{k+1}).$$

345 Summing up (3.18) and (3.19), we notice that

$$346 \quad (3.20) \quad \begin{aligned} \tilde{\mathcal{F}}^{d^k}(x^{k+1}, \mu_k) &\leq \tilde{f}(x^k, \mu_k) + \lambda \Phi^{d^k}(x) + \langle x - x^k, \nabla \tilde{f}(x^k, \mu_k) \rangle \\ &\quad + \frac{1}{2} \gamma_k \mu_k^{-1} \|x - x^k\|^2 - \frac{1}{2} \gamma_k \mu_k^{-1} \|x^{k+1} - x\|^2, \quad \forall x \in \mathcal{X}. \end{aligned}$$

347 For a fixed $\mu > 0$, the convexity of $\tilde{f}(x, \mu)$ with respect to x invokes

$$348 \quad (3.21) \quad \tilde{f}(x^k, \mu_k) + \langle x - x^k, \nabla \tilde{f}(x^k, \mu_k) \rangle \leq \tilde{f}(x, \mu_k), \quad \forall x \in \mathcal{X}.$$

349 Combining (3.20) and (3.21) and recalling the definition of $\tilde{\mathcal{F}}^{d^k}$, one has

$$350 \quad (3.22) \quad \tilde{\mathcal{F}}^{d^k}(x^{k+1}, \mu_k) \leq \tilde{\mathcal{F}}^{d^k}(x, \mu_k) + \frac{1}{2}\gamma_k\mu_k^{-1}\|x - x^k\|^2 - \frac{1}{2}\gamma_k\mu_k^{-1}\|x^{k+1} - x\|^2, \quad \forall x \in \mathcal{X}.$$

351 Letting $x = x^k$ in (3.22) and by $d^k = d^{x^k}$, we obtain

$$352 \quad (3.23) \quad \tilde{\mathcal{F}}^{d^k}(x^{k+1}, \mu_k) + \frac{1}{2}\gamma_k\mu_k^{-1}\|x^{k+1} - x^k\|^2 \leq \tilde{\mathcal{F}}(x^k, \mu_k).$$

353 Thanks to $\tilde{\mathcal{F}}^{d^k}(x^{k+1}, \mu_k) \geq \tilde{\mathcal{F}}(x^{k+1}, \mu_k)$, (3.23) leads to (3.17).

354 Since $\tilde{\mathcal{F}}(x^k, \mu_k) \leq \tilde{\mathcal{F}}(x^k, \mu_{k-1}) + \kappa(\mu_{k-1} - \mu_k)$, by (3.17), we obtain

$$355 \quad (3.24) \quad \tilde{\mathcal{F}}(x^{k+1}, \mu_k) + \kappa\mu_k + \frac{1}{2}\gamma_k\mu_k^{-1}\|x^{k+1} - x^k\|^2 \leq \tilde{\mathcal{F}}(x^k, \mu_{k-1}) + \kappa\mu_{k-1},$$

356 which implies the non-increasing property of $\{\tilde{\mathcal{F}}(x^{k+1}, \mu_k) + \kappa\mu_k\}$. Together this result with (3.16)

357 ensures the existence of $\lim_{k \rightarrow \infty} \tilde{\mathcal{F}}(x^{k+1}, \mu_k) + \kappa\mu_k$. By virtue of $\lim_{k \rightarrow \infty} \mu_k = 0$ and Definition

358 3.1-(i), we get $\lim_{k \rightarrow \infty} \tilde{\mathcal{F}}(x^{k+1}, \mu_k) = \lim_{k \rightarrow \infty} \mathcal{F}(x^k)$.

Recalling the non-increasing property of $\{\tilde{\mathcal{F}}(x^{k+1}, \mu_k) + \kappa\mu_k\}$ again, we see that

$$\mathcal{F}(x^{k+1}) \leq \tilde{\mathcal{F}}(x^{k+1}, \mu_k) + \kappa\mu_k \leq \tilde{\mathcal{F}}(x^1, \mu_0) + \kappa\mu_0 < \infty.$$

We then obtain the boundedness of $\{x^k\}$ from $\{x^k\} \subseteq \mathcal{X}$ and the level bounded assumption of \mathcal{F} on \mathcal{X} . Observe that

$$\mathcal{F}_{\ell_0}(x) \geq \mathcal{F}(x) = \mathcal{F}_{\ell_0}(x) - \lambda \sum_{|x_i| < \nu} (1 - |x_i|/\nu) \geq \mathcal{F}_{\ell_0}(x) - \lambda n, \quad \forall x \in \mathbb{R}^n.$$

359 Then, it is easy to verify the level boundedness of \mathcal{F} by the level boundedness of \mathcal{F}_{ℓ_0} on \mathcal{X} . Hence,

360 the same results in Lemma 3.3 hold when \mathcal{F}_{ℓ_0} is level bounded on \mathcal{X} . \square

361 The following proposition shows that there exists a subsequence of $\{x^k\}$ converging to a lifted

362 stationary point of (1.6), which lays a foundation for the sequence convergence of $\{x^k\}$.

363 PROPOSITION 3.4. *Any accumulation point of $\{x^k : k \in \mathcal{N}^s\}$ is a lifted stationary point of*

364 (1.6).

365 *Proof.* When \mathcal{F} (or \mathcal{F}_{ℓ_0}) is level bounded on \mathcal{X} , by Lemma 3.3, $\{x^k\}$ is bounded. Suppose \bar{x} is

366 an accumulation point of $\{x^k\}_{k \in \mathcal{N}^s}$ with the convergence of subsequence $\{x^{k_i}\}_{k_i \in \mathcal{N}^s}$.

367 Since (3.12) fails for $k_i \in \mathcal{N}^s$, by rearranging (3.24), we obtain that $\gamma_{k_i}\mu_{k_i}^{-1}\|x^{k_i+1} - x^{k_i}\|^2 \leq$

368 $2\alpha\mu_{k_i}^2$, which gives $\|x^{k_i+1} - x^{k_i}\| \leq \sqrt{2\alpha\gamma_{k_i}^{-1}\mu_{k_i}^3}$. Thus, $\gamma_{k_i}\mu_{k_i}^{-1}\|x^{k_i+1} - x^{k_i}\| \leq \sqrt{2\alpha\gamma_{k_i}\mu_{k_i}}$, which

369 together with $\lim_{i \rightarrow \infty} \mu_{k_i} = 0$ and $\{\gamma_{k_i}\} \subseteq [\gamma, \max\{\bar{\gamma}, \rho L\}]$ implies

$$370 \quad (3.25) \quad \lim_{i \rightarrow \infty} \gamma_{k_i}\mu_{k_i}^{-1}\|x^{k_i+1} - x^{k_i}\| = 0 \quad \text{and} \quad \lim_{i \rightarrow \infty} x^{k_i+1} = \bar{x}.$$

371

372 Recalling $x^{k_i+1} = \hat{x}^{k_i+1}$ defined in (3.9) and by its first order necessary optimality condition,
373 we have

$$374 \quad (3.26) \quad \langle \nabla \tilde{f}(x^{k_i}, \mu_{k_i}) + \gamma_{k_i} \mu_{k_i}^{-1} (x^{k_i+1} - x^{k_i}) + \lambda \zeta^{k_i}, x - x^{k_i+1} \rangle \geq 0, \quad \forall \zeta^{k_i} \in \partial \Phi^{d^{k_i}}(x^{k_i+1}), x \in \mathcal{X}.$$

375 Since the elements in $\{d^{k_i} : i \in \mathbb{N}\}$ are finite and $\lim_{i \rightarrow \infty} x^{k_i+1} = \bar{x}$, there exists a subsequence
376 of $\{k_i\}$, denoted as $\{k_{i_j}\}$, and $\bar{d} \in \mathcal{D}(\bar{x})$ such that $d^{k_{i_j}} = \bar{d}, \forall j \in \mathbb{N}$. By the upper semicontinuity
377 of $\partial \Phi^{\bar{d}}$ and $\lim_{j \rightarrow \infty} x^{k_{i_j}+1} = \bar{x}$, it gives

$$378 \quad (3.27) \quad \left\{ \lim_{j \rightarrow \infty} \zeta^{k_{i_j}} : \zeta^{k_{i_j}} \in \partial \Phi^{d^{k_{i_j}}}(x^{k_{i_j}+1}) \right\} \subseteq \partial \Phi^{\bar{d}}(\bar{x}).$$

379 Along with the subsequence $\{k_{i_j}\}$ and letting $j \rightarrow \infty$ in (3.26), from Definition 3.1-(iii), (3.25) and
380 (3.27), we obtain that there exist $\bar{\xi} \in \partial f(\bar{x})$ and $\bar{\zeta} \in \partial \Phi^{\bar{d}}(\bar{x})$ such that

$$381 \quad (3.28) \quad \langle \bar{\xi} + \lambda \bar{\zeta}, x - \bar{x} \rangle \geq 0, \quad \forall x \in \mathcal{X}.$$

382 By $\bar{d} \in \mathcal{D}(\bar{x})$, the definition of $\Phi^{\bar{d}}$ in (2.5) and the convexity of \mathcal{X} , (3.28) implies that \bar{x} is a lifted
383 stationary point of (1.6). \square

384 **Remark 3.1.** *The convexity of Φ^d plays an important role in the analysis of the SPG algo-*
385 *rithm. It is easy to check that all the results in subsection 3.4 are true when the penalty can be*
386 *described by the min of a class of simple convex functions whose proximal operators can be calculated*
387 *effectively.*

388 **3.5. Global sequence convergence of the SPG algorithm for problem (1.6).** It is
389 interesting that the proposed SPG algorithm for this kind of nonconvex nonsmooth optimization
390 problem owns the global sequence convergence without the K-L condition or error bound condition
391 on the objective function, while the special structure of the continuous relaxation for $\|x\|_0$ and the
392 updating rule for μ_k are the key points. Throughout this subsection, the analysis uses the same
393 assumptions in subsection 3.4.

394 We begin this subsection by giving some preliminary analysis, which are Lemma 3.5, Lemma
395 3.6 and Proposition 3.7. Based on these results, we present the two main results for the SPG
396 algorithm: the sequence convergence of $\{x^k\}$ in Theorem 3.8; the local convergence rate of $\{\mathcal{F}(x^k)\}$
397 and the finite-iteration identification of $\mathcal{A}(x^k)$ in Theorem 3.9.

398 **LEMMA 3.5.** *The following statements hold:*

- 399 (i) $\sum_{k=0}^{\infty} \gamma_k \mu_k^{-1} \|x^{k+1} - x^k\|^2 \leq 2 (\mathcal{F}(x^0, \mu_{-1}) + \kappa \mu_{-1} - \min_{\mathcal{X}} \mathcal{F});$
400 (ii) $\sum_{k=0}^{\infty} \mu_k^2 \leq \Lambda$ with $\Lambda = \frac{1}{\alpha} \left(\tilde{\mathcal{F}}(x^0, \mu_{-1}) + \kappa \mu_{-1} - \min_{x \in \mathcal{X}} \mathcal{F}(x) \right) + \frac{2\mu_0^2 \sigma}{2\sigma - 1} < \infty;$
401 (iii) $\mathcal{A}(x^{k+1}) \subseteq \mathcal{A}(x^k).$

402 *Proof.* (i). Recalling (3.24), for all $k \in \mathbb{N}$, we obtain

$$403 \quad (3.29) \quad \gamma_k \mu_k^{-1} \|x^{k+1} - x^k\|^2 \leq 2 \left(\tilde{\mathcal{F}}(x^k, \mu_{k-1}) + \kappa \mu_{k-1} - \tilde{\mathcal{F}}(x^{k+1}, \mu_k) - \kappa \mu_k \right).$$

404 Summing up the above inequality over $k = 0, \dots, K$, it gives

$$405 \quad (3.30) \quad \sum_{k=0}^K \gamma_k \mu_k^{-1} \|x^{k+1} - x^k\|^2 \leq 2 \left(\tilde{\mathcal{F}}(x^0, \mu_{-1}) + \kappa \mu_{-1} - \tilde{\mathcal{F}}(x^{K+1}, \mu_K) - \kappa \mu_K \right).$$

406 By letting K in (3.30) tend to infinity and along with (3.16), we obtain (i).
 407 (ii). From (3.14), we have

$$408 \quad (3.31) \quad \sum_{k \in \mathcal{N}^s} \mu_k^2 = \sum_{r=1}^{\infty} \mu_0^2 \frac{1}{(n_r^s + 1)^{2\sigma}} \leq \sum_{k=1}^{\infty} \frac{\mu_0^2}{k^{2\sigma}} \leq \frac{2\mu_0^2\sigma}{2\sigma - 1},$$

409 where n_r^s is the r th smallest element in \mathcal{N}^s .

410 When $k \notin \mathcal{N}^s$, (3.12) gives $\alpha\mu_k^2 \leq \tilde{\mathcal{F}}(x^k, \mu_{k-1}) + \kappa\mu_{k-1} - \tilde{\mathcal{F}}(x^{k+1}, \mu_k) - \kappa\mu_k$, which together
 411 with the non-increasing property of $\{\tilde{\mathcal{F}}(x^{k+1}, \mu_k) + \kappa\mu_k\}$ and (3.16) implies

$$412 \quad (3.32) \quad \sum_{k \notin \mathcal{N}^s} \mu_k^2 \leq \frac{1}{\alpha} \left(\tilde{\mathcal{F}}(x^0, \mu_{-1}) + \kappa\mu_{-1} - \min_{\mathcal{X}} \mathcal{F} \right).$$

413 Combining (3.31) and (3.32), we finish the proof for the estimation in item (ii).

(iii). We only need to prove that if $x_i^k = 0$, then $x_i^{k+1} = 0$. If $x_i^k = 0$, we get $d_i^k = 1$. From
 (3.7) and $\nu < \lambda/L_f$, we have

$$\left| x_i^k - \gamma_k^{-1} \mu_k \nabla_i \tilde{f}(x^k, \mu_k) \right| \leq \gamma_k^{-1} \mu_k \left\| \nabla \tilde{f}(x^k, \mu_k) \right\| \leq (\lambda \gamma_k^{-1} \mu_k) / \nu.$$

414 By (3.7), we obtain $x_i^{k+1} = 0$, which completes the proof of this statement. \square

415 For $\{x^k\}$, denote

$$416 \quad (3.33) \quad \mathcal{N}_1 = \{k \in \mathbb{N} : \text{there exists } i \in \{1, \dots, n\} \text{ such that } 0 < |x_i^k| < \nu\}.$$

417 Next lemma gives some estimation on $\{x^k\}$ and $\{\mu_k\}$ when k is sufficiently large.

418 LEMMA 3.6. *There exists $K \in \mathbb{N}$ such that for all $k \geq K$, it holds*

- 419 (i) $\left\| \nabla \tilde{f}(x^k, \mu_k) \right\| < \frac{1}{2} (\lambda/\nu + L_f)$;
 420 (ii) $\|x^{k+1} - x^k\| \leq 3(\lambda/\nu) \sqrt{n} \gamma^{-1} \mu_k$;
 421 (iii) for any $k \in \mathcal{N}_1$, either $\|x^{k+1}\|_0 \leq \|x^k\|_0 - 1$ or $\|x^{k+1} - x^k\| \geq \frac{1}{2} (\lambda/\nu - L_f) \gamma_k^{-1} \mu_k$;
 422 (iv) $\sum_{k \in \mathcal{N}_1, k \geq K} \|x^{k+1} - x^k\| < \infty$ and $\sum_{k \in \mathcal{N}_1, k \geq K} \mu_k < \infty$.

423 *Proof.* (i). We argue it by contradiction. Suppose there is a subsequence of $\{x^k\}$, denoted by
 424 $\{x^{k_i}\}$, such that

$$425 \quad (3.34) \quad \left\| \nabla \tilde{f}(x^{k_i}, \mu_{k_i}) \right\| \geq \frac{1}{2} (\lambda/\nu + L_f) > L_f, \quad \forall i \in \mathbb{N}.$$

426 Since $\{x^{k_i}\}$ is bounded, which is proved in Lemma 3.3, there exists a subsequence of $\{x^{k_i}\}$ (also
 427 denoted by $\{x^{k_i}\}$ for simplicity) and $\bar{x} \in \mathcal{X}$ such that $\lim_{i \rightarrow \infty} x^{k_i} = \bar{x}$. Due to $\lim_{i \rightarrow \infty} \mu_{k_i} = 0$,
 428 the property of \tilde{f} in Definition 3.1-(iii), λ/ν and (3.34) imply the existence of $\bar{\xi} \in \partial f(\bar{x})$ such that
 429 $\|\bar{\xi}\| > L_f$, which leads to a contradiction to the definition of L_f given in Assumption 1. Hence, we
 430 establish result (i) in this lemma.

(ii). For any $i \in \{1, 2, \dots, n\}$, by (3.7) and $L_f < \lambda/\nu$, we have

$$\left| x_i^{k+1} - x_i^k \right| \leq 2(\lambda/\nu) \gamma_k^{-1} \mu_k + \gamma_k^{-1} \mu_k \left| \nabla_i \tilde{f}(x^k, \mu_k) \right| \leq 3(\lambda/\nu) \gamma_k^{-1} \mu_k,$$

431 which completes the proof for item (ii).

432 (iii). Denote $w^k = x^k - \gamma_k^{-1} \mu_k \nabla \tilde{f}(x^k, \mu_k)$. For a fixed $k \in \mathcal{N}_1$ and $k \geq K$, there exists j
 433 such that $0 < |x_j^k| < \nu$. Then, $d_j^k = 1$ by (2.4). Next, we will prove that either $x_j^{k+1} = 0$ or
 434 $|x_j^{k+1} - x_j^k| \geq \frac{1}{2}(\lambda/\nu - L_f) \gamma_k^{-1} \mu_k$. We split the proof into three cases.

435 Case 1. If $|w_j^k| \leq (\lambda/\nu) \gamma_k^{-1} \mu_k$, by (3.7), we get $x_j^{k+1} = 0$, which together with $\mathcal{A}(x^{k+1}) \subseteq \mathcal{A}(x^k)$
 436 implies $\|x^{k+1}\|_0 \leq \|x^k\|_0 - 1$.

Case 2. If $w_j^k > (\lambda/\nu) \gamma_k^{-1} \mu_k$, by (3.7) and result (i) of this lemma, we obtain that

$$|x_j^{k+1} - x_j^k| \geq (\lambda/\nu) \gamma_k^{-1} \mu_k - \left| \gamma_k^{-1} \mu_k \nabla_i \tilde{f}(x^k, \mu_k) \right| \geq \frac{1}{2} (\lambda/\nu - L_f) \gamma_k^{-1} \mu_k,$$

437 which implies

$$438 \quad (3.35) \quad \|x^{k+1} - x^k\| \geq \frac{1}{2} (\lambda/\nu - L_f) \gamma_k^{-1} \mu_k.$$

439 Case 3. If $w_j^k < -(\lambda/\nu) \gamma_k^{-1} \mu_k$, similar to the analysis in Case 1, we see that (3.35) holds. Thus,
 440 we complete the proof of statement (iii).

441 (iv). We introduce the notations $\mathcal{N}_{11} = \{k \in \mathcal{N}_1 : k \geq K, \|x^{k+1}\|_0 \leq \|x^k\|_0 - 1\}$ and
 442 $\mathcal{N}_{12} = \{k : k \geq K, k \in \mathcal{N}_1 \setminus \mathcal{N}_{11}\}$. By Lemma 3.5-(iii), \mathcal{N}_{11} has at most n elements. From result (iii)
 443 of this lemma, we have $\gamma_k \mu_k^{-1} \|x^{k+1} - x^k\| \geq \frac{1}{2}(\lambda/\nu - L_f)$, $\forall k \in \mathcal{N}_{12}$. Then, we have

$$444 \quad (3.36) \quad \frac{1}{2}(\lambda/\nu - L_f) \sum_{k \in \mathcal{N}_{12}} \|x^{k+1} - x^k\| \leq \sum_{k \in \mathcal{N}_{12}} \gamma_k \mu_k^{-1} \|x^{k+1} - x^k\|^2 \leq 2 \left(\tilde{\mathcal{F}}(x^0, \mu_{-1}) + \kappa \mu_{-1} - \min_{\mathcal{X}} \mathcal{F} \right),$$

where the second inequality follows from Lemma 3.5-(i). (3.36) implies $\sum_{k \in \mathcal{N}_{12}} \|x^{k+1} - x^k\| < \infty$,
 which together with the finiteness of the elements in \mathcal{N}_{11} gives $\sum_{k \in \mathcal{N}_1, k \geq K} \|x^{k+1} - x^k\| < \infty$.
 Moreover,

$$\sum_{k \in \mathcal{N}_{12}} \gamma_k \mu_k^{-1} \|x^{k+1} - x^k\|^2 = \sum_{k \in \mathcal{N}_{12}} (\gamma_k \mu_k^{-1} \|x^{k+1} - x^k\|)^2 \gamma_k^{-1} \mu_k \geq \frac{1}{4} (\lambda/\nu - L_f)^2 \sum_{k \in \mathcal{N}_{12}} \gamma_k^{-1} \mu_k,$$

445 which together with the second inequality of (3.36) and Lemma 3.2-(i) implies $\sum_{k \in \mathcal{N}_{12}} \mu_k < \infty$.
 446 By $\sum_{k \in \mathcal{N}_{11}} \mu_k \leq n \mu_0$, we conclude that $\sum_{k \in \mathcal{N}_1, k \geq K} \mu_k < \infty$. \square

447 The next proposition explores that all accumulation points of $\{x^k\}$ own a common support set
 448 and a unified lower bound, which provides the main technical support for the forthcoming Theorem
 449 3.8.

PROPOSITION 3.7. Denote $\bar{\mathcal{X}} = \{\bar{x} \in \mathcal{X} : \bar{x} \text{ is an accumulation point of } \{x^k\}\}$, then there
 exists $\mathcal{A}(\bar{\mathcal{X}}) \subseteq \{1, 2, \dots, n\}$ such that for any $\bar{x} \in \bar{\mathcal{X}}$, it holds that

$$|\bar{x}_i| \geq \nu \text{ for any } i \in \mathcal{A}(\bar{\mathcal{X}}) \text{ and } \bar{x}_i = 0 \text{ for any } i \notin \mathcal{A}(\bar{\mathcal{X}}).$$

450 *Proof.* We first prove the following result:

$$451 \quad (3.37) \quad \text{for any } \bar{x} \in \bar{\mathcal{X}} \text{ and any } i \in \{1, \dots, n\}, \text{ either } \bar{x}_i = 0 \text{ or } |\bar{x}_i| \geq \nu.$$

452 If (3.37) does not hold, there exists $\hat{x} \in \bar{\mathcal{X}}$ with the convergence sequence $\{x^{k_j}\}$ and $\iota \in \{1, \dots, n\}$
 453 such that $0 < |\hat{x}_\iota| < \nu$. In what follows, without loss of generality, we suppose $\hat{x}_\iota > 0$.

454 Since any accumulation point of $\{x^k\}_{k \in \mathcal{N}^s}$ is an accumulation point of $\{x^k\}$, there exists $\bar{x} \in \bar{\mathcal{X}}$
 455 and a subsequence of $\{x^k\}$, denoted by $\{x^{t_j}\}$, such that $\lim_{j \rightarrow \infty} x^{t_j} = \bar{x}$. By taking subsequences of
 456 $\{x^{k_j}\}$ and $\{x^{t_j}\}$ if necessary, we assume for the simplicity of notation that $k_j < t_j < k_{j+1}, \forall j \in \mathbb{N}$.
 457 Combining Proposition 2.5, Lemma 2.3 and Proposition 3.4, either $\bar{x}_l = 0$ or $|\bar{x}_l| \geq \nu$.

Let $\varepsilon = \min \left\{ \frac{\nu - \hat{x}_l}{2}, \frac{\hat{x}_l}{4} \right\} > 0$. If $\bar{x}_l = 0$, there exists $J \in \mathbb{N}$ such that

$$|x_l^{k_j} - \hat{x}_l| \leq \varepsilon \quad \text{and} \quad |x_l^{t_j}| \leq \varepsilon, \quad \forall j \geq J,$$

458 which implies

$$459 \quad (3.38) \quad \frac{3}{4}\hat{x}_l \leq \hat{x}_l - \varepsilon \leq x_l^{k_j} \leq \varepsilon + \hat{x}_l \leq \frac{\nu + \hat{x}_l}{2} < \nu \quad \text{and} \quad -\frac{1}{4}\hat{x}_l \leq x_l^{t_j} \leq \frac{1}{4}\hat{x}_l, \quad \forall j \geq J.$$

460 Then, $x_l^{k_j} - x_l^{t_j} \geq \frac{1}{2}\hat{x}_l, \forall j \geq J$. Thus,

$$461 \quad (3.39) \quad \sum_{j=J}^{\infty} |x_l^{t_j} - x_l^{k_j}| = +\infty.$$

If there exists $r \geq J$, such that $x_l^{t_r} = 0$, Lemma 3.5-(iii) gives $x_l^{k_{j+1}} = 0, \forall j \geq r$, which leads to a contradiction to the first inequality in (3.38). Thus, (3.38) gives $0 < |x_l^{k_j}| < \nu$ and $0 < |x_l^{t_j}| < \nu$, which implies $\{x^{k_j}, x^{t_j} : j \geq J\} \subseteq \mathcal{N}_1$ with \mathcal{N}_1 defined in (3.33). Together this with Lemma 3.6-(ii), (iv) and $\lim_{k \rightarrow \infty} \mu_k = 0$, there exists $J_1 \geq J$ such that

$$\sum_{j=J_1}^{\infty} |x_l^{t_j} - x_l^{k_j}| \leq \sum_{k \in \mathcal{N}_1, k \geq K} \|x^{k+1} - x^k\| < \infty,$$

462 which leads to a contradiction to (3.39). Likewise, we can obtain a similar contradiction when
 463 $|\bar{x}_l| \geq \nu$. Therefore, the above analysis ensures the validity of statement (3.37). Together (3.37)
 464 with $\lim_{k \rightarrow \infty} \|x^{k+1} - x^k\| = 0$, we complete the proof of this proposition. \square

465 We next prove the global sequence convergence of iterates $\{x^k\}$.

466 **THEOREM 3.8.** *The iterates $\{x^k\}$ generated by the SPG algorithm is globally convergent to*
 467 *a lifted stationary point of (1.6), i.e. there exists a lifted stationary point \bar{x} of (1.6) such that*
 468 *$\lim_{k \rightarrow \infty} x^k = \bar{x}$.*

469 *Proof.* Let K be a positive integer such that the estimations in Lemma 3.6 hold and \bar{x} be an
 470 accumulation point of $\{x^k\}_{k \in \mathcal{N}^s}$. Suppose $\{x^{k_j}\}$ is a subsequence of $\{x^k\}$ such that

$$471 \quad (3.40) \quad \lim_{j \rightarrow \infty} x^{k_j} = \bar{x}.$$

472 By Proposition 3.4, \bar{x} is a lifted stationary point of (1.6).

From Lemma 2.3, for any $i \in \{1, \dots, n\}$, either $\bar{x}_i = 0$ or $|\bar{x}_i| \geq \nu$. Denote

$$\mathcal{N}(\bar{x}) = \{k \in \mathbb{N} : d_i^k \in \mathcal{D}(\bar{x}_i), \forall i = 1, \dots, n\},$$

473 where $\mathcal{D}(\bar{x}_i)$ is defined in (2.1). We then evaluate $\|x^{k+1} - \bar{x}\|^2$ by considering two cases.

474 Case 1: In this case, we consider the iteration for $k \in \mathcal{N}(\bar{x})$, which implies that $\tilde{\mathcal{F}}^{d^k}(\bar{x}, \mu_k) =$
 475 $\tilde{\mathcal{F}}(\bar{x}, \mu_k)$. Letting $x = \bar{x}$ in (3.22), we have

$$476 \quad \tilde{\mathcal{F}}^{d^k}(x^{k+1}, \mu_k) - \tilde{\mathcal{F}}(\bar{x}, \mu_k) \leq \frac{1}{2}\gamma_k \mu_k^{-1} \|x^k - \bar{x}\|^2 - \frac{1}{2}\gamma_k \mu_k^{-1} \|x^{k+1} - \bar{x}\|^2,$$

477 combining which with (3.1) and (3.4), we obtain

$$478 \quad (3.41) \quad 2\gamma_k^{-1}\mu_k \left(\tilde{\mathcal{F}}(x^{k+1}, \mu_k) + \kappa\mu_k - \mathcal{F}(\bar{x}) \right) \leq \|x^k - \bar{x}\|^2 - \|x^{k+1} - \bar{x}\|^2 + 4\kappa\gamma_k^{-1}\mu_k^2.$$

479 Due to the non-increasing property of $\left\{ \tilde{\mathcal{F}}(x^{k+1}, \mu_k) + \kappa\mu_k \right\}$ and $\lim_{k \rightarrow \infty} \tilde{\mathcal{F}}(x^{k+1}, \mu_k) + \kappa\mu_k = \mathcal{F}(\bar{x})$,
480 we obtain

$$481 \quad (3.42) \quad \|x^{k+1} - \bar{x}\|^2 \leq \|x^k - \bar{x}\|^2 + 4\kappa\gamma_k^{-1}\mu_k^2, \quad \forall k \in \mathcal{N}(\bar{x}).$$

Case 2: In this case, we consider the iteration for $k \notin \mathcal{N}(\bar{x})$. From Proposition 3.7, there exists $K_1 \geq K$ such that for any $k \geq K_1$, it holds

$$|x_i^k| < \nu/2 \text{ for } i \notin \mathcal{A}(\bar{\mathcal{X}}) \text{ and } |x_i^k| \geq \nu/2 \text{ for } i \in \mathcal{A}(\bar{\mathcal{X}}),$$

482 where $\mathcal{A}(\bar{\mathcal{X}})$ is defined in Proposition 3.7.

483 Hence, for $k \notin \mathcal{N}(\bar{x})$ and $k \geq K_1$, there exists $i^k \in \mathcal{A}(\bar{\mathcal{X}})$ such that $\nu/2 \leq |x_{i^k}^k| < \nu$, which
484 means that $k \in \mathcal{N}_1$ with \mathcal{N}_1 defined in (3.33). Then,

$$485 \quad (3.43) \quad \begin{aligned} \|x^{k+1} - \bar{x}\|^2 &= \|x^k - \bar{x}\|^2 + \|x^{k+1} - x^k\|^2 + 2\langle x^{k+1} - x^k, x^k - \bar{x} \rangle \\ &\leq \|x^k - \bar{x}\|^2 + c_1\mu_k^2 + 4R\|x^{k+1} - x^k\|, \quad \forall k \notin \mathcal{N}(\bar{x}), \end{aligned}$$

486 where $c_1 = 9(\lambda/\nu)^2 n\gamma^{-2}$ follows from Lemma 3.6-(ii), and R comes from Lemma 3.3.

487 By (3.42) and (3.43), for any $t \geq K_1$ and $s \in \mathbb{N}$, we have

$$488 \quad (3.44) \quad \|x^{t+s+1} - \bar{x}\|^2 \leq \|x^t - \bar{x}\|^2 + c_2 \sum_{k=t}^{t+s} \mu_k^2 + 4R \sum_{\substack{k=t, \\ k \notin \mathcal{N}(\bar{x})}}^{t+s} \|x^{k+1} - x^k\|,$$

489 where $c_2 = \max\{4\kappa\gamma^{-1}, c_1\}$.

490 Fix an $\epsilon > 0$. There exists $K_2 \geq K_1$ such that when $k_j \geq K_2$, it holds that

$$491 \quad (3.45) \quad \|x^{k_j} - \bar{x}\|^2 \leq \epsilon^2/3, \quad \sum_{k=k_j}^{\infty} \mu_k^2 \leq \epsilon^2/3c_2, \quad \sum_{\substack{k=k_j, \\ k \notin \mathcal{N}(\bar{x})}}^{\infty} \|x^{k+1} - x^k\| \leq \epsilon^2/12R,$$

492 where the first inequality follows from (3.40), the second inequality follows from Lemma 3.5-(ii),
493 and the third inequality follows from and $\{k : k \geq K_1, k \notin \mathcal{N}(\bar{x})\} \subseteq \mathcal{N}_1$ and Lemma 3.6-(iv).

494 Letting $t = k_j$ in (3.44) with $k_j \geq K_2$, from (3.45), we obtain $\|x^k - \bar{x}\| \leq \epsilon$, $\forall k \geq K_3$, where
495 $K_3 = \min\{k_j : k_j \geq K_2\}$. Due to the arbitrariness of $\epsilon > 0$, we get $\lim_{k \rightarrow \infty} x^k = \bar{x}$. \square

496 The lower bound property is used to prove the estimation in Lemma 3.6-(iii), which is the key
497 point to guarantee the global sequence convergence of $\{x^k\}$. Without this lower bound property,
498 due to the nonconvexity of the objective function in (1.6), it is almost impossible to propose a global
499 sequence convergence algorithm without the regularity conditions. Among the existing penalties,
500 only capped- ℓ_1 penalty can be expressed by the min of a class of simple convex functions and make
501 the stationary points of the corresponding minimization problem own a unified lower bound. This
502 is the main motivation of this paper on studying the cardinality penalty problem by the capped- ℓ_1

503 relaxation. Moreover, from the proof of Theorem 3.9, we find that the descent criterion and the
 504 updating method for μ_k are also important to guarantee the global sequence convergence of $\{x^k\}$,
 505 since it needs that $\sum_{k=1}^{\infty} \mu_k^2 < +\infty$.

506 The limit point of $\{x^k\}$ is most likely different with different initial iterates x^0 and μ_0 . The
 507 zero vector is a trivial ν -strong local minimizer of (1.2), which is not we want. By property (iii)
 508 of Lemma 3.5, our theoretical results hold for any initial iterate $x^0 \in \mathcal{X}$. To find interesting ν -
 509 strong local minimizers, we chose x^0 without zero component in the numerical experiments. How to
 510 choose an initial point such that the accumulation point of $\{x^k\}$ is a global minimizer (or an oracle
 511 solution) of (1.2) is an interesting work. To the best of our knowledge, it is still an open problem.
 512 [24, Theorem 1] gave some discussion on this topic for the linear approximation algorithm to solve
 513 the sparsity problem with SCAD penalty. Similar results can be expected for the SPG algorithm.

514 The following theorem gives a local convergence rate of the SPG algorithm on the objective
 515 function values of problem (1.6) and the finite iteration convergence of $\{x^k\}$ in a subspace.

516 **THEOREM 3.9.** *There exist $c > 0$ and $K \in \mathbb{N}$ such that, for $k \geq K$, we have*

$$517 \quad (3.46) \quad \mathcal{F}(x^{k+1}) - \mathcal{F}(\bar{x}) \leq ck^{-(1-\sigma)} \quad \text{and} \quad \left\| x_{\mathcal{A}(\bar{x})^c}^k - \bar{x}_{\mathcal{A}(\bar{x})^c} \right\| = 0,$$

518 where \bar{x} is the limit of $\{x^k\}$.

519 *Proof.* Denote $\epsilon = \min\{\nu, \min\{|\bar{x}_i| - \nu : |\bar{x}_i| > \nu, i = 1, \dots, n\}\}$. From Theorem 3.8, there ex-
 520 ists $K_1 \in \mathbb{N}$ such that $\|x^k - \bar{x}\| < \epsilon, \forall k \geq K_1$. Then, $k \in \mathcal{N}(\bar{x}), \forall k \geq K_1$.

521 From the proof of Theorem 3.8, (3.41) holds for any $k \geq K_1$. Summing up (3.41) for $k =$
 522 $K_1, K_1 + 1, \dots, K_1 + t$, we have

$$523 \quad (3.47) \quad \begin{aligned} & 2t \max\{\bar{\gamma}, \rho L\}^{-1} \mu_{K_1+t} \left(\tilde{\mathcal{F}}(x^{K_1+t+1}, \mu_{K_1+t}) + \kappa \mu_{K_1+t} - \mathcal{F}(\bar{x}) \right) \\ & \leq \|x^{K_1} - \bar{x}\|^2 - \|x^{K_1+t+1} - \bar{x}\|^2 + 4\kappa \sum_{k=K_1}^{K_1+t} \gamma_k^{-1} \mu_k^2, \end{aligned}$$

524 where we use $\tilde{\mathcal{F}}(x^{k+1}, \mu_k) + \kappa \mu_k \geq \mathcal{F}(\bar{x}), \{\gamma_k\} \subseteq [\underline{\gamma}, \max\{\bar{\gamma}, \rho L\}]$, and the non-increasing property
 525 of $\{\mu_k\}$ and $\{\tilde{\mathcal{F}}(x^{k+1}, \mu_k) + \kappa \mu_k\}$.

526 We first consider the right hand side of (3.47). We observe that $4\kappa \sum_{k=K_1}^{K_1+t} \gamma_k^{-1} \mu_k^2 \leq 4\kappa \underline{\gamma}^{-1} \Lambda$,
 527 where Λ is defined in Lemma 3.5-(ii). Then,

$$528 \quad (3.48) \quad \|x^{K_1} - \bar{x}\|^2 - \|x^{K_1+t+1} - \bar{x}\|^2 + 4\kappa \sum_{k=K_1}^{K_1+t} \gamma_k^{-1} \mu_k^2 \leq 4R^2 + 4\kappa \underline{\gamma}^{-1} \Lambda, \quad \forall t \in \mathbb{N},$$

529 with R defined in Lemma 3.3.

By $\mu_{K_1+t} \geq \mu_0(K_1+t)^{-\sigma}$ and $\tilde{\mathcal{F}}(x^{K_1+t+1}, \mu_{K_1+t}) + \kappa \mu_{K_1+t} \geq \mathcal{F}(x^{K_1+t+1}), \forall t \in \mathbb{N}$, we observe
 from (3.47) and (3.48) that

$$\mathcal{F}(x^{K_1+t+1}) - \mathcal{F}(\bar{x}) \leq \left(\frac{(4R^2 + 4\kappa \underline{\gamma}^{-1} \Lambda) \max\{\bar{\gamma}, \rho L\}}{2\mu_0} \right) \left(\frac{(K_1+t)^\sigma}{t} \right).$$

Therefore, letting $c = (4R^2 + 4\kappa \underline{\gamma}^{-1} \Lambda) \max\{\bar{\gamma}, \rho L\} / \mu_0$, we obtain

$$\mathcal{F}(x^{k+1}) - \mathcal{F}(\bar{x}) \leq \frac{c}{2} \left(\frac{k^\sigma}{k - K_1} \right) \leq ck^{-(1-\sigma)}, \quad \forall k \geq 2K_1.$$

530 To prove the second statement in (3.46), we argue it by contradiction. If there is no $K \in \mathbb{N}$
 531 such that $x_i^k = 0$ for all $i \in \mathcal{A}(\bar{x})^c$ and $k \geq K$, then there is a subsequence of $\{x^k\}$, denoted by
 532 $\{x^{k_j}\}$, and $\hat{i} \in \mathcal{A}(\bar{x})^c$ such that $|x_{\hat{i}}^{k_j}| \neq 0$. Since $\mathcal{A}(x^{k+1}) \subseteq \mathcal{A}(x^k)$ and $\lim_{k \rightarrow \infty} x^k = \bar{x}$, the above
 533 assumption implies that there exists $K_1 \in \mathbb{N}$ such that $0 < |x_{\hat{i}}^k| < \nu$, $\forall k \geq K_1$. Thus, for all
 534 $k \geq K_1$, it gives $k \in \mathcal{N}_1$ with \mathcal{N}_1 given in (3.33). Recalling Lemma 3.6-(iv), we get $\sum_{k=K_1}^{\infty} \mu_k < \infty$.
 535 However, due to $\mu_k \geq \mu_0 k^{-\sigma}$ with $\sigma < 1$, we have $\sum_{k=K_1}^{\infty} \mu_k = \infty$, which leads to a contradiction.
 536 Therefore, the second statement in (3.46) holds. \square

537 Following the proof of Theorem 3.9, the local convergence rate of $\mathcal{F}(x^k) - \mathcal{F}(\bar{x})$ is $O(\frac{1}{k\mu_k})$.
 538 Moreover, thanks to the lower bound property, the SPG algorithm owns the finite iteration identi-
 539 fication on the support set of the limit point of $\{x^k\}$, which inspires us that the local convergence
 540 rate can be improved when f satisfies some proper conditions. For example, when f is strongly
 541 convex with modulus $\delta > 0$, then the local convergence rate can be exponential; when f satisfies
 542 the K-L inequality on \mathcal{X} with exponent $\alpha \in [0, 1)$, then $\{x^k\}$ is convergent finitely if $\alpha = 0$, linearly
 543 if $\alpha \in (0, \frac{1}{2}]$ and sublinearly if $\alpha \in (\frac{1}{2}, 1)$.

544 **4. Numerical experiments.** To verify and illustrate the performance of the continuous re-
 545 laxation (1.6) and the SPG algorithm, we use a test example and generate two examples randomly
 546 with normal distribution. All experiments are performed in MATLAB 2016a on a Lenovo PC
 547 (3.00GHz, 2.00GB of RAM). In the following examples, the stopping criterion is set as

$$548 \quad (4.1) \quad \text{number of iterations} \leq \text{Maxiter} \quad \text{or} \quad \mu_k \leq \epsilon.$$

549 Denote \bar{x} the output of iterate x^k , **Iter** the number of running iterations and **Time** the CPU time
 550 of the SPG algorithm by the criterion in (4.1). Examples 4.1 and 4.2 are for the under-determined
 551 linear regression problems. Moreover, Example 4.1 is a typical under-determined linear regression
 552 problem, which shows that the proposed method in this paper can find a global solution with certain
 553 sparsity. The aim of Example 4.2 is to solve a random generated under-determined sparse linear
 554 regression problem, while Example 4.3 is to solve a over-determined censored regression problem.

555 **Example 4.1. (A test example)** We consider the problem in Example 2.1 to verify the
 556 validity of the theoretical results and the efficiency of SPG algorithm. Problem (2.9) is an example
 557 of problem (1.2) with the ℓ_1 loss function given in (1.3), where $m = 1$, $\mathbf{A} = (1 \ 1)$ and $b = 1$.

Let the smoothing function of f be defined by (3.2). Some fixed parameters in the SPG algorithm are given as follows:

$$\underline{\gamma} = \bar{\gamma} = \sqrt{2}, \alpha = 1, \sigma = 0.8, \rho = 1.1, \text{Maxiter} = 10^4, \epsilon = 10^{-3}, \kappa = 1/2, L_f = \sqrt{2}.$$

Let \mathcal{LM} , $\nu - \mathcal{LM}$ and \mathcal{GM} denote the sets of local minimizers, ν -strong local minimizers and global minimizers of (2.9), respectively. When $\nu < \lambda/L_f$,

$$\nu - \mathcal{LM} = \{x : x_1 + x_2 = 1, \nu \leq x_1, x_2 \leq 1\} \cup \{(1, 0)^T, (0, 1)^T, (0, 0)^T\}.$$

558 Set $\mu_0 = 0.1$ and $x^0 = (1, 0.8)^T$. The other parameters and the numerical results are listed in
 559 Table 4.1, where the global minimizers are same for the cases in one line. For problem (2.9), many
 560 different values of λ and the corresponding ν if $\nu < \lambda/L_f$ are given in Table 4.1, which shows that
 561 \bar{x} is always a ν -strong local minimizer and sometimes a global minimizer of (2.9). In particular,
 562 when $\lambda = 0.7$, $\nu = 0.4$ and $x^0 = (1, 0.8)^T$, the SPG algorithm finds a global solution of (2.9).
 563 Moreover, we consider the influence of the values of ν on the SPG algorithm for solving (2.9) in

564 Table 4.1. When $\lambda = 1$, $\bar{\nu}$ as defined in Assumption 2 is 0.7071. From Table 4.1, we find that the
 565 SPG algorithm finds different ν -strong local minimizer for different values of ν satisfying $\nu < \bar{\nu}$.
 566 And it is interesting that when $\nu \geq 0.5$, the SPG algorithm converges to a global minimizer. We
 567 notice that when ν is a lower bound for the global minimizers, it holds that

568 (4.2)
$$\mathcal{GM} \subseteq \nu_1\text{-}\mathcal{LM} \subseteq \nu_2\text{-}\mathcal{LM}, \forall \nu_2 \leq \nu_1 \leq \nu.$$

569 Hence when ν is a lower bound for the global minimizers, the larger ν is likely to let the SPG
 570 algorithm converge to a global minimizer with higher possibility.

571 The updating rule for μ_k in the SPG algorithm is to ensure its global sequence convergence.
 572 How to improve the local convergence rate with the guarantee of global sequence convergence is an
 interesting work for further research.

λ	\mathcal{GM}	ν	Iter	\bar{x}
0.7/0.8/0.9	(1, 0), (0, 1)	0.4/0.5/0.6	18/19/10	(1, 0)/(1, 0)/(1, 0)
1/1/1	(1, 0), (0, 1), (0, 0)	0.7/0.5/0.3	21/11/5	(0, 0)/(1, 0)/(0.6, 0.4)
1.1/1.2/1.3	(0, 0)	0.7/0.9/1	18/17/16	(0, 0)/(0, 0)/(0, 0)

Table 4.1: Numerical results of the SPG algorithm for problem (2.9) with different λ and ν

573
 574 Using the same parameters and initial point, the IRL1 and IRTight algorithms in [43] may
 575 generate

576 (4.3)
$$x^k = \arg \min_{0 \leq x_1, x_2 \leq 1} |x_1 + x_2 - 1|$$

577 for $k \geq 0$ with $x^k \equiv (\alpha, \beta) > 0$ and $\alpha + \beta = 1$. Obviously, x^k is not a global minimizer of (2.9). Hence
 578 almost surely the reweighted algorithms in [43] cannot find a global minimizer (2.9). In fact, at any
 579 point $x^k > 0$, the derivative of $\|x^k\|_0$ is $(0, 0)^T$ and $x^{k+1} = (\alpha, \beta) > 0$ with $\alpha + \beta = 1$ is an optimal
 580 solution of the subproblem $\min_{0 \leq x_1, x_2 \leq 1} |x_1 + x_2 - 1|$ in the algorithms. Hence the SPG algorithm
 581 has better performance than the algorithms in [43] for solving the nonsmooth optimization problem
 582 with cardinality penalty (1.2).

583 **Example 4.2. (Linear regression problem)** Linear regression problem is the most repre-
 584 sentative problem in sparse regression, which has been widely used in information theory [12],
 585 image restoration[5, 10, 41], signal precessing [10, 41] and variable selection [23, 24] problems. Le-
 586 ast square function is the most frequently used loss function due to its convexity and differentiability
 587 [24, 29, 31, 32, 54]. However, the ℓ_1 loss function often owns the stronger outlier-resistant property
 588 than the least square loss function [23]. So, in this example, we consider the following cardinality
 589 penalty problem with ℓ_1 loss function:

590 (4.4)
$$\min_{0 \leq x \leq 10\mathbf{1}_n} \mathcal{F}_{\ell_0}(x) := \frac{1}{m} \|\mathbf{A}x - b\|_1 + \lambda \|x\|_0,$$

591 where $b \in \mathbb{R}^m$ and $A \in \mathbb{R}^{m \times n}$ with $m < n$.

592 **Generating data and setting parameters.** For positive integers m, n and s , we generate the
 593 original signal x^* with $\|x^*\|_0 = s$, sensing matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ and observation $b \in \mathbb{R}^m$ as follows:

594 `index=randperm(n); index=index(1:s); x*=zeros(n,1); B=randn(n,m);`

595 $x^*(\text{index})=\text{unifrnd}(2,10,[s,1]); \mathbf{A}=\text{orth}(\mathbf{B})'; \mathbf{b} = \mathbf{A}*x^*+ 0.01*\text{randn}(\text{size}(\mathbf{b})).$
 596 In the proposed SPG algorithm, we use the smoothing function of f in (3.2) and set the parameters
 597 as below

$$\gamma = \bar{\gamma} = 1, \alpha = 1, \mu_0 = 50, \rho = 1.1, \sigma = 0.9, \text{Maxiter} = 10^4, \kappa = 1/2.$$

598 It is not hard to show that all assumptions in sections 2 and 3 hold. Thus, the sequence $\{x^k\}$
 599 of the SPG algorithm should be convergent to a ν -strong local minimizer of (4.4).

600 Generate A, b and x^* with $m = 80, n = 160$ and $s = 16$, set $\lambda = 18.8$ in (4.4) and $\epsilon = 10^{-3}$ in the
 601 stopping criterion (4.1). We calculate that $L_f = 10.6168$ and define $\nu = 1.77, x^0 = 1.97*\text{ones}(n, 1)$.
 602 The numerical results are shown in Fig. 4.1. Fig. 4.1(a) plots x^* and \bar{x} . From Fig 4.1(a), we see
 603 that the output of x^k is very close to the original generated signal and satisfies the lower bound
 property in (2.8). Fig. 4.1(b) exhibits the convergence of μ^k and $\mathcal{F}(x^k) - \mathcal{F}(\bar{x})$.

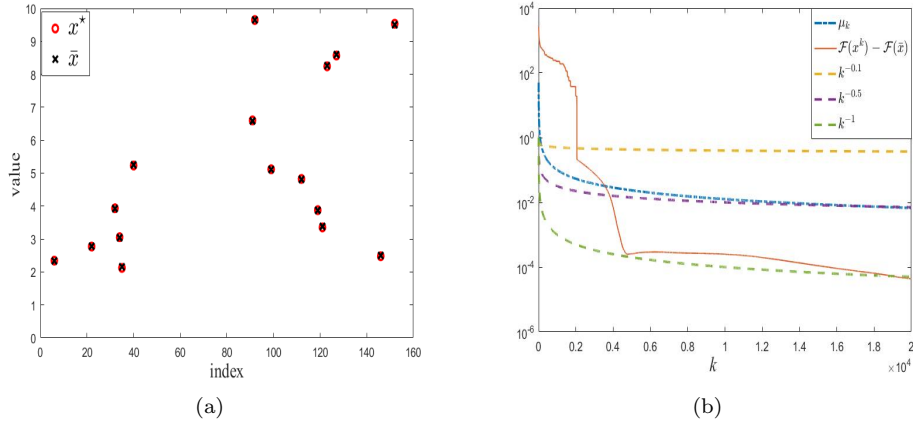


Fig. 4.1: Numerical results of the SPG algorithm for Example 4.2

604

Example 4.3. (Censored regression problem) A typical class of censored regression problem is the linear regression model with left-censoring (or right-censoring) at zero, i.e.

$$\max\{A_i x - c_i, 0\} \approx b_i, \quad i = 1, \dots, m,$$

605 where A_i, b_i and c_i are defined as in (1.4). This class of problems have wide applications in wireless
 606 communication [38], machine learning [21], variable selection[23, 53], economics [9], etc. To solve
 607 it, the loss function is often defined by (1.4), which is nonsmooth for any $p \in [1, 2]$. So the censored
 608 regression problem is a typical class of sparse regression problems with nonsmooth convex loss
 609 functions [53]. Different from the case considered in Example 4.2, we let $m \gg n$ in this example,
 610 which comes from the stochastic optimization models in the portfolio management.

In this example, we let $l = \mathbf{0}$ and $u = \mathbf{1}_n$ in (1.2), and define the loss function f by (1.4) with $c_i = 0, i = 1, \dots, m$ and $p = 1$. The aim of this model is to find a sparse signal $x^* \in [\mathbf{0}, \mathbf{1}_n]$ for the nonlinear system $\max\{\mathbf{A}x^*, 0\} \approx b$ with some unobservable noise, where $\mathbf{A} = (A_1^T, \dots, A_m^T)^T$ and $b = (b_1, \dots, b_m)^T$. We use the relative error (**rel-err**), sparsity regression rate (**spa-rat**) and successful rate (**suc-rat**) to judge the performance of the continuous relaxation model for (1.2)

and the proposed SPG algorithm. Here, the relative error (**rel-err**) and sparsity regression rate (**spa-rat**) of \bar{x} with respect to x^* are defined by

$$\text{rel-err} := \frac{\|\bar{x} - x^*\|}{\|\bar{x}\|}, \quad \text{spa-rat} := \frac{|\mathcal{A}(x^*) \cap \mathcal{A}(\bar{x})|}{\max\{|\mathcal{A}(\bar{x})|, |\mathcal{A}(x^*)|\}},$$

611 where $|\Xi|$ means the cardinality of set Ξ with finite elements. The running regression test is regarded
 612 as a successful one, if the relative error is smaller than 10^{-2} and $\mathcal{A}(\bar{x}) = \mathcal{A}(x^*)$.

613 For the given positive integers m, n and s , the data are generated by

```
614     index = randperm(n); index = index(1:s); x*=zeros(n,1);
615     x*(index)=unifrnd(0,0.9,[s,1]); x*=sign(x)*(abs(x)+0.1)
616     A=randn(m,n); b =max{A*x+0.01*randn(size(b)),0},
```

617 which let x^* satisfies $|x_i^*| \geq 0.1, \forall i \in \mathcal{A}(x^*)$.

618 We use the smoothing function of f in (3.3). Let $L_f = \|\mathbf{A}\|_\infty$ and set $\nu = \min\{\lambda/L_f, 1\}$. Set
 619 $x^0 = 0.1 * \text{ones}(n, 1)$, $\mu_0 = 1$ and $\epsilon = 0.01$. Let the other parameters in the SPG algorithm be the
 620 same as in Example 4.2.

621 For each group of given numbers m, n and s , we generate the codes with 100 independent trials,
 622 and the results displayed in Table 4.2 are the average values for these 100 independent tests. For
 623 each test, regarding the lower bound of the true solution x^* , we run the SPG algorithm for problem
 624 (1.2) with $\lambda := \delta L_f$ for $\delta \in [0.001 : 0.001 : 0.1]$, and report the result with the smallest **rel-err**
 625 for this test. From the displayed results in Table 4.2, we see that the the proposed SPG algorithm
 626 can find the true solution with high possibility, and all the sparsity regression rates are more than
 627 90%. In particular, when $m = 2000$ and $n = 400$, the SPG algorithm can identify almost all the
 628 locations of $\mathcal{A}(x^*)$ when the sparsity levels of x^* are 10%, 20% and 30%. Correctly identifying the
 629 zero and nonzero locations of the true solution is the most important thing in solving the variable
 630 selection and classification problems. When $m = 1000$ and $n = 200$, the values of relative error and
 631 sparsity regression rates by the 100 tests are plotted in Fig. 4.2 for $s = 20, 40$ and 60 , respectively.

m	n	s	Time	Iter	rel-err	$\mathcal{A}(\bar{x})$	spa-rat	suc-rat
1000	200	20	0.612	166	173e-3	19.99	100%	99%
1000	200	40	0.659	178	5.73e-3	39.96	99.7%	89%
1000	200	60	0.708	204	9.12e-3	59.94	92.7%	69%
2000	400	40	2.079	181	1.96e-3	40	100%	96%
2000	400	80	2.686	217	6.93e-3	79.89	99.7%	83%
2000	400	120	3.658	291	9.34e-3	119.91	99.3%	64%

Table 4.2: Average numerical results of the SPG algorithm for the censored regression problem

632

633 **5. Conclusions.** Problem (1.2) includes a class of constrained optimization problems with the
 634 objective function defined by the sum of a nonsmooth convex function and a cardinality function.
 635 Using the capped- ℓ_1 penalty, we propose a continuous relaxation (1.6) of problem (1.2). We prove
 636 that the sets of global minimizers of problems (1.2) and (1.6) are same, and local minimizers of
 637 (1.6) are local minimizers of (1.2) with the lower bound property. Moreover, \bar{x} is a local minimizer
 638 of (1.2) satisfying a desired lower bound property if and only if it is a lifted stationary point of
 639 the continuous relaxation problem (1.6). Though problem (1.6) is a nonsmooth and nonconvex

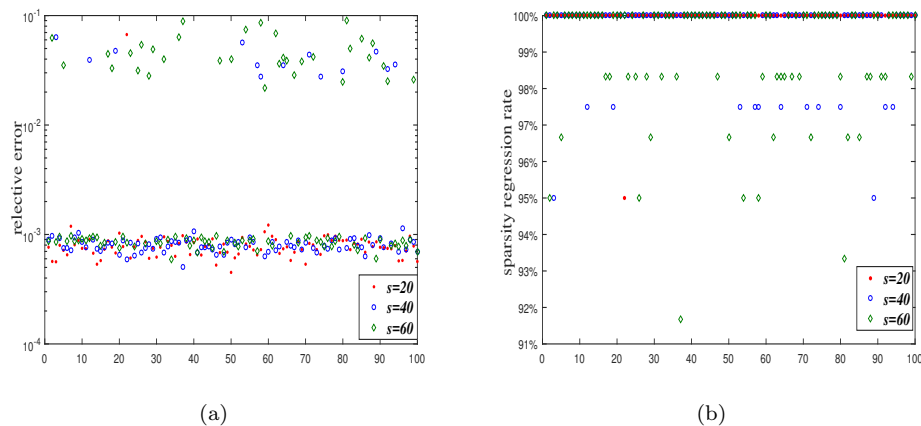


Fig. 4.2: The values of relative error and sparsity regression rate for the 100 tests with $m = 1000$ and $n = 200$

640 optimization problem, its piecewise linear penalty offers us the opportunity to solve it efficiently.
 641 Following this idea, we propose the SPG algorithm based on the smoothing method and the proximal
 642 gradient algorithm to solve problem (1.6), which can find a “good” local minimizer of (1.2) that
 643 satisfies the desired lower bound. The proposed algorithm is simple, whose subproblem has a
 644 closed form solution, and can be run efficiently. We prove the global sequence convergence without
 645 using the K-L condition. Another interesting result is that the local convergence rate of the SPG
 646 algorithm on the objective function value is $o(k^{-\tau})$ with $\tau \in (0, \frac{1}{2})$ and the zero entries of a lifted
 647 stationary point of (1.6) can be identified in finite iterations.

648

REFERENCES

- 649 [1] M. AHN, J.-S. PANG, AND J. XIN, *Difference-of-convex learning: directional stationarity, optimality, and*
 650 *sparsity*, SIAM J. Optim., 27 (2017), pp. 1637–1655.
 651 [2] H. ATTOUCH, J. BOLTE, AND B. F. SVAITER, *Convergence of descent methods for semi-algebraic and tame*
 652 *problems*, Math. Program., 137 (2013), pp. 961–978.
 653 [3] A. BECK AND N. HALLAK, *Proximal mapping for symmetric penalty and sparsity*, SIAM J. Optim., 28 (2018),
 654 pp. 496–527.
 655 [4] A. BECK AND N. HALLAK, *Optimization problems involving group sparsity terms*, Math. Program., 178 (2019),
 656 pp. 39–67.
 657 [5] W. BIAN AND X. CHEN, *Linearly constrained non-Lipschitz optimization for image restoration*, SIAM J. Imaging
 658 Sci., 8 (2015), pp. 2294–2322.
 659 [6] W. BIAN AND X. CHEN, *Optimality and complexity for constrained optimization problems with nonconvex*
 660 *regularization*, Math. Oper. Res., 42 (2017), pp. 1063–1084.
 661 [7] W. BIAN, X. CHEN, AND Y. YE, *Complexity analysis of interior point algorithms for non-Lipschitz and non-*
 662 *convex minimization*, Math. Program., 149 (2015), pp. 301–327.
 663 [8] T. BLUMENSATH AND M. E. DAVIES, *Iterative thresholding for sparse approximations*, J. Fourier Anal. Appl.,
 664 14 (2008), pp. 629–654.
 665 [9] R. BLUNDELL AND J. L. POWELL, *Censored regression quantiles with endogenous regressors*, J. Econometrics,
 666 141 (2007), pp. 65–83.
 667 [10] A. M. BRUCKSTEIN, D. L. DONOHO, AND M. ELAD, *From sparse solutions of systems of equations to sparse*

- 668 *modeling of signals and images*, SIAM Rev., 51 (2009), pp. 34–81.
- 669 [11] P. BUHLMANN, M. KALISCH, AND L. MEIER, *High-dimensional statistics with a view toward applications in*
670 *biology*, Ann. Rev. Stat. Appl., 1 (2014), pp. 255–278.
- 671 [12] E. J. CANDÈS, J. ROMBERG, AND T. TAO, *Robust uncertainty principles: exact signal reconstruction from*
672 *highly incomplete frequency information*, IEEE Trans. Inf. Theory, 52 (2006), pp. 489–509.
- 673 [13] E. J. CANDÈS, M. B. WALKIN, AND S. P. BOYD, *Enhancing sparsity by reweighted ℓ_1 minimization*, J. Fourier
674 Anal. Appl., 14 (2008), pp. 877–905.
- 675 [14] C. CARTIS, N. I. M. GOULD, AND P. L. TOINT, *Adaptive cubic overestimation methods for unconstrained*
676 *optimization. Part I: motivation, convergence and numerical results*, Math. Program., 127 (2011), pp. 245–
677 295.
- 678 [15] R. CHARTRAND AND V. STANEVA, *Restricted isometry properties and nonconvex compressive sensing*, Inverse
679 Probl., 24 (2008), pp. 1–14.
- 680 [16] X. CHEN, *Smoothing methods for nonsmooth, nonconvex minimization*, Math. Program., 134 (2012), pp. 71–99.
- 681 [17] X. CHEN, D. GE, Z. WANG, AND Y. YE, *Complexity of unconstrained ℓ_2 - ℓ_p minimization*, Math. Program., 143
682 (2014), pp. 371–383.
- 683 [18] X. CHEN, L. NIU, AND Y. YUAN, *Optimality conditions and smoothing trust region Newton method for non-*
684 *Lipschitz optimization*, SIAM J. Optim., 23 (2013), pp. 1528–1552.
- 685 [19] E. CHOUZENOUX, A. JEZERSKA, J.-C. PESQUET, AND H. TALBOT, *A majorize-minimize subspace approach for*
686 *ℓ_2 - ℓ_0 image regularization*, SIAM J. Imaging Sci., 6 (2013), pp. 563–591.
- 687 [20] F. H. CLARKE, *Optimization and Nonsmooth Analysis*, SIAM, New York, 1990.
- 688 [21] C. CORTES AND VAPNIK, *Support-vector networks*, Mach. Learn., 20 (1995), pp. 273–297.
- 689 [22] D. L. DONOHO, *Compressed sensing*, IEEE Trans. Inf. Theory, 52 (2006), pp. 1289–1306.
- 690 [23] J. FAN AND R. LI, *Variable selection via nonconcave penalized likelihood and its oracle properties*, J. Amer.
691 Statist. Assoc., 9 (2001), pp. 1348–1360.
- 692 [24] J. FAN, L. XUE, AND H. ZOU, *Strong oracle optimization of folded concave penalized estimation*, Ann. Statist.,
693 42 (2014), pp. 819–849.
- 694 [25] S. FOUcart AND M.-J. LAI, *Sparsest solutions of underdetermined linear system via ℓ_q -minimization for*
695 *$0 < q \leq 1$* , Appl. Comput. Harmon. Anal., 26 (2009), pp. 395–407.
- 696 [26] G. M. FUNG AND O. L. MANGASARIAN, *Equivalence of minimal ℓ_0 - and ℓ_p - norm solutions of linear equations,*
697 *inequalities and linear programs for sufficiently small p* , J. Optim. Theory Appl., 153 (2011), pp. 1–10.
- 698 [27] D. GHILLI AND K. KUNISCH, *On monotone and primal-dual active set schemes for ℓ_p -type problems, $p \in (0, 1]$,*
699 *Comput. Optim. Appl.*, 72 (2019), pp. 45–85.
- 700 [28] J. HUANG, J. L. HOROWITZ, AND S. MA, *Asymptotic properties of bridge estimators in sparse high-dimensional*
701 *regression models*, Ann. Statist., 36 (2008), pp. 587–613.
- 702 [29] J. HUANG, Y. JIAO, B. JIN, J. LIU, X. LU, AND C. YANG, *A unified primal dual active set algorithm for*
703 *nonconvex sparse recovery*.
- 704 [30] P. J. HUBER, *Robust Estimation*, Wiley, 1981.
- 705 [31] K. ITO AND K. KARL, *A variational approach to sparsity optimization based on lagrange multiplier theory*,
706 *Inverse Probl.*, 30 (2014), 015001 (23 pages).
- 707 [32] Y. JIAO, B. JIN, AND X. LU, *A primal dual active set with continuation algorithm for the ℓ^0 -regularized*
708 *optimization problem*, Appl. Comput. Harmon. Anal., 39 (2015), pp. 400–426.
- 709 [33] R. KOENKER, *Quantile Regression*, Cambridge Univ. Press, 2005.
- 710 [34] J. KREIMER AND R. Y. RUBINSTEIN, *Nondifferentiable optimization via smooth approximation: general analy-*
711 *tical approach*, Ann. Oper. Res., 39 (1992), pp. 97–119.
- 712 [35] M. LAI AND J. WANG, *An unconstrained ℓ_q minimization with $0 < q \leq 1$ for sparse solution of under-determined*
713 *linear systems*, SIAM J. Optim., 21 (2011), pp. 82–101.
- 714 [36] M. LAI, Y. XU, AND W. YIN, *Improved iteratively reweighted least squares for unconstrained smoothed ℓ_q*
715 *minimization*, SIAM J. Numer. Anal., 51 (2013), pp. 927–957.
- 716 [37] H. A. LE THI, T. PHAM DINH, H. M. LE, AND X. T. VO, *DC approximation approaches for sparse optimization*,
717 *Eur. J. Oper. Res.*, 244 (2015), pp. 26–46.
- 718 [38] Y. LIU, S. MA, Y. DAI, AND S. ZHANG, *A smoothing SQP framework for a class of composite ℓ_q minimization*
719 *over polyhedron*, Math. Program., 158 (2016), pp. 467–500.
- 720 [39] Y. LIU AND Y. WU, *Variable selection via a combination of the ℓ_0 and ℓ_1 penalties*, J. Comput. Graph. Statist.,
721 16 (2007), pp. 4036–4048.
- 722 [40] Z. LU, *Iterative hard thresholding methods for ℓ_0 regularized convex programming*, Math. Program., 147 (2014),
723 pp. 125–154.
- 724 [41] M. NIKOLOVA AND M. K. NG, *Analysis of half-quadratic minimization methods for signal and image recovery*,
725 *SIAM J. Sci. Comput.*, 27 (2005), pp. 937–966.

- 726 [42] J. NOCEDAL AND S. J. WRIGHT, *Numerical Optimization*, Springer, 2006.
- 727 [43] P. OCHS, A. DOSOVITSKIY, T. BROX, AND T. POCK, *On iteratively reweighted algorithms for nonsmooth*
728 *nonconvex optimization in computer vision*, SIAM J. Imaging Sci., 8 (2015), pp. 331–372.
- 729 [44] B. A. OLSHAUSEN AND D. J. FIELD, *Emergence of simple-cell receptive field properties by learning a sparse*
730 *code for natural images*, Nature, 381 (1996), pp. 607–609.
- 731 [45] C. S. ONG AND L. T. H. AN, *Learning sparse classifiers with difference of convex functions algorithms*, Optim.
732 *Method Softw.*, 28 (2013), pp. 830–854.
- 733 [46] J.-S. PANG, M. RAZAVIYAYN, AND A. ALVARADO, *Computing B-stationary points of nonsmooth DC programs*,
734 *Math. Oper. Res.*, 42 (2017), pp. 95–118.
- 735 [47] D. PELEG AND R. MEIR, *A bilinear formulation for vector sparsity optimization*, Signal Process., 88 (2008),
736 pp. 375–389.
- 737 [48] R. T. ROCKAFELLAR, *Convex Analysis*, Princeton University Press, USA, 1972.
- 738 [49] R. T. ROCKAFELLAR AND R. J.-B. WETS, *Variational Analysis*, Springer-Verlag, Berlin, Germany, 1998.
- 739 [50] E. SOUBIES, L. BLANC-FÉRAUD, AND G. AUBERT, *A unified view of exact continuous penalties for ℓ_2 - ℓ_0 mini-*
740 *mization*, SIAM J. Optim., 27 (2017), pp. 2034–2060.
- 741 [51] E. SOUBIES, L. BLANC-FRAUD, AND G. AUBERT, *A continuous exact ℓ_0 penalty (CEL0) for least squares*
742 *regularized problem*, SIAM J. Imaging Sci., 8 (2015), pp. 1607–1639.
- 743 [52] R. TIBSHIRANI, *Regression shrinkage and selection via the Lasso*, J. Roy. Statist. Soc. B, 58 (1996), pp. 267–288.
- 744 [53] L. WANG, Y. WU, AND L. RUNZE, *Quantile regression for analyzing heterogeneity in ultra-high dimension*,
745 *Ann. Stat.*, 107 (2012), pp. 214–222.
- 746 [54] C. ZHANG, *Nearly unbiased variable selection under minimax concave penalty*, Ann. Stat., 38 (2010), pp. 894–
747 942.
- 748 [55] T. ZHANG, *Multi-stage convex relaxation for feature selection*, Bernoulli, 19 (2013), pp. 2277–2293.
- 749 [56] Z. ZHANG, Y. FAN, AND J. LV, *High dimensional thresholded regression and shrinkage effect*, J. R. Stat. Soc.
750 *Ser. B. Sta. Methodol.*, 76 (2014), pp. 627–649.