

1                    **SMOOTHING FAST ITERATIVE HARD THRESHOLDING**  
2                    **ALGORITHM FOR  $\ell_0$  REGULARIZED NONSMOOTH CONVEX**  
3                    **REGRESSION PROBLEM\***

4                    FAN WU<sup>†</sup>, WEI BIAN<sup>\*‡</sup>, AND XIAOPING XUE<sup>‡</sup>

5                    **Abstract.** We investigate a class of constrained sparse regression problem with cardinality pen-  
6                    alty, where the feasible set is defined by box constraint, and the loss function is convex, but not  
7                    necessarily smooth. First, we put forward a smoothing fast iterative hard thresholding (SFIHT)  
8                    algorithm for solving such optimization problems, which combines smoothing approximations, ex-  
9                    trapolation techniques and iterative hard thresholding methods. The extrapolation coefficients can  
10                    be chosen to satisfy  $\sup_k \beta_k = 1$  in the proposed algorithm. We discuss the convergence behavior of  
11                    the algorithm with different extrapolation coefficients, and give sufficient conditions to ensure that  
12                    any accumulation point of the iterates is a local minimizer of the original cardinality penalized prob-  
13                    lem. In particular, for a class of fixed extrapolation coefficients, we discuss several different update  
14                    rules of the smoothing parameter and obtain the convergence rate of  $O(\ln k/k)$  on the loss and objec-  
15                    tive function values. Second, we consider the case in which the loss function is Lipschitz continuously  
16                    differentiable, and develop a fast iterative hard thresholding (FIHT) algorithm to solve it. We prove  
17                    that the iterates of FIHT converge to a local minimizer of the problem that satisfies a desirable lower  
18                    bound property. Moreover, we show that the convergence rate of loss and objective function values  
19                    are  $o(k^{-2})$ . Finally, some numerical examples are presented to illustrate the theoretical results.

20                    **Key words.** cardinality penalty, smoothing method, accelerated algorithm, extrapolation, con-  
21                    vergence rate, local minimizer

22                    **AMS subject classifications.** 90C30, 65K05, 49J52, 49M37

23                    **1. Introduction.** In this paper, we consider the following minimization prob-  
24                    lem:

25                    (1.1)                    
$$\begin{aligned} \min \quad & F(x) := f(x) + \lambda \|x\|_0 \\ \text{s.t.} \quad & x \in \mathcal{X} := \{x \in \mathbb{R}^n : l \leq x \leq u\}, \end{aligned}$$

26                    for some  $\lambda > 0$ ,  $l \in \overline{\mathbb{R}}_-^n := \{x \in \mathbb{R}^n : -\infty \leq x_i \leq 0, 1 \leq i \leq n\}$ ,  $u \in \overline{\mathbb{R}}_+^n := \{x \in \mathbb{R}^n : 0 \leq x_i \leq +\infty, 1 \leq i \leq n\}$  with  $l < u$ . In (1.1), we call that  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  is the loss  
27                    function to characterize the data fitting, and  $\|x\|_0$  is the penalty function to control  
28                    the sparsity of solutions. Penalty parameter  $\lambda$  is to coordinate the trade-off between  
29                    the data fitting and sparsity. Throughout this paper, we assume that  $f$  in (1.1) is  
30                    convex but not necessarily smooth, and we focus on the case that  $f$  is nonsmooth.  
31                    Such nonsmooth convex regression problems with cardinality penalty arise from

32                    many important applications including compressed sensing [12, 19], variable selection  
33                    [23], signal and image processing [32, 11], pattern recognition [8] and regression [34],  
34                    etc. The purpose of these problems is to find the sparse solutions, most of whose  
35                    elements are zeros. Owing to the existence of  $\ell_0$  function, optimization problem (1.1)  
36                    is typical NP-hard in general. When  $f$  in (1.1) is a smooth convex function, a variety  
37                    of first-order algorithms have been proposed. One important application is the least  
38                    squares problem, i.e.,  $f(x) = \|Ax - b\|^2$ , with sensing matrix  $A \in \mathbb{R}^{m \times n}$  and observa-  
39                    tion vector  $b \in \mathbb{R}^m$ . Greedy methods have been proposed to seek the solutions of the  
40                   

---

\*Corresponding author.

**Funding:** This work is funded by the National Science Foundation of China (Nos: 11871178, 61773136)

<sup>†</sup>School of Mathematics, Harbin Institute of Technology, Harbin, China ([wufanmath@163.com](mailto:wufanmath@163.com)).

<sup>‡</sup>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](mailto:bianweilvse520@163.com), [xiaopingxue@hit.edu.cn](mailto:xiaopingxue@hit.edu.cn)).

41  $\ell_0$  penalized least squares problems in the early stage, such as matching pursuit (MP)  
 42 [26], orthogonal matching pursuit (OMP) [30], subspace pursuit (SP) [16], and so on.  
 43 With the development of compressed sensing, Donoho [19] and Candès, Romberg,  
 44 Tao [12] confirmed the equivalence of the  $\ell_0$  problem and  $\ell_1$  problem when  $A$  satisfies  
 45 some proper conditions. Continuous convex relaxation methods is to replace the  $\ell_0$   
 46 function by a continuous convex function. Even though  $\ell_1$  penalty can be used to find  
 47 the sparse solutions effectively, recent theoretical analysis shows that it often leads  
 48 to an over-penalized problem or a biased estimator. Therefore, there occur various  
 49 continuous nonconvex relaxation functions for  $\ell_0$  function, such as smoothly clipped  
 50 absolute deviation (SCAD) function [20], hard thresholding function [37], capped- $\ell_1$   
 51 function [31], transformed  $\ell_1$  function [29], etc. It is proved that these continuous  
 52 nonconvex penalty functions can not only bring accurate sparse solutions, but also  
 53 reduce the deviation of nonzero elements with respect to the true estimator. Never-  
 54 theless, in [6], it has been proved that finding global minimizers of these nonconvex  
 55 relaxation problems are also NP-hard in general.

56 Blumensath and Davies [9] proposed an iterative hard thresholding (IHT) algo-  
 57 rithm for solving the unconstrained and constrained  $\ell_0$  penalized problems, respec-  
 58 tively. And they proved that the iterates converge to a local minimizer when  $\|A\|_2 < 1$ .  
 59 In [10], they also verified that the IHT algorithm can obtain an approximated solu-  
 60 tion if  $A$  has restricted isometry property. Moreover, Lu and Zhang [25] presented a  
 61 penalty decomposition method for general  $\ell_0$ -penalized and  $\ell_0$ -constrained minimiza-  
 62 tion problem, and established that any accumulation point of the iterates satisfies  
 63 the first-order optimality conditions. Lu [24] also studied an IHT algorithm and its  
 64 variant for solving (1.1) when  $f$  is Lipschitz continuously differentiable.

65 The proximal forward-backward splitting algorithm [13, 15, 17, 21] is a classi-  
 66 cal first-order splitting method, which is also called the proximal gradient algorithm.  
 67 When this method is used to solve the  $\ell_1$  penalized convex regression problem, it is  
 68 often called the iterative shrinkage thresholding algorithm (ISTA). As we know, for  
 69 the case that the loss function is Lipschitz continuously differentiable and convex, the  
 70 convergence rate of the objective function values generated by ISTA is  $O(k^{-1})$ , where  
 71  $k$  is the iteration counter. Based on ISTA and Nesterov's acceleration scheme, Beck  
 72 and Teboulle [4] proposed a fast iterative shrinkage thresholding algorithm (FISTA),  
 73 which not only keeps the simplicity and computation cheapness of ISTA, but also  
 74 improves the convergence rate of the objective function values to  $O(k^{-2})$ . More-  
 75 over, Nesterov [28] independently proposed an accelerated gradient algorithm for the  
 76 same problem with the same convergence rate as FISTA. It's noteworthy that Su,  
 77 Boyd and Candès [33] studied the relationships between a second-order ordinary dif-  
 78 ferential equation and the Nesterov's accelerated gradient method. Inspired by the  
 79 analysis in [33], Attouch and Peypouquet [3] investigated a proximal gradient method  
 80 with extrapolation coefficients  $\beta_k = \frac{k-1}{k+\alpha-1}$  for minimizing the sum of a Lipschitz con-  
 81 tinuously differentiable convex function and a proper closed convex function. In [3],  
 82 it is proved that the convergence rate of the objective function values is  $o(k^{-2})$ , and  
 83 the iterates  $\{x^k\}$  converge to a minimizer of this problem as  $\alpha > 3$ . Recently, Doikov  
 84 and Nesterov [18] presented a new accelerated algorithm for solving such problem, in  
 85 which they used the high-order tensor methods to solve the inner subproblems and  
 86 gave a complexity estimate under some assumptions. For the case with Lipschitz con-  
 87 tinuously differentiable but nonconvex loss function, Wen, Chen and Pong [35] proved  
 88 that the iterates and objective function values generated by the proximal gradient  
 89 algorithm with extrapolation are R-linearly convergent under the error bound condi-  
 90 tion. Later, Adly and Attouch [1] proposed an inertial proximal gradient algorithm

91 with Hessian damping and dry friction to obtain the finite convergence under some  
92 certain conditions.

93 A class of direct methods for solving nonsmooth convex minimization is the sub-  
94 gradient methods. For an  $\epsilon > 0$ , if  $x^\epsilon$  satisfies  $f(x^\epsilon) - \min f \leq \epsilon$ , then  $x^\epsilon$  is called an  
95  $\epsilon$ -approximation solution of  $\min f$ . It has been reported that the complexity of most  
96 subgradient methods for finding an  $\epsilon$ -approximation solution is of the order  $O(\epsilon^{-2})$ .  
97 For a class of nonsmooth functions with max operator, Nesterov [27] gave a smooth  
98 convex function with Lipschitz continuous gradient of factor  $\epsilon^{-1}$  to approximate it,  
99 where  $\epsilon$  is a given and fixed positive parameter. Then, Nesterov [27] proved that  
100 the complexity for finding an  $\epsilon$ -approximation solution of this nonsmooth problem  
101 can be improved to  $O(\epsilon^{-1})$  when the accelerated gradient method is applied to solve  
102 the approximate smooth function. Later, the authors in [22] applied the Nesterov's  
103 smoothing technique to the Nash equilibria problem and gave a first-order method  
104 with the same complexity as [27]. Notably, Chen [14] presented a smoothing gradi-  
105 ent method for solving the constrained nonsmooth nonconvex minimization problem  
106 and demonstrated how to update the smoothing parameter so that the algorithm  
107 converges to a stationary point of the problem. Recently, Bian and Chen [7] uti-  
108 lized an exact continuous relaxation problem to solve optimization problem (1.1) and  
109 presented a smoothing proximal gradient algorithm, whose iterates are globally con-  
110 vergent to a local minimizer of problem (1.1) and convergence rate on the objective  
111 function values is  $o(k^{-\tau})$  with any  $\tau \in (0, 1/2)$ . In [2], the authors used the proximal  
112 regularized inertial Newton algorithm to solve the nonsmooth convex optimization  
113 problem, and proved that the convergence rate of the Moreau envelope values of the  
114 objective function is  $o(k^{-2})$  when the index satisfies an updating rule.

115 Up to now, very few studies investigated accelerated algorithm for solving prob-  
116 lem (1.1) in any systematic way. Inspired by the good performance of the accelerated  
117 algorithm with extrapolation and the smoothing method, we present a smoothing fast  
118 iterative hard thresholding (SFIHT) algorithm for solving problem (1.1). It is worth  
119 emphasizing that the SFIHT algorithm is used to minimize the sum of a nonsmooth  
120 convex function and a discontinuous nonconvex function. The strategy of accelera-  
121 tion is to adopt extrapolation on the iterates. As we know, the larger the range of  
122 the extrapolation coefficients, the better. It is worth noting that the extrapolation  
123 coefficients in our algorithm can satisfy  $\sup_k \beta_k = 1$ . A key technique in the SFIHT  
124 algorithm is that we split the range of the extrapolation coefficients into three cases,  
125 which are divided in line with the relationships among the  $\ell_0$  norms of the newest  
126 three adjacent iterates. Besides, though the subproblem in SFIHT algorithm is a  
127 nonconvex minimization problem, it has a closed-form solution and can be calculated  
128 exactly due to the special structure of  $\ell_0$  norm. So, the proposed SFIHT algorithm  
129 is well-defined. We show that the  $\ell_0$  norms of the iterates will not change after a  
130 finite number of iterations. Then, we discuss the convergence behavior of the SFIHT  
131 algorithm with different extrapolation coefficients. Also, we study the case that the  
132 loss function in problem (1.1) is smooth, and show a derived algorithm of SFIHT  
133 algorithm (called FIHT algorithm) for solving it. We prove that the iterates of FIHT  
134 algorithm converge to a local minimizer of this problem with an important lower  
135 bound property. Moreover, we also substantiate that the convergence rate of FIHT  
136 algorithm for the corresponding objective and loss function values is  $o(k^{-2})$ .

137 **Contents.** The rest of this paper is organized as follows. In Section 2, we first  
138 review some preliminary results on smoothing method, then we present the SFIHT  
139 algorithm for solving nonsmooth nonconvex problem (1.1). Next, we analyze the con-  
140 vergence properties of the proposed algorithm with different extrapolation coefficients

141 for solving (1.1). In Section 3, we focus on solving (1.1) with a smooth convex loss  
 142 function. When the extrapolation coefficients in SFIHT algorithm are appropriately  
 143 fixed, we give a better convergence behaviour of the proposed algorithm for solving  
 144 this kind of problems. In Section 4, we apply the proposed algorithms to some practi-  
 145 cal instances, and show the value of acceleration by extrapolation in solving problem  
 146 (1.1).

147 **Notations.** Throughout this paper, we denote  $\mathbb{N} := \{1, 2, \dots\}$ . Let  $\mathbb{R}^n$  be the  
 148 Euclidean space with inner product  $\langle \cdot, \cdot \rangle$  and corresponding Euclidean norm  $\|\cdot\|$ . For  
 149 vectors  $x, y \in \mathbb{R}^n$ ,  $x \geq y$  means that  $x_i \geq y_i$ ,  $i = 1, 2, \dots, n$ . The  $\ell_1$  norm of vector  
 150  $x \in \mathbb{R}^n$  is denoted by  $\|x\|_1$ , and let  $I(x) := \{i : x_i = 0\}$ . For a matrix  $A \in \mathbb{R}^{m \times n}$ , we  
 151 use  $A^T$ ,  $\lambda_{\max}(A)$  and  $\|A\| = \sqrt{\lambda_{\max}(A^T A)}$  to denote its transpose, largest eigenvalue  
 152 and spectral norm, respectively. Given a nonempty closed convex set  $\Omega \subseteq \mathbb{R}^n$  and a  
 153 vector  $x \in \mathbb{R}^n$ ,  $P_\Omega(x) := \arg \min\{\|x - z\| : z \in \Omega\}$ ,  $N_\Omega(x)$  denotes the normal cone of  
 154  $\Omega$  at  $x$  and  $\Omega_J := \{x \in \Omega : x_j = 0, j \in J\}$  for a given index set  $J \subseteq \{1, 2, \dots, n\}$ .

155 **2. Numerical algorithm and its convergence analysis.** In this section, we  
 156 focus on the case that  $f$  is a nonsmooth convex function. In what follows, we assume  
 157 that  $f$  is level bounded on  $\mathcal{X}$ , i.e., set  $\{x \in \mathcal{X} : f(x) \leq r\}$  is bounded for any  $r \in \mathbb{R}$ ,  
 158 which holds naturally if  $\mathcal{X}$  is bounded. Note that function  $F$  in (1.1) is level bounded  
 159 on  $\mathcal{X}$  if and only if  $f$  is level bounded on  $\mathcal{X}$ .

160 **2.1. Smoothing method and basic properties.** To overcome the nondiffer-  
 161 entiability of loss function  $f$  in (1.1), we use a sequence of continuous differentiable  
 162 functions to approximate  $f$ .

163 **DEFINITION 2.1.** [7] We call  $\tilde{f} : \mathbb{R}^n \times (0, \bar{\mu}] \rightarrow \mathbb{R}$  with  $\bar{\mu} > 0$  a smoothing function  
 164 of the convex function  $f$  on  $\mathcal{X}$ , if  $\tilde{f}(x, \mu)$  satisfies the following conditions:

- 165 (i) for any fixed  $\mu > 0$ ,  $\tilde{f}(\cdot, \mu)$  is continuously differentiable in  $\mathbb{R}^n$ ;  
 166 (ii)  $\lim_{z \rightarrow x, \mu \downarrow 0} \tilde{f}(z, \mu) = f(x)$ ,  $\forall x \in \mathcal{X}$ ;  
 167 (iii)  $\tilde{f}(\cdot, \mu)$  is convex on  $\mathcal{X}$  for any fixed  $\mu > 0$ ;  
 168 (iv)  $\left\{ \lim_{z \rightarrow x, \mu \downarrow 0} \nabla_z \tilde{f}(z, \mu) \right\} \subseteq \partial f(x)$ ,  $\forall x \in \mathcal{X}$ ;  
 169 (v) there exists a positive constant  $\kappa$  such that

$$170 \quad |\tilde{f}(x, \mu_2) - \tilde{f}(x, \mu_1)| \leq \kappa |\mu_1 - \mu_2|, \quad \forall x \in \mathcal{X}, \mu_1, \mu_2 \in (0, \bar{\mu}];$$

- 171 (vi) there exists a constant  $L_{\tilde{f}} > 0$  such that for any  $\mu \in (0, \bar{\mu}]$ ,  $\nabla_x \tilde{f}(\cdot, \mu)$  is Lipschitz  
 172 continuous on  $\mathcal{X}$  with Lipschitz constant  $L_{\tilde{f}} \mu^{-1}$ .

173 For the convenience of description, we provide  $\tilde{f}$  a smoothing function of  $f$  with  
 174 the definition in Definition 2.1 in the following analysis and denote  $\nabla \tilde{f}(x, \mu)$  the  
 175 gradient of  $\tilde{f}(x, \mu)$  with respect to  $x$ . By virtue of Definition 2.1-(v), we see that

$$176 \quad (2.1) \quad |\tilde{f}(x, \mu) - f(x)| \leq \kappa \mu, \quad \forall x \in \mathcal{X}, 0 < \mu \leq \bar{\mu}.$$

177 Smooth approximations for nonsmooth optimization problems have been studied  
 178 for decades. The fundamental of smoothing method we use in this paper is as follows.  
 179 We first approximate loss function  $f$  by a smooth function with fixed smoothing  
 180 parameter  $\mu$ . Then, one find an approximate solution of the following problem

$$181 \quad (2.2) \quad \min_{x \in \mathcal{X}} \tilde{F}(x, \mu) := \tilde{f}(x, \mu) + \lambda \|x\|_0.$$

182 Next, by updating the smoothing parameter  $\mu$ , we can find a local minimizer of  
 183 problem (1.1).

184 According to the definition of  $\tilde{f}$ , for any fixed  $\mu \in (0, \bar{\mu}]$ ,  $L \geq L_{\tilde{f}}$ , it holds that

$$185 \quad (2.3) \quad \tilde{f}(x, \mu) - \tilde{f}(y, \mu) \leq \langle \nabla \tilde{f}(y, \mu), x - y \rangle + \frac{L}{2\mu} \|x - y\|^2, \quad \forall x, y \in \mathcal{X}.$$

186 Using the convexity of function  $f$ , we can prove that  $x^* \in \mathcal{X}$  is a local minimizer of  
187 problem (1.1) if and only if  $x^*$  satisfies  $0 \in [\partial f(x^*)]_i + [N_{\mathcal{X}}(x^*)]_i, \forall i \notin I(x^*)$ , which is  
188 equivalent to

$$189 \quad (2.4) \quad x^* \in \arg \min \{f(x) : x \in \mathcal{X}_{I(x^*)}\}.$$

190 From (2.4), we can easily find that any local minimizer of problem (1.1) has the oracle  
191 property [20].

## 192 2.2. Smoothing fast iterative hard thresholding (SFIHT) algorithm.

193 In this subsection, we combine the smoothing method, extrapolation technique and  
194 iterative hard thresholding algorithm to present a fast scheme for solving problem  
195 (1.1). We name it smoothing fast iterative hard thresholding algorithm and denote it  
196 by SFIHT algorithm for short.

197 In order to find an approximate solution of problem (2.2) with a fixed  $\mu > 0$ , we  
198 introduce an approximation of  $\tilde{F}(x, \mu)$  around the given point  $y$  as follows

$$199 \quad (2.5) \quad Q(x, y, \mu) := \tilde{f}(y, \mu) + \langle \nabla \tilde{f}(y, \mu), x - y \rangle + \frac{L}{2\mu} \|x - y\|^2 + \lambda \|x\|_0$$

200 with a constant  $L > L_{\tilde{f}}$ . Further, we solve the following optimization problem

$$201 \quad (2.6) \quad \min_{x \in \mathcal{X}} Q(x, y, \mu)$$

202 to find an approximate solution of problem (2.2). Although  $\ell_0$  function is nonsmooth  
203 and nonconvex, the objective function and the constraint set in (2.6) are both sepa-  
204 rable with respect to all elements of  $x$ . Using this fact, it has been proved in [24] that  
205 optimization problem (2.6) has the closed-form solution denoted by  $\bar{x}$  and expressed  
206 by, for  $i = 1, 2, \dots, n$ ,

$$207 \quad (2.7) \quad \bar{x}_i = \begin{cases} [P_{\mathcal{X}}(S_L(y, \mu))]_i & \text{if } [S_L(y, \mu)]_i^2 - [q(y, \mu)]_i^2 > \frac{2\lambda\mu}{L}, \\ 0 & \text{if } [S_L(y, \mu)]_i^2 - [q(y, \mu)]_i^2 < \frac{2\lambda\mu}{L}, \\ [P_{\mathcal{X}}(S_L(y, \mu))]_i \text{ or } 0 & \text{otherwise,} \end{cases}$$

208 where  $S_L(y, \mu) := y - \frac{\mu}{L} \nabla \tilde{f}(y, \mu)$ ,  $q(y, \mu) := P_{\mathcal{X}}(S_L(y, \mu)) - S_L(y, \mu)$ . We take problem  
209 (2.6) as the unique subproblem of the proposed SFIHT algorithm. See [Algorithm 2.1](#).  
210 Upon the above fact, we know that SFIHT algorithm is well-defined.

211 In each iteration of the SFIHT algorithm, the accelerated iterative hard thresh-  
212 olding method is used to find an approximate solution of problem (2.2).  $\beta_k$  in Step  
213 1 satisfies  $\beta_k \in \left[0, \sqrt{\frac{\mu_k}{\mu_{k-1}}}\right)$ , which is a basic condition for the convergence analysis  
214 of the SFIHT algorithm. The new iterate depends on the two previous computed  
215 iterates. In order to improve the effect of extrapolation, we divide the extrapolation  
216 coefficients into three cases in the algorithm. We adjust the range of extrapolation  
217 coefficients according to the relationships among the  $\ell_0$  norms of the new computed  
218 iterate and two previous iterates. Step 4 is to update the smoothing parameter  $\mu_k$ ,  
219 which ensures that  $\mu_k$  is decreasing and tends to zero.

**Algorithm 2.1** Smoothing Fast Iterative Hard Thresholding (SFIHT) algorithm

**Initialization:** Take  $x^1 = x^0 \in \mathcal{X}$ ,  $L > L_{\bar{f}}$ ,  $\mu_1 = \mu_0 \in (0, \bar{\mu}]$  and  $\sigma \in (0, 2)$ . Set  $k = 1$ .

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

**Step 1.** Choose  $\beta_k \in \left[0, \sqrt{\frac{\mu_k}{\mu_{k-1}}}\right)$ .

**Step 2.** Compute

$$(2.8) \quad y^k = x^k + \beta_k(x^k - x^{k-1}),$$

$$(2.9) \quad \bar{x}^{k+1} \in \arg \min\{Q(x, y^k, \mu_k) : x \in \mathcal{X}\}.$$

**Step 3.** **(3a)** If  $I(x^{k-1}) = I(x^k) = I(\bar{x}^{k+1})$ , let

$$x^{k+1} = \bar{x}^{k+1}$$

and go to **Step 4**.

**(3b)** Otherwise, choose  $\beta_k \in \left[0, \sqrt{\frac{L-L_{\bar{f}}}{4L} \frac{\mu_k}{\mu_{k-1}}}\right]$ , compute **Step 2** to

obtain  $\bar{x}^{k+1}$ .

**(3b-1)** If  $I(x^k) = I(\bar{x}^{k+1})$ , let

$$x^{k+1} = \bar{x}^{k+1}$$

and go to **Step 4**.

**(3b-2)** Otherwise, choose  $\beta_k \in \left[0, \sqrt{\frac{L-L_{\bar{f}}}{8L-4L_{\bar{f}}} \frac{\mu_k}{\mu_{k-1}}}\right]$ , compute **Step**

**2** to obtain  $\bar{x}^{k+1}$  and set  $x^{k+1} = \bar{x}^{k+1}$ .

**Step 4.** Set

$$\mu_{k+1} = \frac{\mu_0}{(k+2)^\sigma}.$$

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

**end while**

**Output**  $x^k$ ,  $\mu_k$  and  $\beta_k$ .

220 **2.3. Convergence analysis.** Let  $\{x^k\}$ ,  $\{y^k\}$  and  $\{\mu_k\}$  be the output iterates  
 221 generated by the SFIHT algorithm. For  $\kappa > 0$  in [Definition 2.1-\(v\)](#), we introduce the  
 222 following sequence

$$223 \quad H(x^k, \mu_k, \tau_k) := \tilde{F}(x^k, \mu_k) + \kappa\mu_k + \tau_k \|x^k - x^{k-1}\|^2,$$

224 where  $\tau_k > 0$ . Specially, we give a way of choosing  $\tau_k$  as follows. For all  $k \in \mathbb{N}$ ,

$$225 \quad (2.10) \quad \tau_k = \begin{cases} \frac{L}{4}\mu_{k-1}^{-1} + \frac{L}{4}\beta_k^2\mu_k^{-1} & \text{if } I(x^{k-1}) = I(x^k) = I(x^{k+1}), \\ \frac{L-L_{\bar{f}}}{8}\mu_{k-1}^{-1} & \text{otherwise.} \end{cases}$$

226 In the following, we will use the above  $\tau_k$  to analyze the convergence of the SFIHT  
 227 algorithm, which is the key point for all the following analysis.

228 First of all, we show that sequence  $\{\tau_k\}$  guarantees that sequence  $H(x^k, \mu_k, \tau_k)$   
 229 is nonincreasing and convergent, i.e., it can be used as an energy function of the  
 230 SFIHT algorithm. To do this, we first state some basic properties and analyze the  
 231 relationships between  $H(x^k, \mu_k, \tau_k)$  and  $H(x^{k+1}, \mu_{k+1}, \tau_{k+1})$ .

232 LEMMA 2.2. *The following statements hold.*

- 233 (i) For every  $k \in \mathbb{N}$ ,  $x^k \in \mathcal{X}$ ,  $\{\mu_k\}$  is monotone decreasing and  $\lim_{k \rightarrow \infty} \mu_k = 0$ .  
 234 (ii) When  $I(x^k) \neq I(x^{k+1})$ , we have

$$\begin{aligned} & H(x^{k+1}, \mu_{k+1}, \tau_{k+1}) - H(x^k, \mu_k, \tau_k) \\ 235 \quad (2.11) \quad & \leq \left[ \tau_{k+1} - \frac{L - L_{\tilde{f}}}{4\mu_k} \right] \|x^{k+1} - x^k\|^2 + \left[ \frac{2L - L_{\tilde{f}}}{2\mu_k} \beta_k^2 - \tau_k \right] \|x^k - x^{k-1}\|^2. \end{aligned}$$

- 236 (iii) When  $I(x^k) = I(x^{k+1})$ , we have

$$\begin{aligned} & H(x^{k+1}, \mu_{k+1}, \tau_{k+1}) - H(x^k, \mu_k, \tau_k) \\ 237 \quad (2.12) \quad & \leq \left[ \tau_{k+1} - \frac{L}{2\mu_k} \right] \|x^{k+1} - x^k\|^2 + \left[ \frac{L}{2\mu_k} \beta_k^2 - \tau_k \right] \|x^k - x^{k-1}\|^2. \end{aligned}$$

238 *Proof.* (i). By the proposed SFIHT algorithm, it's easy to verify this statement.

239 (ii). Using (2.3) with  $x = x^{k+1}$ ,  $y = y^k$  and  $\mu = \mu_k$ , we obtain

$$240 \quad (2.13) \quad \tilde{f}(x^{k+1}, \mu_k) \leq \tilde{f}(y^k, \mu_k) + \langle \nabla \tilde{f}(y^k, \mu_k), x^{k+1} - y^k \rangle + \frac{L_{\tilde{f}}}{2\mu_k} \|x^{k+1} - y^k\|^2.$$

241 Since  $\tilde{f}(x, \mu)$  is convex with respect to  $x$  for any fixed  $\mu \in (0, \bar{\mu}]$ , it holds that

$$242 \quad (2.14) \quad \tilde{f}(y^k, \mu_k) + \langle \nabla \tilde{f}(y^k, \mu_k), x - y^k \rangle \leq \tilde{f}(x, \mu_k), \quad \forall x \in \mathcal{X}.$$

243 According to (2.13), for any  $x \in \mathcal{X}$ ,  $L > L_{\tilde{f}}$ , we have

$$\begin{aligned} & \tilde{F}(x^{k+1}, \mu_k) = \tilde{f}(x^{k+1}, \mu_k) + \lambda \|x^{k+1}\|_0 \\ & \leq \tilde{f}(y^k, \mu_k) + \langle \nabla \tilde{f}(y^k, \mu_k), x^{k+1} - y^k \rangle + \frac{L}{2\mu_k} \|x^{k+1} - y^k\|^2 + \lambda \|x^{k+1}\|_0 \\ 244 \quad & + \frac{L_{\tilde{f}} - L}{2\mu_k} \|x^{k+1} - y^k\|^2. \end{aligned}$$

245 This, combined with the definition of  $x^{k+1}$ , yields that for  $x^k \in \mathcal{X}$ ,

$$\begin{aligned} & \tilde{F}(x^{k+1}, \mu_k) \leq \tilde{f}(y^k, \mu_k) + \langle \nabla \tilde{f}(y^k, \mu_k), x^k - y^k \rangle + \frac{L}{2\mu_k} \|x^k - y^k\|^2 + \lambda \|x^k\|_0 \\ & + \frac{L_{\tilde{f}} - L}{2\mu_k} \|x^{k+1} - y^k\|^2 \\ 246 \quad & \leq \tilde{F}(x^k, \mu_k) + \frac{L}{2\mu_k} \|x^k - y^k\|^2 + \frac{L_{\tilde{f}} - L}{2\mu_k} \|x^{k+1} - y^k\|^2, \end{aligned}$$

247 where the last inequality holds by (2.14). The above inequality together with  $y^k =$   
 248  $x^k + \beta_k(x^k - x^{k-1})$ , gives

$$\begin{aligned} & \tilde{F}(x^{k+1}, \mu_k) \leq \tilde{F}(x^k, \mu_k) + \frac{L_{\tilde{f}}}{2\mu_k} \beta_k^2 \|x^k - x^{k-1}\|^2 + \frac{L_{\tilde{f}} - L}{2\mu_k} \|x^{k+1} - x^k\|^2 \\ 249 \quad (2.15) \quad & + \frac{L - L_{\tilde{f}}}{\mu_k} \langle x^{k+1} - x^k, \beta_k(x^k - x^{k-1}) \rangle. \end{aligned}$$

250 Using the algebraic inequality, it holds that

$$251 \quad \frac{L - L_{\tilde{f}}}{\mu_k} \langle x^{k+1} - x^k, \beta_k(x^k - x^{k-1}) \rangle \leq \frac{L - L_{\tilde{f}}}{4\mu_k} \|x^{k+1} - x^k\|^2 + \frac{L - L_{\tilde{f}}}{\mu_k} \beta_k^2 \|x^k - x^{k-1}\|^2.$$

252 Combining this with (2.15), we obtain

$$253 \quad (2.16) \quad \tilde{F}(x^{k+1}, \mu_k) \leq \tilde{F}(x^k, \mu_k) - \frac{L - L_{\tilde{f}}}{4\mu_k} \|x^{k+1} - x^k\|^2 + \frac{2L - L_{\tilde{f}}}{2\mu_k} \beta_k^2 \|x^k - x^{k-1}\|^2.$$

254 By Definition 2.1-(v) and the monotone decreasing of  $\{\mu_k\}$ , we have

$$255 \quad (2.17) \quad \tilde{F}(x^{k+1}, \mu_{k+1}) + \kappa\mu_{k+1} - \kappa\mu_k \leq \tilde{F}(x^{k+1}, \mu_k),$$

256 plugging (2.17) in (2.16), then we get

$$257 \quad (2.18) \quad \begin{aligned} & \tilde{F}(x^{k+1}, \mu_{k+1}) + \kappa\mu_{k+1} \\ & \leq \tilde{F}(x^k, \mu_k) + \kappa\mu_k - \frac{L - L_{\tilde{f}}}{4\mu_k} \|x^{k+1} - x^k\|^2 + \frac{2L - L_{\tilde{f}}}{2\mu_k} \beta_k^2 \|x^k - x^{k-1}\|^2. \end{aligned}$$

258 Using (2.18) and the definition of  $H(x^k, \mu_k, \tau_k)$ , we obtain (2.11) immediately.  
(iii). Let

$$G(x, y, \mu) := \tilde{f}(y, \mu) + \langle \nabla \tilde{f}(y, \mu), x - y \rangle + \frac{L}{2\mu} \|x - y\|^2.$$

259 For each fixed  $y$  and  $\mu$ ,  $G(x, y, \mu)$  is differentiable and strongly convex with respect  
260 to  $x$  with modulus  $L\mu^{-1}$ . Hence, for arbitrary  $x \in \mathcal{X}$ , we have

$$261 \quad \begin{aligned} & \langle \nabla_x G(x^{k+1}, y^k, \mu_k), x - x^{k+1} \rangle \\ & = \langle \nabla \tilde{f}(y^k, \mu_k), x - x^{k+1} \rangle + \frac{L}{2\mu_k} \|x - y^k\|^2 - \frac{L}{2\mu_k} \|x^{k+1} - x\|^2 - \frac{L}{2\mu_k} \|x^{k+1} - y^k\|^2. \end{aligned}$$

262 When  $I(x^k) = I(x^{k+1})$ , in view of the construction of the SFIHT algorithm, we find

$$263 \quad (2.19) \quad x^{k+1} = \arg \min_{x \in \mathcal{X}} \{G(x, y^k, \mu_k) + \lambda \|x^k\|_0\}.$$

264 By the above fact, we obtain that  $\langle \nabla_x G(x^{k+1}, y^k, \mu_k), x - x^{k+1} \rangle \geq 0$ ,  $\forall x \in \mathcal{X}$ , which  
265 indicates that, for all  $x \in \mathcal{X}$ ,

$$266 \quad (2.20) \quad 0 \leq \langle \nabla \tilde{f}(y^k, \mu_k), x - x^{k+1} \rangle - \frac{L}{2\mu_k} \|x^{k+1} - y^k\|^2 + \frac{L}{2\mu_k} \|x - y^k\|^2 - \frac{L}{2\mu_k} \|x^{k+1} - x\|^2.$$

267 Summing up (2.13) and (2.20), and by  $L > L_{\tilde{f}}$ , we have

$$268 \quad (2.21) \quad \begin{aligned} & \tilde{f}(x^{k+1}, \mu_k) \\ & \leq \tilde{f}(y^k, \mu_k) + \langle \nabla \tilde{f}(y^k, \mu_k), x - y^k \rangle + \frac{L}{2\mu_k} \|x - y^k\|^2 - \frac{L}{2\mu_k} \|x^{k+1} - x\|^2. \end{aligned}$$

269 Substituting (2.14) into the right side of (2.21), one has

$$270 \quad (2.22) \quad \tilde{f}(x^{k+1}, \mu_k) \leq \tilde{f}(x, \mu_k) + \frac{L}{2\mu_k} \|x - y^k\|^2 - \frac{L}{2\mu_k} \|x^{k+1} - x\|^2.$$



271 Letting  $x = x^k$  in (2.22), and using the definition of  $y^k$ ,  $\|x^k\|_0 = \|x^{k+1}\|_0$  and (2.17),  
 272 we obtain

$$\begin{aligned} & \tilde{F}(x^{k+1}, \mu_{k+1}) + \kappa\mu_{k+1} \\ 273 \quad (2.23) \quad & \leq \tilde{F}(x^k, \mu_k) + \kappa\mu_k + \frac{L}{2\mu_k}\beta_k^2\|x^k - x^{k-1}\|^2 - \frac{L}{2\mu_k}\|x^{k+1} - x^k\|^2. \end{aligned}$$

274 By the definition of  $H(x^k, \mu_k, \tau_k)$  and (2.23), we see the statement in item (iii).  $\square$

275 For simplicity, we define the notations

$$276 \quad (2.24) \quad \gamma := \min \left\{ \frac{L}{4}, \frac{L - L_{\bar{f}}}{8} \right\}, \quad \mathcal{K} := \{k : I(x^k) \neq I(x^{k+1})\}$$

277 and

$$278 \quad (2.25) \quad \nu := \min \{l_i^2\mu_0^{-1}, u_j^2\mu_0^{-1}, 2\lambda L^{-1} : l_i \neq 0, u_j \neq 0, i, j \in \{1, 2, \dots, n\}\}.$$

279 LEMMA 2.3. *The following statements hold.*

- 280 (i) When  $k \in \mathcal{K}$ ,  $\frac{2L-L_{\bar{f}}}{2\mu_k}\beta_k^2 \leq \tau_k$ ; otherwise,  $\frac{L}{2\mu_k}\beta_k^2 \leq \tau_k$ .  
 281 (ii)  $\{H(x^k, \mu_k, \tau_k)\}$  is nonincreasing and convergent, i.e.,  $\lim_{k \rightarrow \infty} H(x^k, \mu_k, \tau_k) =$   
 282  $H_\infty < \infty$ .

283 *Proof.* (i). For  $k \in \mathcal{K}$ , there must be  $\beta_k^2 \leq \frac{L-L_{\bar{f}}}{8L-4L_{\bar{f}}}\frac{\mu_k}{\mu_{k-1}}$ , which means that  $\tau_k =$

$$284 \quad \frac{L-L_{\bar{f}}}{8}\mu_{k-1}^{-1} \geq \frac{2L-L_{\bar{f}}}{2}\beta_k^2\mu_k^{-1}.$$

285  $k \notin \mathcal{K}$  implies  $k \in \mathcal{N}_1 := \{k : I(x^{k-1}) = I(x^k) = I(x^{k+1})\}$  or  $k \in \mathcal{N}_2 := \{k :$   
 286  $I(x^{k-1}) \neq I(x^k) = I(x^{k+1})\}$ . From the SFIHT algorithm, we know  $\beta_k^2 < \frac{\mu_k}{\mu_{k-1}}$  for  
 287  $k \in \mathcal{N}_1$  and  $\beta_k^2 \leq \frac{L-L_{\bar{f}}}{4L}\frac{\mu_k}{\mu_{k-1}}$  for  $k \in \mathcal{N}_2$ . Then, for  $k \in \mathcal{N}_1$ ,  $\mu_{k-1}^{-1} > \beta_k^2\mu_k^{-1}$ , and by  
 288 (2.10), we have

$$289 \quad \tau_k = \frac{L}{4}\mu_{k-1}^{-1} + \frac{L}{4}\beta_k^2\mu_k^{-1} > \frac{L}{4}\beta_k^2\mu_k^{-1} + \frac{L}{4}\beta_k^2\mu_k^{-1} = \frac{L}{2}\beta_k^2\mu_k^{-1};$$

290 for  $k \in \mathcal{N}_2$ ,  $\mu_{k-1}^{-1} \geq \frac{4L}{L-L_{\bar{f}}}\beta_k^2\mu_k^{-1}$ , and by (2.10), we have  $\tau_k = \frac{L-L_{\bar{f}}}{8}\mu_{k-1}^{-1} \geq \frac{L}{2}\beta_k^2\mu_k^{-1}$ .

291 Then, we establish result (i).

292 (ii). We first analyze the nonincreasing of  $\{H(x^k, \mu_k, \tau_k)\}$ , and divide the proof  
 293 into two cases as follows.

294 Case 1. If  $k \in \mathcal{K}$ , by (2.10), we have  $\tau_{k+1} = \frac{L-L_{\bar{f}}}{8}\mu_k^{-1}$ . This together with (2.11)  
 295 and (i) of this lemma, it holds that

$$296 \quad (2.26) \quad H(x^{k+1}, \mu_{k+1}, \tau_{k+1}) - H(x^k, \mu_k, \tau_k) \leq -\frac{L-L_{\bar{f}}}{8}\mu_k^{-1}\|x^{k+1} - x^k\|^2, \quad \forall k \in \mathcal{K}.$$

297 Case 2. If  $k \notin \mathcal{K}$ , we know that  $k \in \{I(x^k) = I(x^{k+1}) = I(x^{k+2})\}$  or  $k \in \{I(x^k) =$   
 298  $I(x^{k+1}) \neq I(x^{k+2})\}$ .

299 When  $k \in \{I(x^k) = I(x^{k+1}) = I(x^{k+2})\}$ ,  $\tau_{k+1} = \frac{L}{4}\mu_k^{-1} + \frac{L}{4}\beta_{k+1}^2\mu_{k+1}^{-1}$ . By (2.12)  
 300 and statement (i) in this lemma, we obtain

$$\begin{aligned} & H(x^{k+1}, \mu_{k+1}, \tau_{k+1}) - H(x^k, \mu_k, \tau_k) \\ 301 \quad (2.27) \quad & \leq \left[ \tau_{k+1} - \frac{L}{2\mu_k} \right] \|x^{k+1} - x^k\|^2 = -\frac{L}{4} \left[ 1 - \beta_{k+1}^2 \frac{\mu_k}{\mu_{k+1}} \right] \mu_k^{-1} \|x^{k+1} - x^k\|^2. \end{aligned}$$

302 By  $\beta_{k+1} \in \left[0, \sqrt{\frac{\mu_{k+1}}{\mu_k}}\right)$  in this case, we obtain  $H(x^{k+1}, \mu_{k+1}, \tau_{k+1}) \leq H(x^k, \mu_k, \tau_k)$ .

303 When  $k \in \{I(x^k) = I(x^{k+1}) \neq I(x^{k+2})\}$ , by (2.10), we have  $\tau_{k+1} = \frac{L-L_{\bar{f}}}{8}\mu_k^{-1}$ .

304 According to (2.12) and result (i) of this lemma, it yields that

$$305 \quad (2.28) \quad H(x^{k+1}, \mu_{k+1}, \tau_{k+1}) - H(x^k, \mu_k, \tau_k) \leq -\frac{3L + L_{\bar{f}}}{8}\mu_k^{-1}\|x^{k+1} - x^k\|^2.$$

306 Hence,  $\{H(x^k, \mu_k, \tau_k)\}$  is nonincreasing. Since  $f$  is bounded from below on  $\mathcal{X}$   
 307 and  $\{x^k\} \subseteq \mathcal{X}$ ,  $\{H(x^k, \mu_k, \tau_k)\}$  is also bounded from below on  $\mathcal{X}$ . This together with  
 308 the nonincreasing of  $\{H(x^k, \mu_k, \tau_k)\}$  implies that  $\{H(x^k, \mu_k, \tau_k)\}$  is convergent.  $\square$

309 The next lemma explores the boundedness of sequence  $\{x^k\}$ , and gives an estimate  
 310 on  $\{x^k\}$  and  $\{\mu_k\}$ , which lays a foundation for the analysis of  $I(x^k)$ .

311 LEMMA 2.4. *The following statements hold:*

312 (i) *sequence  $\{x^k\}$  is bounded;*

313 (ii) *for any  $k \in \mathcal{K}$ ,  $\|x^{k+1} - x^k\|^2 \geq \nu\mu_k$ , where  $\nu$  is defined as in (2.25).*

314 *Proof.* (i). By (2.1) and the definition of  $\{H(x^k, \mu_k, \tau_k)\}$ , we have

$$315 \quad F(x^k) \leq \tilde{F}(x^k, \mu_k) + \kappa\mu_k \leq H(x^k, \mu_k, \tau_k) \leq H(x^1, \mu_1, \tau_1) = \tilde{F}(x^1, \mu_1) + \kappa\mu_1,$$

316 which together with the level boundedness of  $F$  on  $\mathcal{X}$  gives result (i).

(ii). For any  $k \in \mathcal{K}$ , from (2.7), there exists some  $i \in \{1, 2, \dots, n\}$  such that

$$x_i^k = [P_{\mathcal{X}}(S_L(y^{k-1}, \mu_{k-1}))]_i \neq 0, \quad x_i^{k+1} = 0$$

or

$$x_i^k = 0, \quad x_i^{k+1} = [P_{\mathcal{X}}(S_L(y^k, \mu_k))]_i \neq 0.$$

317 When  $x_i^k = [P_{\mathcal{X}}(S_L(y^{k-1}, \mu_{k-1}))]_i \neq 0$  and  $x_i^{k+1} = 0$ , by the definition of  $P_{\mathcal{X}}(\cdot)$ , there  
 318 exist three cases: (a)  $x_i^k = l_i < 0$ ; (b)  $x_i^k = u_i > 0$ ; (c)  $x_i^k = [S_L(y^{k-1}, \mu_{k-1})]_i$ . For  
 319 case (a), we have

$$320 \quad (2.29) \quad \|x^{k+1} - x^k\|^2 \geq |x_i^{k+1} - x_i^k|^2 = |x_i^k|^2 = |l_i|^2 \geq \nu\mu_0.$$

321 For case (b), similar to the analysis in case (a), we see that (2.29) also holds. For case  
 322 (c), according to (2.7), we have  $[S_L(y^{k-1}, \mu_{k-1})]_i^2 \geq 2\lambda L^{-1}\mu_{k-1}$ . Hence,

$$323 \quad \|x^{k+1} - x^k\|^2 \geq |x_i^{k+1} - x_i^k|^2 = |x_i^k|^2 = [S_L(y^{k-1}, \mu_{k-1})]_i^2 \geq \nu\mu_{k-1}.$$

324 These together with the decreasing of  $\{\mu_k\}$ , yields that

$$325 \quad (2.30) \quad \|x^{k+1} - x^k\|^2 \geq \nu\mu_k.$$

326 When  $x_i^k = 0$  and  $x_i^{k+1} = [P_{\mathcal{X}}(S_L(y^k, \mu_k))]_i \neq 0$ , by the similar arguments, we can  
 327 get the same result as above. Then, we have thus proved this statement.  $\square$

328 Now we show that the support set of  $\{x^k\}$  generated by the SFIHT algorithm will  
 329 no longer change after finitely many iterations, and give some other useful estimates.

330 LEMMA 2.5. *The sequences generated by the SFIHT algorithm own the following*  
 331 *properties:*

332 (i)  *$I(x^k) = \{i : x_i^k = 0\}$  changes finite times at most;*

333 (ii) *the output extrapolation coefficients can be chosen to satisfy  $\sup_k \beta_k = 1$ ;*

334 (iii)  $\sum_{k=1}^{\infty} \left[ 1 - \beta_{k+1}^2 \frac{\mu_k}{\mu_{k+1}} \right] \mu_k^{-1} \|x^{k+1} - x^k\|^2 < \infty$ ;

335 (iv)  $\lim_{k \rightarrow \infty} [f(x^k) + \tau_k \|x^k - x^{k-1}\|^2]$  exists.

336 *Proof.* (i). Recalling  $\mathcal{K} = \{k : I(x^k) \neq I(x^{k+1})\}$ , then, we only need to show that  
 337 set  $\mathcal{K}$  has at most finite elements. We argue it by contradiction and suppose there  
 338 are infinite elements in  $\mathcal{K}$ . This, together with (2.26) and (ii) of Lemma 2.3, we have

$$\begin{aligned}
 & 0 \leq \sum_{k \in \mathcal{K}} \frac{1}{8} (L - L_{\bar{f}}) \mu_k^{-1} \|x^{k+1} - x^k\|^2 \\
 339 \quad (2.31) \quad & \leq \sum_{k=1}^{\infty} [H(x^k, \mu_k, \tau_k) - H(x^{k+1}, \mu_{k+1}, \tau_{k+1})] = H(x^1, \mu_1, \tau_1) - H_{\infty} < \infty.
 \end{aligned}$$

340 On the basis of Lemma 2.4-(ii), we obtain

$$341 \quad \sum_{k \in \mathcal{K}} \frac{1}{8} (L - L_{\bar{f}}) \mu_k^{-1} \|x^{k+1} - x^k\|^2 \geq \sum_{k \in \mathcal{K}} \frac{1}{8} (L - L_{\bar{f}}) \nu = \infty.$$

342 This leads to a contradiction to (2.31). Hence, set  $\mathcal{K}$  has at most finite elements, we  
 343 see further that  $I(x^k) = \{i : x_i^k = 0\}$  changes finite times at most.

344 (ii). In view of result (i) of this lemma, we know that the SFIHT algorithm will  
 345 continue to run (3a) in Step 3 after finite iterations. Then, by the update rule of  $\mu_k$   
 346 in Step 4, the statement in (ii) holds.

347 (iii). According to (2.26), (2.27) and (2.28), and by (2.24), for  $\forall k \in \mathbb{N}$ , we find

$$348 \quad H(x^{k+1}, \mu_{k+1}, \tau_{k+1}) - H(x^k, \mu_k, \tau_k) \leq -\gamma \left( 1 - \beta_{k+1}^2 \frac{\mu_k}{\mu_{k+1}} \right) \mu_k^{-1} \|x^{k+1} - x^k\|^2.$$

349 Summing up the above inequality from 1 to  $\infty$  and by Lemma 2.3-(ii), we obtain

$$\begin{aligned}
 & 0 \leq \sum_{k=1}^{\infty} \gamma \left( 1 - \beta_{k+1}^2 \frac{\mu_k}{\mu_{k+1}} \right) \mu_k^{-1} \|x^{k+1} - x^k\|^2 \\
 350 \quad & \leq \sum_{k=1}^{\infty} [H(x^k, \mu_k, \tau_k) - H(x^{k+1}, \mu_{k+1}, \tau_{k+1})] = H(x^1, \mu_1, \tau_1) - H_{\infty} < \infty.
 \end{aligned}$$

351 In view of the definition of  $\gamma$  in (2.24), we get the desired result (iii).

352 (iv). Taking  $x = x^k$  and  $\mu = \mu_k$  in (2.1), and by direct computation, it yields  
 353 that

$$354 \quad H(x^k, \mu_k, \tau_k) - 2\kappa\mu_k \leq f(x^k) + \lambda \|x^k\|_0 + \tau_k \|x^k - x^{k-1}\|^2 \leq H(x^k, \mu_k, \tau_k).$$

355 Letting  $k$  tend to infinity in the above inequality, along with  $\lim_{k \rightarrow \infty} \mu_k = 0$ , we get

$$356 \quad \lim_{k \rightarrow \infty} f(x^k) + \lambda \|x^k\|_0 + \tau_k \|x^k - x^{k-1}\|^2 = \lim_{k \rightarrow \infty} H(x^k, \mu_k, \tau_k).$$

357 This, combined with the fact that  $\lim_{k \rightarrow \infty} \|x^k\|_0$  exists, we deduce the existence of  $\square$

$$358 \quad \lim_{k \rightarrow \infty} [f(x^k) + \tau_k \|x^k - x^{k-1}\|^2].$$

359 Refs. [24] and [36] also study the IHT algorithm for solving the constrained  $\ell_0$   
 360 penalized convex regression problem modeled by (1.1). In terms of problem, the main  
 361 difference is that the loss functions studied by them are smooth, while it can be  
 362 nonsmooth in this paper. In terms of algorithm, [24] considers the IHT algorithm,  
 363 while both [36] and this paper focus on the IHT algorithm with extrapolation. It's  
 364 worth stressing that the extrapolation coefficients in [36] need satisfy  $\sup_k \beta_k \leq \frac{\sqrt{2}}{2}$ ,  
 365 but the SFIHT algorithm proposed in this paper expands the range of extrapolation  
 366 coefficients in a significant way, which can be seen clearly by the following results on  
 367 convergence rate.

368 *Remark 2.6.* If  $\beta_k \in \left[0, \sqrt{(1 - a\mu_{k-1}) \frac{\mu_k}{\mu_{k-1}}}\right]$  with  $a > 0$  in Step 1, by Lemma 2.5-  
 369 (iii), we obtain  $\sum_{k=1}^{\infty} \|x^{k+1} - x^k\|^2 < \infty$ .

370 Combining Lemma 2.5-(i) and the framework of the SFIHT algorithm, we know  
 371 that the algorithm only runs (3a) in Step 3 after a finite number of iterations, which  
 372 means that the SFIHT algorithm solves subproblem (2.9) only one time and is used  
 373 to solve a nonsmooth convex optimization after finite iterations. Hence, in order to  
 374 improve the convergence behavior of the iterates generated by the SFIHT algorithm,  
 375 we will consider two different choices of  $\beta_k$  in Step 1 of the SFIHT algorithm.

376 **2.3.1. A sufficient condition for the convergence of  $\{x^k\}$ .** In this sub-  
 377 subsection, we will analyze the convergence of the iterates generated by the SFIHT  
 378 algorithm for solving (1.1) when

$$379 \quad (2.32) \quad \sigma \in \left[\frac{1}{2}, 1\right]$$

380 and

$$381 \quad (2.33) \quad \beta_k \in \left[0, \sqrt{\left(1 - \frac{1}{2k^{1-\sigma}}\right) \frac{\mu_k}{\mu_{k-1}}}\right]$$

382 in Step 1. It's easy to find that  $\beta_k < \sqrt{\frac{\mu_k}{\mu_{k-1}}}$ , hence, the previous results are still hold  
 383 in this case. Based on these results, we first give the following estimations.

LEMMA 2.7. *When  $\sigma$  and  $\beta_k$  in Step 1 are chosen as in (2.32) and (2.33), it holds that*

$$\sum_{k=1}^{\infty} \frac{1}{k+1} \mu_k^{-2} \|x^{k+1} - x^k\|^2 < \infty \quad \text{and} \quad \lim_{k \rightarrow \infty} \|x^{k+1} - x^k\| = 0.$$

*Proof.* From Lemma 2.5-(i) and (2.33), there exists  $K \in \mathbb{N}$ , such that

$$\beta_{k+1}^2 \leq \left(1 - \frac{1}{2(k+1)^{1-\sigma}}\right) \frac{\mu_{k+1}}{\mu_k}, \quad \forall k \geq K.$$

384 This together with Lemma 2.5-(iii), we obtain  $\sum_{k=1}^{\infty} \frac{1}{2(k+1)^{1-\sigma}} \mu_k^{-1} \|x^{k+1} - x^k\|^2 < \infty$ .  
 385 According to  $\mu_k = \frac{\mu_0}{(k+1)^\sigma}$ , we see that

$$386 \quad \sum_{k=1}^{\infty} \frac{1}{k+1} \mu_k^{-2} \|x^{k+1} - x^k\|^2 < \infty \quad \text{and} \quad \sum_{k=1}^{\infty} (k+1)^{2\sigma-1} \|x^{k+1} - x^k\|^2 < \infty.$$

387 Thanks to  $\sigma \in [\frac{1}{2}, 1]$ , we get  $\sum_{k=1}^{\infty} \|x^{k+1} - x^k\|^2 < \infty$ . Then, we have got all the  
 388 estimates in this lemma.  $\square$

389 Next, we give another preliminary result.

390 **PROPOSITION 2.8.** *Let  $\{a_k\}$  and  $\{b_k\}$  be nonnegative sequences with  $\sum_{k=1}^{\infty} a_k b_k <$   
 391  $\infty$ . If  $\{a_k\}$  is a nonincreasing sequence and satisfies  $\sum_{k=1}^{\infty} a_k = \infty$ , then there exists a  
 392 subsequence of  $\{b_k\}$ , denoted by  $\{b_{k_i}\}$ , satisfying  $\lim_{i \rightarrow \infty} b_{k_i-1} = 0$  and  $\lim_{i \rightarrow \infty} b_{k_i} =$   
 393  $0$ .*

394 *Proof.* Since sequence  $\{a_k\}$  is nonincreasing, we have

$$395 \quad \sum_{k=2}^{\infty} a_k (b_{k-1} + b_k) \leq \sum_{k=2}^{\infty} a_{k-1} b_{k-1} + \sum_{k=2}^{\infty} a_k b_k < \infty.$$

396 This together with  $\sum_{k=2}^{\infty} a_k = \infty$  implies that there exists a subsequence of  $\{b_{k-1} +$   
 397  $b_k\}$ , denoted by  $\{b_{k_i-1} + b_{k_i}\}$ , such that  $\lim_{i \rightarrow \infty} (b_{k_i-1} + b_{k_i}) = 0$ . Due to the nonneg-  
 398 ativity of  $\{b_k\}$ , it follows that  $\lim_{i \rightarrow \infty} b_{k_i-1} = 0$  and  $\lim_{i \rightarrow \infty} b_{k_i} = 0$ .  $\square$

399 **THEOREM 2.9.** *Any accumulation point of sequence  $\{x^k\}$  generated by the SFIHT  
 400 algorithm is a local minimizer of problem (1.1).*

401 *Proof.* From [Lemma 2.5-\(i\)](#), we know that there exist some  $K \in \mathbb{N}$  and  $I \subseteq$   
 402  $\{1, 2, \dots, n\}$  such that  $I(x^k) = I$  for all  $k \geq K$ . This, combined with the SFIHT  
 403 algorithm, we have  
 (2.34)

$$404 \quad x^{k+1} = \arg \min_{x \in \mathcal{X}_I} \left\{ \tilde{f}(y^k, \mu_k) + \langle \nabla \tilde{f}(y^k, \mu_k), x - y^k \rangle + \frac{L}{2\mu_k} \|x - y^k\|^2 \right\}, \quad \forall k \geq K.$$

405 By [Lemma 2.7](#), we know

$$406 \quad \sum_{k=1}^{\infty} \frac{1}{k+1} [\mu_k^{-2} \|x^{k+1} - x^k\|^2] < \infty.$$

407 Using the above result, [Proposition 2.8](#) with  $a_k = \frac{1}{k+1}$  and  $b_k = \mu_k^{-2} \|x^{k+1} - x^k\|^2$ ,  
 408 there exists a subsequence of  $\{\mu_k^{-2} \|x^{k+1} - x^k\|^2\}$  satisfying

$$409 \quad (2.35) \quad \lim_{i \rightarrow \infty} \mu_{k_i-1}^{-2} \|x^{k_i} - x^{k_i-1}\|^2 = 0 \quad \text{and} \quad \lim_{i \rightarrow \infty} \mu_{k_i}^{-2} \|x^{k_i+1} - x^{k_i}\|^2 = 0.$$

410 We see further that  $\lim_{i \rightarrow \infty} \mu_{k_i-1}^{-1} (x^{k_i} - x^{k_i-1}) = 0$  and  $\lim_{i \rightarrow \infty} \mu_{k_i}^{-1} (x^{k_i+1} - x^{k_i}) = 0$ ,  
 411 which implies

$$412 \quad (2.36) \quad \lim_{i \rightarrow \infty} \frac{x^{k_i+1} - y^{k_i}}{\mu_{k_i}} = \lim_{i \rightarrow \infty} \left[ \frac{x^{k_i+1} - x^{k_i}}{\mu_{k_i}} - \beta_{k_i} \frac{\mu_{k_i-1}}{\mu_{k_i}} \frac{x^{k_i} - x^{k_i-1}}{\mu_{k_i-1}} \right] = 0.$$

413 From [Lemma 2.4-\(i\)](#) and [Lemma 2.5-\(i\)](#), we know that there exists a subsequence of  
 414  $\{x^{k_i}\}$  (also denoted by  $\{x^{k_i}\}$  for simplicity) and  $\bar{x} \in \mathcal{X}_I$  such that  $\lim_{i \rightarrow \infty} x^{k_i} = \bar{x}$ .  
 415 By the triangle inequality, it holds that

$$416 \quad \|x^{k_i+1} - \bar{x}\| \leq \|x^{k_i+1} - x^{k_i}\| + \|x^{k_i} - \bar{x}\|,$$

417 then we immediately deduce that  $\lim_{i \rightarrow \infty} x^{k_i+1} = \bar{x}$  by  $\lim_{i \rightarrow \infty} \|x^{k_i+1} - x^{k_i}\| = 0$   
 418 proved in [Lemma 2.7](#). According to the definitions of  $y^{k_i}$  and  $\beta_{k_i}$ , we can also obtain  
 419  $\lim_{i \rightarrow \infty} y^{k_i} = \bar{x}$ , and by [Definition 2.1-\(iv\)](#), we know that

$$420 \quad (2.37) \quad \left\{ \lim_{i \rightarrow \infty} \nabla \tilde{f}(y^{k_i}, \mu_{k_i}) \right\} \subseteq \partial f(\bar{x}).$$

421 For any  $k_i \geq K$ , recalling the definition of  $x^{k_i+1}$  in (2.34), we have

$$422 \quad (2.38) \quad \left\langle \nabla \tilde{f}(y^{k_i}, \mu_{k_i}) + L\mu_{k_i}^{-1}(x^{k_i+1} - y^{k_i}), x - x^{k_i+1} \right\rangle \geq 0, \quad \forall x \in \mathcal{X}_I, k_i \geq K.$$

423 Letting  $i \rightarrow \infty$  in (2.38), by (2.36), (2.37) and  $\lim_{i \rightarrow \infty} x^{k_i+1} = \bar{x}$ , we notice that there  
424 exists  $\bar{\xi} \in \partial f(\bar{x})$  satisfying

$$425 \quad (2.39) \quad \langle \bar{\xi}, x - \bar{x} \rangle \geq 0, \quad \forall x \in \mathcal{X}_I.$$

Since  $f$  is a convex function and  $\mathcal{X}_I$  is a nonempty closed convex set, (2.39) implies that  $\bar{x}$  is a global minimizer of  $f$  on  $\mathcal{X}_I$ . By (2.35) and  $\lim_{k \rightarrow \infty} \mu_k = 0$ , we have  $\lim_{i \rightarrow \infty} \mu_{k_i-1}^{-1} \|x^{k_i} - x^{k_i-1}\|^2 = 0$ , which together with the choice of  $\tau_k$  in (2.10), we deduce

$$\lim_{i \rightarrow \infty} [f(x^{k_i}) + \tau_{k_i} \|x^{k_i} - x^{k_i-1}\|^2] = \lim_{i \rightarrow \infty} f(x^{k_i}) = f(\bar{x}).$$

This, combined with Lemma 2.5-(iv), we have

$$\lim_{k \rightarrow \infty} [f(x^k) + \tau_k \|x^k - x^{k-1}\|^2] = f(\bar{x}).$$

426 Assume  $\hat{x}$  is an accumulation point of  $\{x^k\}$  with the convergence of subsequence  $\{x^{l_j}\}$ .

427 Then, it holds that

$$428 \quad (2.40) \quad f(\hat{x}) = \lim_{j \rightarrow \infty} f(x^{l_j}) \leq \lim_{j \rightarrow \infty} [f(x^{l_j}) + \tau_{l_j} \|x^{l_j} - x^{l_j-1}\|^2] = f(\bar{x}).$$

429 Since  $\hat{x} \in \mathcal{X}_I$  and  $\bar{x} \in \arg \min_{x \in \mathcal{X}_I} f(x)$ , by (2.40), we obtain  $\hat{x} \in \arg \min_{x \in \mathcal{X}_I} f(x)$ .

430 Hence, any accumulation point of  $\{x^k\}$  is a global minimizer of  $f$  on  $\mathcal{X}_I$ . Equivalently,  
431 any accumulation point of  $\{x^k\}$  is a local minimizer of problem (1.1).  $\square$

Combining the construction and convergence analysis of the SFIHT algorithm, we know that the extrapolation coefficients can be considered into two cases for simplicity, i.e.,

$$\begin{cases} \beta_k \in \left[ 0, \sqrt{\left(1 - \frac{1}{2k^{1-\sigma}}\right) \frac{\mu_k}{\mu_{k-1}}} \right] & \text{if } I(x^{k-1}) = I(x^k) = I(\bar{x}^{k+1}), \\ \beta_k \in \left[ 0, \sqrt{\frac{L - L_{\bar{f}}}{8L - 4L_{\bar{f}}} \frac{\mu_k}{\mu_{k-1}}} \right] & \text{otherwise.} \end{cases}$$

432 for which the theoretical results in this paper are still valid. It's well known that  
433 the condition of the extrapolation parameters is weaker is better and  $\sup_k \beta_k = 1$  is  
434 the most important. Thus, we divide the extrapolation coefficients into three cases,  
435 mainly to expand the range of extrapolation coefficients so as to greatly improve  
436 the convergence performance of the SFIHT algorithm. Though this behavior will  
437 increase the computation of the algorithm, it is surprising by Lemma 2.5-(i) that the  
438 computational increase is finite.

439 **2.3.2. Convergence rate on objection function values.** In this subsection,  
440 we discuss a specific  $\beta_k$  in Step 1 of the SFIHT algorithm. It not only keeps the  
441 validity for finding the local minimizers of problem (1.1), but also obtains the conver-  
442 gence rate of the loss and objective function values. Based on the FISTA algorithm in  
443 [4] and the smoothing method, Bian in [5] proposed a smoothing fast iterative shrink-  
444 age thresholding algorithm (FISTA<sub>S</sub>) for solving the constrained nonsmooth convex

445 optimization problem. We set extrapolation coefficient  $\beta_k$  in Step 1 to be the same  
446 as that in [5], namely, the following recurrence relation:

$$447 \quad (2.41) \quad \begin{cases} t_k = \frac{1 + \sqrt{1 + 4 \left(\frac{\mu_{k-1}}{\mu_k}\right) t_{k-1}^2}}{2}, \\ \beta_k = \frac{t_{k-1} - 1}{t_k}, \end{cases}$$

448 with  $t_0 = 1$ . Then,  $\beta_k < \sqrt{\frac{\mu_k}{\mu_{k-1}}}$ , which satisfies the condition in Step 1. Since the  
449 SFIHT algorithm always runs (3a) in Step 3 after finite iterations, by [5, Theorem 1],  
450 we can directly get the following corollary.

451 **COROLLARY 2.10.** *Let  $\{x^k\}$  be the sequence generated by the SFIHT algorithm,*  
452 *in which the  $\beta_k$  in Step 1 is defined by (2.41). Then, the following statements hold:*

- 453 (i) *any accumulation point of  $\{x^k\}$  is a local minimizer of problem (1.1);*  
454 (ii)  *$\lim_{k \rightarrow \infty} F(x^k) = F_\infty < \infty$  exists and*

$$455 \quad F(x^{k+1}) - F_\infty = \begin{cases} O(k^{-\sigma}) & \text{if } 0 < \sigma < 1, \\ O(k^{-1} \ln k) & \text{if } \sigma = 1, \\ O(k^{\sigma-2}) & \text{if } 1 < \sigma < 2. \end{cases}$$

456 Since the  $\ell_0$  norms of the iterates change only a finite number of times, the above  
457 convergence rate also holds for loss function  $f$ . From the estimation in [Corollary 2.10](#),  
458 we can see that the convergence rate on the objective function values generated by  
459 the SFIHT algorithm greatly improves the rate given by the algorithm in [7].

460 **3. A case on smooth loss function.** In this section, we discuss the case that  
461 loss function  $f$  in problem (1.1) is Lipschitz continuously differentiable. In summary,  
462 throughout this section, we require  $f$  in (1.1) to satisfy the following assumptions

- $$463 \quad \begin{cases} \bullet f \text{ is a smooth convex function on } \mathcal{X}; \\ \bullet f \text{ is level bounded on } \mathcal{X}; \\ \bullet \nabla f \text{ is Lipschitz continuous with Lipschitz constant } L_f. \end{cases}$$

In the context of this problem, we select a special extrapolation coefficient  $\beta_k$  in the  
Step 1 of the SFIHT algorithm. Since the smoothing method is not needed in this case,  
we call it fast iterative hard thresholding (FIHT) algorithm. Please see [Algorithm 3.1](#).  
For the sake of brevity, we also define an approximation of  $F$  at a given point  $y$  as  
follows:

$$Q(x, y) := f(y) + \langle \nabla f(y), x - y \rangle + \frac{L}{2} \|x - y\|^2 + \lambda \|x\|_0,$$

464 where  $L > L_f$ .

465 Let  $\{x^k\}$  and  $\{y^k\}$  be the iterates generated by the FIHT algorithm. Similarly, the  
466 subproblem in Step 2 of FIHT algorithm has a closed-form solution. Hence, the FIHT  
467 algorithm is also well-defined. The following lemma shows a lower bound property of  
468  $\{x^k\}$ , which can be easily obtained by [36, Lemma 3.2], so we omit its proof here.

469 **LEMMA 3.1.** *The following statements hold.*

- 470 (i) *When  $x_j^k \neq 0$  for some  $j \in \{1, 2, \dots, n\}$ , it holds that*

$$471 \quad (3.1) \quad |x_j^k| \geq \delta := \min_{i=1,2,\dots,n} \delta_i > 0,$$

**Algorithm 3.1** Fast Iterative Hard Thresholding (FIHT) algorithm

**Initialization:** Take  $x^1 = x^0 \in \mathcal{X}$ ,  $\alpha > 0$  and  $L > L_f$ . Set  $k = 1$ .

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

**Step 1.** Take  $\beta_k = \frac{k-1}{k+\alpha-1}$ .

**Step 2.** Compute

$$\begin{aligned} y^k &= x^k + \beta_k(x^k - x^{k-1}), \\ \bar{x}^{k+1} &\in \arg \min\{Q(x, y^k) : x \in \mathcal{X}\}. \end{aligned}$$

**Step 3.** **(3a)** If  $I(x^{k-1}) = I(x^k) = I(\bar{x}^{k+1})$ , let

$$x^{k+1} = \bar{x}^{k+1},$$

increment  $k$  by one and return to **Step 1**.

**(3b)** Otherwise, choose  $\beta_k \in \left[0, \sqrt{\frac{L-L_f}{4L}}\right]$ , compute **Step 2** to obtain  $\bar{x}^{k+1}$ .

**(3b-1)** If  $I(x^k) = I(\bar{x}^{k+1})$ , let

$$x^{k+1} = \bar{x}^{k+1},$$

increment  $k$  by one and return to **Step 1**.

**(3b-2)** Otherwise, choose  $\beta_k \in \left[0, \sqrt{\frac{L-L_f}{8L-4L_f}}\right]$ , compute **Step 2** to obtain  $\bar{x}^{k+1}$  and let  $x^{k+1} = \bar{x}^{k+1}$ .

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

**end while**

**Output**  $x^k$ ,  $\mu_k$  and  $\beta_k$ .

472 *where*

$$473 \quad \delta_i = \begin{cases} \min\left(u_i, \sqrt{2\lambda/L}\right) & \text{if } l_i = 0, \\ \min\left(-l_i, \sqrt{2\lambda/L}\right) & \text{if } u_i = 0, \\ \min\left(-l_i, u_i, \sqrt{2\lambda/L}\right) & \text{otherwise.} \end{cases}$$

474 **(ii)** For every  $k \in \mathbb{N}$ , if  $I(x^k) \neq I(x^{k+1})$ , then  $\|x^{k+1} - x^k\| \geq \delta$ .

475 We begin the convergence analysis of the FIHT algorithm by defining the following  
476 important auxiliary sequence

$$477 \quad W(x^k, \zeta_k) := F(x^k) + \zeta_k \|x^k - x^{k-1}\|^2,$$

478 where  $\zeta_k > 0$ . Again, for all  $k \in \mathbb{N}$ , we give a choice of  $\zeta_k$  as follows,

$$479 \quad (3.2) \quad \zeta_k = \begin{cases} \frac{L}{4}(1 + \beta_k^2) & \text{if } I(x^{k-1}) = I(x^k) = I(x^{k+1}), \\ \frac{L - L_f}{8} & \text{otherwise.} \end{cases}$$

480 By a similar analysis to sequence  $H(x^k, \mu_k, \tau_k)$  in [Section 2](#), we can get some basic



481 results on the sequence  $W(x^k, \zeta_k)$  for the FIHT algorithm. For easy of reading, we  
 482 only list them in the following lemma, but omit their proofs.

483 **LEMMA 3.2.** *The following properties are satisfied:*

- 484 (i) for every  $k \in \mathbb{N}$ ,  $x^k \in \mathcal{X}$ ;  
 485 (ii) when  $I(x^k) \neq I(x^{k+1})$ , we have

$$486 \quad \begin{aligned} & W(x^{k+1}, \zeta_{k+1}) - W(x^k, \zeta_k) \\ & \leq \left[ \zeta_{k+1} - \frac{L - L_f}{4} \right] \|x^{k+1} - x^k\|^2 + \left[ \frac{2L - L_f}{2} \beta_k^2 - \zeta_k \right] \|x^k - x^{k-1}\|^2; \end{aligned}$$

- 487 (iii) when  $I(x^k) = I(x^{k+1})$ , we have

$$488 \quad \begin{aligned} & W(x^{k+1}, \zeta_{k+1}) - W(x^k, \zeta_k) \\ & \leq \left[ \zeta_{k+1} - \frac{L}{2} \right] \|x^{k+1} - x^k\|^2 + \left[ \frac{L}{2} \beta_k^2 - \zeta_k \right] \|x^k - x^{k-1}\|^2; \end{aligned}$$

- 489 (iv) when  $I(x^k) = I(x^{k+1})$ ,  $\frac{L}{2} \beta_k^2 \leq \zeta_k$ ; otherwise,  $\frac{2L - L_f}{2} \beta_k^2 \leq \zeta_k$ ;  
 490 (v)  $\{W(x^k, \zeta_k)\}$  is nonincreasing and  $\lim_{k \rightarrow \infty} W(x^k, \zeta_k) = W_\infty < \infty$  exists.

491 Based on the above results, the next lemma shows that the iterate sequence  $\{x^k\}$   
 492 generated by the FIHT algorithm is bounded and  $I(x^k)$  only changes finite times.

493 **LEMMA 3.3.** *The following statements hold:*

- 494 (i) there exists a positive constant  $R$  such that  $\|x^k\| \leq R$ ,  $\forall k \in \mathbb{N}$ ;  
 495 (ii) the  $\ell_0$  norm of sequence  $\{x^k\}$  changes only finitely often.

496 *Proof.* (i). By virtue of the nonincreasing of sequence  $\{W(x^k, \zeta_k)\}$ , we have

$$497 \quad F(x^k) \leq W(x^k, \zeta_k) \leq W(x^1, \zeta_1) = F(x^1).$$

498 Combining this and the level boundedness of  $f$  on  $\mathcal{X}$ , it is easy to obtain the bound-  
 499 edness of  $\{x^k\}$ .

500 (ii). Also denote  $\mathcal{K} = \{k : I(x^k) \neq I(x^{k+1})\}$ , and we prove the finiteness of set  
 501  $\mathcal{K}$  by contradiction. Let's assume that  $\mathcal{K}$  contains infinite elements. By (ii) and (iv)  
 502 of **Lemma 3.2** and  $\zeta_{k+1} = \frac{L - L_f}{8}$  for  $k \in \mathcal{K}$ , we obtain

$$503 \quad W(x^{k+1}, \zeta_{k+1}) - W(x^k, \zeta_k) \leq -\frac{L - L_f}{8} \|x^{k+1} - x^k\|^2, \quad \forall k \in \mathcal{K}.$$

504 Summing up the above inequality over  $k \in \mathcal{K}$  and using **Lemma 3.2**-(v), we find

$$505 \quad \begin{aligned} (3.3) \quad & \sum_{k \in \mathcal{K}} \frac{L - L_f}{8} \|x^{k+1} - x^k\|^2 \leq \sum_{k \in \mathcal{K}} [W(x^k, \zeta_k) - W(x^{k+1}, \zeta_{k+1})] \\ & \leq \sum_{k=1}^{\infty} [W(x^k, \zeta_k) - W(x^{k+1}, \zeta_{k+1})] \\ & = W(x^1, \zeta_1) - W_\infty < \infty. \end{aligned}$$

In addition, from **Lemma 3.1**-(ii), we have

$$\sum_{k \in \mathcal{K}} \|x^{k+1} - x^k\|^2 \geq \sum_{k \in \mathcal{K}} \delta^2 = \infty,$$

506 which is inconsistent with (3.3). Hence,  $I(x^k)$  only changes finite times, as claimed.  $\square$

507 From the above lemma, we can easily find that the FIHT algorithm is reduced to  
 508 the proximal gradient algorithm with extrapolation  $\beta_k = \frac{k-1}{k+\alpha-1}$  for solving  $\min_{\mathcal{X}_J} f$   
 509 after a finite number of iterations, where  $J$  is a fixed index set by [Lemma 3.3](#). Liter-  
 510 ature [\[3\]](#) studied the forward-backward method with extrapolation coefficient  $\frac{k-1}{k+\alpha-1}$   
 511 for solving the sum of a convex function with Lipschitz continuous gradient and a  
 512 proper closed convex function. In the context of this section, the algorithm can be  
 513 written as follows:

$$514 \quad (3.4) \quad \begin{cases} y^k = x^k + \frac{k-1}{k+\alpha-1}(x^k - x^{k-1}), \\ x^{k+1} = \arg \min_{x \in \Omega} \left\{ \langle \nabla f(y^k), x - y^k \rangle + \frac{L}{2} \|x - y^k\|^2 \right\}, \end{cases}$$

515 where  $\Omega$  is a nonempty closed convex set of  $\mathbb{R}^n$ . The convergence results of this  
 516 algorithm in [\[3\]](#) are as follows.

517 **LEMMA 3.4.** [\[3\]](#) *Let  $\{x^k\}$  be a sequence generated by (3.4). When  $\alpha > 3$ , the*  
 518 *following statements hold:*

- 519 (i) *the iterates  $\{x^k\}$  converge to a global minimizer of  $\min_{\Omega} f$ ;*  
 520 (ii)  *$\lim_{k \rightarrow \infty} k^2(f(x^k) - \min f) = 0$  and  $\lim_{k \rightarrow \infty} k\|x^{k+1} - x^k\| = 0$ .*

521 Using the above lemma, we obtain the following significant conclusions for the  
 522 proposed FIHT algorithm.

523 **THEOREM 3.5.** *When  $\alpha > 3$ , sequence  $\{x^k\}$  generated by the FIHT algorithm*  
 524 *satisfies*

- 525 (i)  *$\lim_{k \rightarrow \infty} x^k = x^*$ , where  $x^*$  is a local minimizer of problem (1.1) satisfying the*  
 526 *lower bound property:*

$$527 \quad (3.5) \quad |x_i^*| \geq \delta \quad \text{or} \quad x_i^* = 0, \quad \forall i = 1, 2, \dots, n,$$

528 where  $\delta$  is defined as in [\(3.1\)](#). Moreover,  $\lim_{k \rightarrow \infty} f(x^k) = f_{\infty} < \infty$  and  
 529  $\lim_{k \rightarrow \infty} F(x^k) = F_{\infty} < \infty$  exist;

- 530 (ii)  *$\lim_{k \rightarrow \infty} k\|x^{k+1} - x^k\| = 0$ ;*  
 531 (iii)  *$\lim_{k \rightarrow \infty} k^2(f(x^k) - f_{\infty}) = 0$  and  $\lim_{k \rightarrow \infty} k^2(F(x^k) - F_{\infty}) = 0$ .*

532 *Proof.* (i). Statement (ii) of [Lemma 3.3](#) implies that there exist  $\bar{K} \in \mathbb{N}$  and  
 533  $J \subset \{1, 2, \dots, n\}$  such that  $I(x^k) = J, \forall k \geq \bar{K}$ . Therefore, we have

$$534 \quad (3.6) \quad x^{k+1} = \arg \min_{x \in \mathcal{X}_J} \left\{ f(y^k) + \langle \nabla f(y^k), x - y^k \rangle + \frac{L}{2} \|x - y^k\|^2 \right\}, \quad \forall k \geq \bar{K}.$$

535 Observing the above relation and [Lemma 3.4](#), we know  $\lim_{k \rightarrow \infty} x^k = x^*$ , where  $x^*$  is a  
 536 global minimizer of  $\min_{\mathcal{X}_J} f$ , which is a local minimizer of problem (1.1). Combining  
 537 this with [Lemma 3.1](#)(i), we deduce that  $x^*$  has the lower bound property in [\(3.5\)](#).  
 538 This together with the continuity of  $f$ , we obtain all the results in this statement.

- 539 (ii). In view of [\(3.6\)](#) and [Lemma 3.4](#)(ii), we have the estimation in item (ii).

- (iii). From [Lemma 3.3](#)(ii) and [\(3.5\)](#), there exists  $\bar{K} \in \mathbb{N}$  such that

$$\|x^k\|_0 = \|x^*\|_0, \quad \forall k \geq \bar{K}.$$

540 Further, it holds that

$$541 \quad (3.7) \quad f(x^k) - f(x^*) = F(x^k) - F(x^*), \quad \forall k \geq \bar{K}.$$

542 By [\(3.6\)](#) and [Lemma 3.4](#)(ii), we obtain the estimation in (iii).  $\square$

543 If the extrapolation coefficients are chosen below some threshold, the authors in  
 544 [35] proved that the iterates and the function value sequence of the proximal gradient  
 545 algorithm with extrapolation are R-linear convergent when the objective function  
 546 satisfies the error bound condition. It's worth emphasizing that when loss function  $f$   
 547 satisfies the error bound condition and  $\beta_k$  in the FIHT algorithm satisfies  $\sup_k \beta_k < 1$ ,  
 548 the R-linear convergence of sequence  $\{x^k\}$  and objective function values can also be  
 549 obtained for the FIHT algorithm. For easy reading, the convergence results obtained  
 in this paper are summarized in Table 1.

TABLE 1  
 Summary of convergence of the SFIHT algorithm

$f$	$\beta_k$ in Step 1	$\sigma$	local minimizer	lower bound	$\{F(x^k) - F_\infty\}$
nonsmooth	$\left(0, \sqrt{\frac{\mu_k}{\mu_k - 1}}\right)$	$(0, 2)$	—	—	—
	(2.33)	$[\frac{1}{2}, 1]$	✓	—	—
	(2.41)	$(0, 1)$	✓	—	$O(k^{-\sigma})$
		1	✓	—	$O(k^{-1} \ln k)$
		$(1, 2)$	✓	—	$O(k^{\sigma-2})$
smooth	$\frac{k-1}{k+\alpha-1}$	—	✓	✓	$o(k^{-2})$

550

551 **4. Experimental results.** The aim of this section is to verify the theoretical  
 552 results and performance of the proposed two algorithms by some numerical experi-  
 553 ments. The SFIHT algorithm without extrapolation is called the smoothing iterative  
 554 hard thresholding (SIHT) algorithm in this paper. Example 4.1 and Example 4.2  
 555 are an under-determined linear regression problem and an over-determined censored  
 556 regression problem, respectively. The purpose of Example 4.1 and Example 4.2 is to  
 557 illustrate the ability of the SFIHT algorithm for solving the problem, and to compare  
 558 the good performance of the SFIHT algorithm with respect to the SIHT algorithm.  
 559 In Example 4.3, we use the FIHT algorithm to solve the under-determined  $\ell_0$  regular-  
 560 ized least squares problem. At the same time, we compare the performance of FIHT  
 561 algorithm and IHT algorithm in Example 4.3. For different problems, we choose ap-  
 562 propriate equilibrium parameter  $\lambda$  to adjust the data fitting and sparsity. The CPU  
 563 time (in seconds) reported here doesn't include the time of data initialization.

564 The numerical experiments are performed in Python 3.7.0 on a 64-bit Lenovo PC  
 565 with an Intel(R) Core(TM) i7-10710U CPU @1.10GHz 1.61GHz and 16GB RAM.

For any given  $\epsilon > 0$ , we call  $x^\epsilon$  an  $\epsilon$  local minimizer of problem (1.1), if it holds

$$\|\nabla \tilde{f}(x^\epsilon, \mu)\|_{I(x^\epsilon)^c} \leq \epsilon \quad \text{and} \quad \mu \leq \epsilon,$$

566 where  $I(x^\epsilon)^c = \{i : x_i^\epsilon \neq 0\}$ . We stop the proposed algorithm if  $x^k$  is an  $\epsilon$  local  
 567 minimizer of problem (1.1) or the number of iteration  $k$  exceeds 15000. Set some fixed  
 568 parameters  $\mu_0 = 0.7$ ,  $L = 2L_{\tilde{f}}$  and  $\alpha = 4$  throughout the numerical experiments.

569 **Example 4.1** We consider the following  $\ell_0$  regularized linear regression problem:

$$570 \quad (4.1) \quad \min_{-1 \leq x \leq 1} F(x) := \|Ax - b\|_1 + \lambda \|x\|_0,$$

571 where  $A \in \mathbb{R}^{m \times n}$  with  $m < n$ ,  $b \in \mathbb{R}^m$ . We choose a smoothing function of the  $\ell_1$  loss

572 function as below, and it satisfies the conditions in [Definition 2.1](#),

$$573 \quad (4.2) \quad \tilde{f}(x, \mu) = \sum_{i=1}^m \tilde{\theta}(A_i x - b_i, \mu) \quad \text{with} \quad \tilde{\theta}(z, \mu) = \begin{cases} |z| & \text{if } |z| > \mu, \\ \frac{z^2}{2\mu} + \frac{\mu}{2} & \text{if } |z| \leq \mu. \end{cases}$$

574 Three choices of  $\beta_k$  in Step 1 and Step 3 of the SFIHT algorithm are set as follows:

$$575 \quad \beta_k = \frac{k-1}{k+\alpha-1} \sqrt{\left(1 - \frac{1}{2k^{1-\sigma}}\right) \frac{\mu_k}{\mu_{k-1}}} \quad (\text{Step 1}), \quad \beta_k = \sqrt{\frac{L-L_{\tilde{f}}}{4L} \frac{\mu_k}{\mu_{k-1}}} \quad (\text{Step 3b}) \quad \text{and} \quad \beta_k =$$

$$576 \quad \sqrt{\frac{L-L_{\tilde{f}}}{8L-4L_{\tilde{f}}} \frac{\mu_k}{\mu_{k-1}}} \quad (\text{Step 3b-2}). \quad \text{Denote } s \text{ the } \ell_0 \text{ norm of true solution } x^*, \text{ i.e., } \|x^*\|_0 = s.$$

577 For positive integers  $m, n$  and  $s$ , the data is generated as follows:

$$578 \quad \bar{x} = \text{zeros}(n, 1); \quad \bar{n} = \text{randperm}(n); \quad \bar{x}(\bar{n}(1:s)) = \text{randn}(s, 1); \quad a = 0.005;$$

$$579 \quad x^* = (\text{median}([\bar{x}'; l'; u'])); \quad A = \text{orth}(\text{randn}(m, n)'); \quad b = A * x^* + a * \text{randn}(m, 1).$$

580 Set  $\epsilon = 10^{-3}$ ,  $L_{\tilde{f}} = \|A^T A\|$ ,  $\sigma = 0.95$  and  $x^0 = \text{zeros}(n, 1)$ . We randomly generate

581  $A, b$  and  $x^*$  with  $(m, n) = (300, 1000)$  and  $(m, n) = (500, 5000)$ . Fig. [1\(a\)](#) shows that

582 the support set of sequences  $\{x^k\}$  generated by the SFIHT and SIHT algorithm only

583 change finite times and are convergent. Observing Fig. [1\(b\)](#), we find that compared

584 with the SIHT algorithm, the SFIHT algorithm can find a better solution with fewer

585 iterations. From Fig. [1](#), we see that the convergence rate of the SFIHT algorithm

586 is faster than the SIHT algorithm, and the sparsity of the solution obtained by the

587 SFIHT algorithm is also closer to the true solution  $x^*$ . For three different stopping

588 criterions, we record the CPU time and iterations of the two algorithms in [Table 2](#).

589 It's clear that the computational cost of the SFIHT algorithm is much less than that of

590 the SIHT algorithm. The two algorithms find  $\epsilon$  local minimizer with  $\epsilon = 10^{-4}$

591 with the same iterations, because  $\mu_k$  doesn't meet the termination condition when

592  $\|[\nabla \tilde{f}(x^k, \mu_k)]_{I(x^k)^c}\|_\infty \leq \epsilon$  holds.

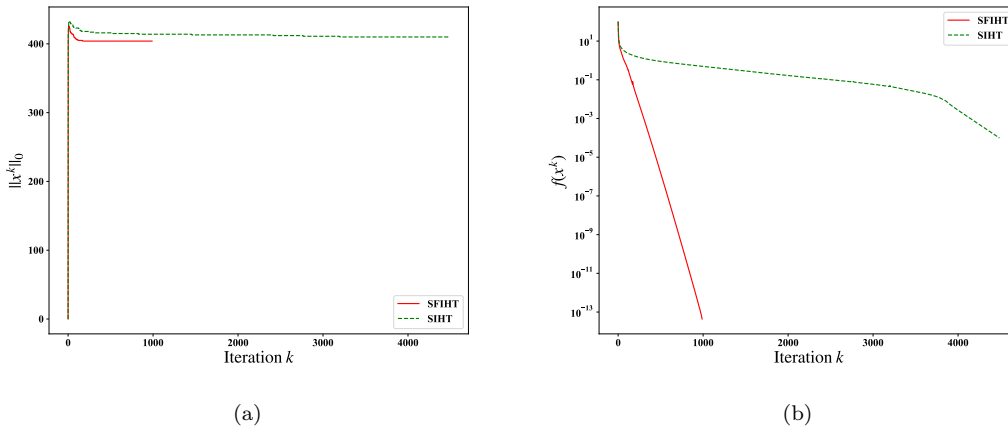


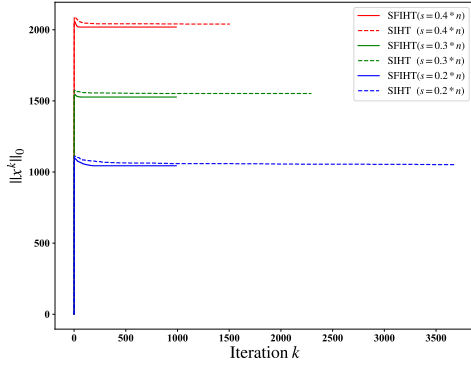
FIG. 1. Convergence of cardinality and loss function values for Example 4.1 with  $m = 300$ ,  $n = 1000$  and  $s = 400$ .

593 When  $m = 500$  and  $n = 5000$ , we generate the problem data with three different  
 594 sparsity levels, which are 20%, 30% and 40%. The results drawn in Fig. [2](#) show that  
 595 the SFIHT algorithm substantially outperforms the SIHT algorithm in terms of the  
 596 solution quality both on the loss function value and the cardinality.

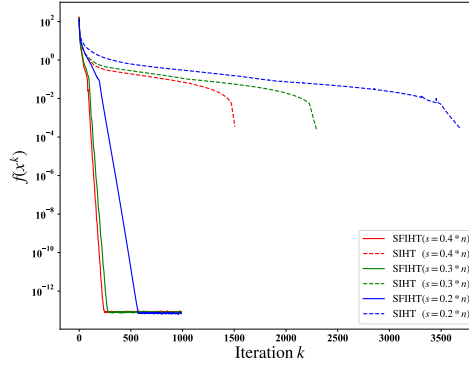
TABLE 2

Computational cost for Example 4.1 with different stopping criterions when  $m = 300$ ,  $n = 1000$  and  $s = 400$ .

$\epsilon$	$10^{-2}$		$10^{-3}$		$10^{-4}$	
Algorithm	SFIHT	SIHT	SFIHT	SIHT	SFIHT	SIHT
Time	<b>0.168</b>	23.388	<b>1.644</b>	27.108	<b>192.258</b>	209.628
Iterations	<b>217</b>	4150	<b>988</b>	4487	11155	11155



(a)



(b)

FIG. 2. Convergence of cardinality and loss function values for Example 4.1 with  $m = 500$ ,  $n = 5000$  and different sparsity levels.

597 **Example 4.2** We consider the following  $\ell_0$  regularized censored regression prob-  
598 lem:

$$599 \quad \min_{\mathbf{0} \leq x \leq \mathbf{1}} F(x) := \frac{1}{m} \|\max\{Ax, 0\} - b\|_1 + \lambda \|x\|_0,$$

600 where  $A \in \mathbb{R}^{m \times n}$  with  $m > n$  and  $b \in \mathbb{R}^m$ . A smoothing function satisfying **Defini-**  
601 **tion 2.1** of the above loss function is defined by

$$602 \quad \tilde{f}(x, \mu) = \frac{1}{m} \sum_{i=1}^m \tilde{\theta}(\tilde{\phi}(A_i x, \mu) - b_i, \mu) \quad \text{with} \quad \tilde{\phi}(s, \mu) = \begin{cases} \max\{s, 0\} & \text{if } |s| > \mu, \\ \frac{(s + \mu)^2}{4\mu} & \text{if } |s| \leq \mu. \end{cases}$$

603 Set  $\epsilon = 10^{-2}$ ,  $L_{\tilde{f}} = \frac{3}{2m} \|A^T A\|$ ,  $\sigma = 0.7$  and  $x^0 = 0.1 * \text{ones}(n, 1)$ .  $\beta_k$  in Step 1 of  
604 the SFIHT algorithm is the same as that in (2.41), and the others are the same as in  
605 Example 4.1. We randomly generate the problem data as follows:

$$606 \quad A = \text{randn}(m, n); \quad \bar{n} = \text{randperm}(n); \quad x^* = \text{zeros}(n, 1);$$

$$607 \quad x^*(\bar{n}(1:s)) = \text{unifrnd}(0.1, 1, [s, 1]); \quad b = \max(A * x^* + 0.01 * \text{randn}(m, 1), 0).$$

608 In this example, we run numerical experiments with  $(m, n, s) = (1000, 200, 60)$  and  
609  $(m, n, s) = (2000, 400, 80)$ . Results recorded in Fig. 3 and Fig. 4 show that the  
610 SFIHT algorithm performs much better than the SIHT algorithm in terms of both

611 the  $\ell_0$  norms and loss function values. Similarly, the SFIHT algorithm needs much  
 612 less iterations to get an  $\epsilon$  local minimizer of problem (1.1) than the SIHT algorithm.

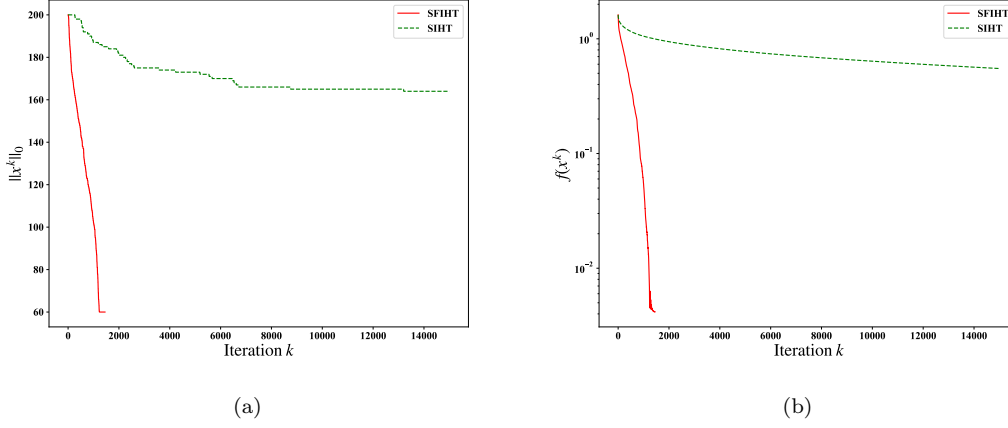


FIG. 3. Convergence of cardinality and loss function values for Example 4.2 with  $m = 1000$ ,  $n = 200$  and  $s = 60$ .

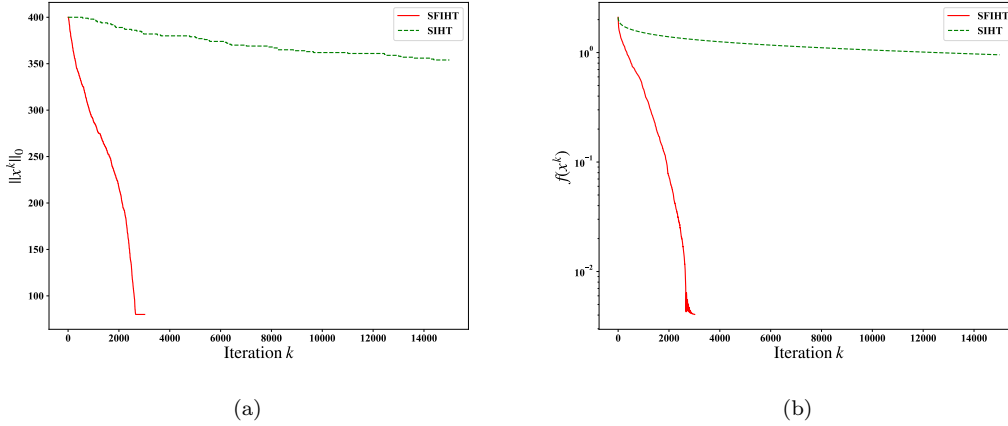


FIG. 4. Convergence of cardinality and loss function values for Example 4.2 with  $m = 2000$ ,  $n = 400$  and  $s = 80$ .

613 **Example 4.3** We consider the following  $\ell_0$  regularized least squares problem:

$$614 \quad (4.3) \quad \min_{0 \leq x \leq 5} F(x) := \frac{1}{2} \|Ax - b\|^2 + \lambda \|x\|_0$$

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

616 We set  $\epsilon = 10^{-4}$ ,  $L_f = \|A^T A\|$  and  $x^0 = \text{zeros}(n, 1)$ . We take the extrapolation

617 coefficients  $\beta_k = \sqrt{\frac{k}{k+1} \frac{L-L_f}{4L}}$  (Step 3b) and  $\beta_k = \sqrt{\frac{k}{k+1} \frac{L-L_f}{8L-4L_f}}$  (Step 3b-2) in the

618 FIHT algorithm. For two cases of  $(m, n)$  and three cases of  $s$ , we randomly generate

619 the problem data as them in Example 4.1. When  $(m, n) = (300, 1000)$ , for the different  
 620 choices of  $s$ , the cardinalities and loss function values versus iteration  $k$  are plotted in  
 621 Fig. 5. From this figure, we can see that the total iterations of the FIHT algorithm  
 622 is much smaller than that of the IHT algorithm, and both the cardinalities and loss  
 623 function values of the final output iterate obtained by the FIHT algorithm are more  
 624 accurate. Fig. 6 illustrates the convergence rate of the FIHT algorithm is also faster  
 625 than that of the IHT algorithm for solving problem (4.3) when the dimension of the  
 626 problems is larger. It's also worth emphasizing that the loss function value at iterative  
 627 point generated by SFIHT algorithm is smaller for all iterations. For different values  
 628 of  $\epsilon$ , the CPU time and iterations for obtaining an  $\epsilon$  local minimizer by the FIHT  
 629 algorithm and the IHT algorithm are recorded in Table 3. From the above results, We  
 630 can observe that the FIHT algorithm performs much better than the IHT algorithm.

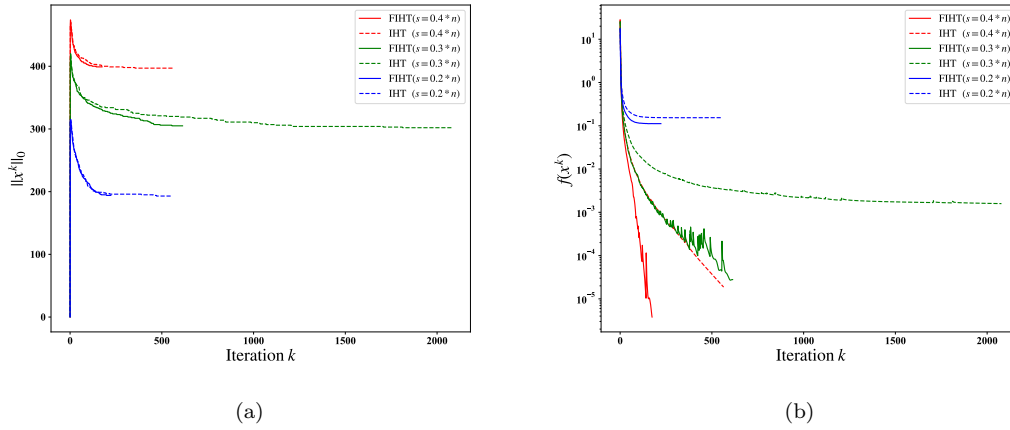


FIG. 5. Convergence of cardinality and loss function values for Example 4.3 with  $m = 300$ ,  $n = 1000$  and different sparsity levels.

TABLE 3

Computational cost for Example 4.3 with different stopping criterions when  $m = 500$ ,  $n = 5000$  and  $s = 1000$ .

$\epsilon$	$10^{-2}$		$10^{-3}$		$10^{-4}$		$10^{-5}$	
	FIHT	IHT	FIHT	IHT	FIHT	IHT	FIHT	IHT
Time	0.180	<b>0.144</b>	<b>0.696</b>	0.738	<b>1.302</b>	2.322	<b>2.280</b>	4.41
Iterations	<b>26</b>	37	<b>92</b>	144	<b>178</b>	346	<b>284</b>	542

631 **5. Conclusions.** The main contribution of this paper is to propose an effective  
 632 and fast algorithm for solving the constrained cardinality penalty problem with a  
 633 continuous convex loss function, and analyze its convergence properties. We first use  
 634 a parametric smoothing approximation of the loss function to generate a cardinality  
 635 penalty problem with smooth loss function. Then, the iterative hard thresholding  
 636 algorithm with extrapolation is used, in which the smoothing parameter is updated  
 637 step by step. The only one subproblem in the proposed algorithm has a closed-form

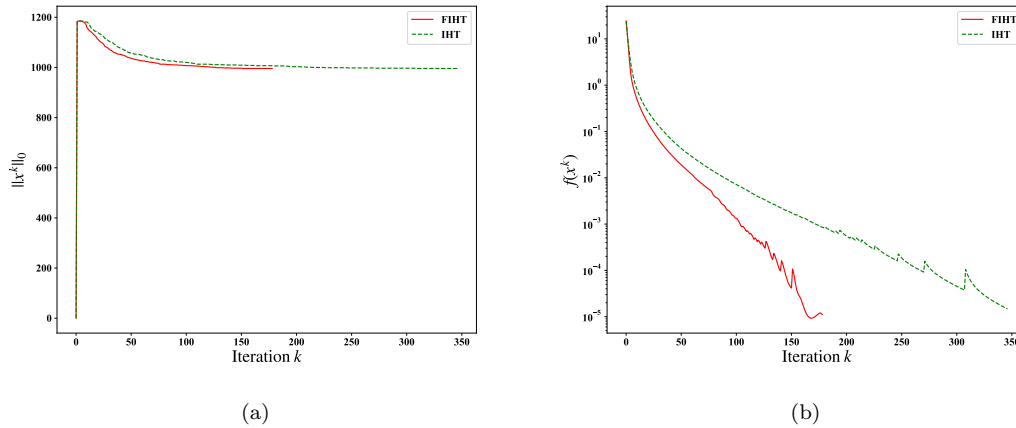


FIG. 6. Convergence of cardinality and loss function values for Example 4.3 with  $m = 500$ ,  $n = 5000$  and  $s = 1000$ .

638 solution, and the extrapolation coefficients can be chosen to satisfy  $\sup_k \beta_k = 1$ . After  
 639 finitely many iterations, the support set of the iterate does not change. Further, we  
 640 give sufficient conditions to guarantee that any accumulation point of  $\{x^k\}$  generated  
 641 by the proposed algorithm is a local minimizer of the considered problem. For a class  
 642 of extrapolation coefficients, we obtain not only the convergence of the sequence,  
 643 but also the convergence rate of  $O(\ln k/k)$  on the loss and objective function values.  
 644 Moreover, we consider in particular the case that the loss function is a Lipschitz  
 645 continuous convex function. To solve it, we provide an algorithm, which can be  
 646 viewed as a specific form of the above proposed algorithm. This algorithm owns both  
 647 sequence convergence on the iterates and convergence rate of  $o(k^{-2})$  on the loss and  
 648 objective function values. Additionally, the limit point not only is a local minimizer  
 649 of the considered problem but also possesses a desirable lower bound.

650

## REFERENCES

- 651 [1] S. ADLY AND H. ATTOUCH, *Finite convergence of proximal-gradient inertial algorithms combin-*  
 652 *ing dry friction with Hessian-driven damping*, SIAM J. Optim., 30 (2020), pp. 2134–2162.  
 653 [2] H. ATTOUCH AND S. LÁSZLÓ, *Newton-like inertial dynamics and proximal algorithms governed*  
 654 *by maximally monotone operators*, SIAM J. Optim., 30 (2020), pp. 3252–3283.  
 655 [3] H. ATTOUCH AND J. PEYPOUQUET, *The rate of convergence of Nesterov’s accelerated forward-*  
 656 *backward method is actually faster than  $1/k^2$* , SIAM J. Optim., 26 (2016), pp. 1824–1834.  
 657 [4] A. BECK AND M. TEOULLE, *A fast iterative shrinkage-thresholding algorithm for linear inverse*  
 658 *problems*, SIAM J. Imaging Sci., 2 (2009), pp. 183–202.  
 659 [5] W. BIAN, *Smoothing accelerated algorithm for constrained nonsmooth convex optimization*  
 660 *problems (in chinese)*, Sci. Sin. Math., 50 (2020), pp. 1651–1666.  
 661 [6] W. BIAN AND X. CHEN, *Optimality and complexity for constrained optimization problems with*  
 662 *nonconvex regularization*, Math. Oper. Res., 42 (2017), pp. 1063–1084.  
 663 [7] W. BIAN AND X. CHEN, *A smoothing proximal gradient algorithm for nonsmooth convex re-*  
 664 *gression with cardinality penalty*, SIAM J. Numer. Anal., 58 (2020), pp. 858–883.  
 665 [8] T. BLUMENSATH AND M. DAVIES, *Sparse and shift-invariant representations of music*, IEEE  
 666 *Trans. Audio Speech Lang. Process.*, 14 (2006), pp. 50–57.  
 667 [9] T. BLUMENSATH AND M. DAVIES, *Iterative thresholding for sparse approximations*, J. Fourier  
 668 *Anal. Appl.*, 14 (2008), pp. 629–654.  
 669 [10] T. BLUMENSATH AND M. DAVIES, *Iterative hard thresholding for compressed sensing*, Appl.



- 670 Comput. Harmon. Anal., 27 (2009), pp. 265–274.
- 671 [11] A. BRUCKSTEIN, D. DONOHO, AND M. ELAD, *From sparse solutions of systems of equations to*  
672 *sparse modeling of signals and images*, SIAM Rev., 51 (2009), pp. 34–81.
- 673 [12] E. CANDÈS, J. ROMBERG, AND T. TAO, *Robust uncertainty principles: exact signal reconstruction*  
674 *from highly incomplete frequency information*, IEEE Trans. Inf. Theory, 52 (2006),  
675 pp. 489–509.
- 676 [13] A. CHAMBOLLE, R. DEVORE, N. LEE, AND B. LUCIER, *Nonlinear wavelet image processing:*  
677 *variational problems, compression, and noise removal through wavelet shrinkage*, IEEE  
678 Trans. Image Process., 7 (1998), pp. 319–335.
- 679 [14] X. CHEN, *Smoothing methods for nonsmooth, nonconvex minimization*, Math. Program., 134  
680 (2012), pp. 71–99.
- 681 [15] P. COMBETTES AND V. WAJS, *Signal recovery by proximal forward-backward splitting*, Multi-  
682 scale Model. Simul., 4 (2005), pp. 1168–1200.
- 683 [16] W. DAI AND O. MILENKOVIC, *Subspace pursuit for compressive sensing signal reconstruction*,  
684 IEEE Trans. Inf. Theory, 55 (2009), pp. 2230–2249.
- 685 [17] I. DAUBECHIES, M. DEFRISE, AND C. DE MOL, *An iterative thresholding algorithm for lin-*  
686 *ear inverse problems with a sparsity constraint*, Commun. Pure Appl. Math., 57 (2004),  
687 pp. 1413–1457.
- 688 [18] N. DOIKOV AND Y. NESTEROV, *Contracting proximal methods for smooth convex optimization*,  
689 SIAM J. Optim., 30 (2020), pp. 3146–3169.
- 690 [19] D. DONOHO, *Compressed sensing*, IEEE Trans. Inf. Theory, 52 (2006), pp. 1289–1306.
- 691 [20] J. FAN AND R. LI, *Variable selection via nonconcave penalized likelihood and its oracle prop-*  
692 *erties*, J. Am. Stat. Assoc., 96 (2001), pp. 1348–1360.
- 693 [21] E. HALE, W. YIN, AND Y. ZHANG, *Fixed-point continuation for  $\ell_1$ -minimization: methodology*  
694 *and convergence*, SIAM J. Optim., 19 (2008), pp. 1107–1130.
- 695 [22] S. HODA, A. GILPIN, J. PENA, AND T. SANDHOLM, *Smoothing techniques for computing Nash*  
696 *equilibria of sequential games*, Math. Oper. Res., 35 (2010), pp. 494–512.
- 697 [23] Y. LIU AND Y. WU, *Variable selection via a combination of the  $\ell_0$  and  $\ell_1$  penalties*, J. Comput.  
698 Graph. Stat., 16 (2007), pp. 782–798.
- 699 [24] Z. LU, *Iterative hard thresholding methods for  $\ell_0$  regularized convex cone programming*, Math.  
700 Program., 147 (2014), pp. 125–154.
- 701 [25] Z. LU AND Y. ZHANG, *Sparse approximation via penalty decomposition methods*, SIAM J.  
702 Optim., 23 (2013), pp. 2448–2478.
- 703 [26] S. MALLAT AND Z. ZHANG, *Matching pursuits with time-frequency dictionaries*, IEEE Trans.  
704 Signal Process., 41 (1993), pp. 3397–3415.
- 705 [27] Y. NESTEROV, *Smooth minimization of non-smooth functions*, Math. Program., 103 (2005),  
706 pp. 127–152.
- 707 [28] Y. NESTEROV, *Gradient methods for minimizing composite functions*, Math. Program., 140  
708 (2013), pp. 125–161.
- 709 [29] M. NIKOLOVA, *Local strong homogeneity of a regularized estimator*, SIAM J. Appl. Math., 61  
710 (2000), pp. 633–658.
- 711 [30] Y. PATI, R. REZAIHAR, AND P. KRISHNAPRASAD, *Orthogonal matching pursuit-recursive func-*  
712 *tion approximation with applications to wavelet decomposition*, in Conference Record of the  
713 Twenty-Seventh Asilomar Conference on Signal, Systems and Computers, vol. 1-2, 1993,  
714 pp. 40–44.
- 715 [31] D. PELEG AND R. MEIR, *A bilinear formulation for vector sparsity optimization*, Signal  
716 Process., 88 (2008), pp. 375–389.
- 717 [32] E. SOUBIES, L. BLANC-FERAUD, AND G. AUBERT, *A continuous exact  $\ell_0$  penalty (CEL0) for*  
718 *least squares regularized problem*, SIAM J. Imaging Sci., 8 (2015), pp. 1607–1639.
- 719 [33] W. SU, S. BOYD, AND E. CANDÈS, *A differential equation for modeling Nesterov’s accelerated*  
720 *gradient method: theory and insights*, J. Mach. Learn. Res., 17 (2016), pp. 1–43.
- 721 [34] R. TIBSHIRANI, *Regression shrinkage and selection via the Lasso*, J. R. Stat. Soc. Ser. B-  
722 Methodol., 58 (1996), pp. 267–288.
- 723 [35] B. WEN, X. CHEN, AND T. PONG, *Linear convergence of proximal gradient algorithm with*  
724 *extrapolation for a class of nonconvex nonsmooth minimization problems*, SIAM J. Optim.,  
725 27 (2017), pp. 124–145.
- 726 [36] F. WU AND W. BIAN, *Accelerated iterative hard thresholding algorithm for  $\ell_0$  regularized re-*  
727 *gression problem*, J. Glob. Optim., 76 (2020), pp. 819–840.
- 728 [37] Z. ZHENG, Y. FAN, AND J. LV, *High dimensional thresholded regression and shrinkage effect*,  
729 J. R. Stat. Soc. Ser. B-Stat. Methodol., 76 (2014), pp. 627–649.