# $\alpha\ell_1 - \beta\ell_2$ regularization for sparse recovery

## Liang Ding[1] and Weimin Han[2]

[1] Department of Mathematics, Northeast Forestry University, Harbin 150040, People's Republic of China
[2] Department of Mathematics, University of Iowa, Iowa City, IA 52242, United States of America

E-mail: dl@nefu.edu.cn and weimin-han@uiowa.edu

CrossMark

## Abstract

This paper presents a novel regularization with a non-convex, non-smooth term of the form $\alpha\|\cdot\|_{\ell_1} - \beta\|\cdot\|_{\ell_2}$ with parameters $\alpha > \beta \geqslant 0$ to solve ill-posed linear problems with sparse solutions. We investigate the existence, stability and convergence of the regularized solution. It is shown that this type of regularization is well-posed and yields sparse solutions. Under an appropriate source condition, we get the convergence rate $O(\delta)$ in the $\ell_2$-norm for *a priori* and *a posteriori* parameter choice rules, respectively. A numerical algorithm is proposed and analyzed based on an iterative threshold strategy with the generalized conditional gradient method. We prove the convergence even though the regularization term is non-smooth and non-convex. The algorithm can easily be implemented because of its simple structure. Some numerical experiments are performed to test the efficiency of the proposed approach. The experiments show that regularization with $\alpha\|\cdot\|_{\ell_1} - \beta\|\cdot\|_{\ell_2}$ performs better in comparison with the classical $\ell_1$ sparsity regularization and can be used as an alternative to the $\ell_p$ $(0 \leqslant p < 1)$ regularizer.

Keywords: sparsity regularization, non-convex, non-smooth, generalized conditional gradient, soft threshold algorithm

(Some figures may appear in colour only in the online journal)

## 1. Introduction

This paper is concerned with solving an ill-posed operator equation of the form

$$Ax = y, \tag{1.1}$$

where $A : \ell_2 \to Y$ is linear and bounded, $x$ is sparse and $Y$ is a Hilbert space with norm $\|\cdot\|_Y$. Throughout this paper, $\langle \cdot, \cdot \rangle$ denotes the inner product in the $\ell_2$ space. In applications, the data $y$ is not given exactly and only its approximation $y^\delta$ is known with $\|y^\delta - y\|_Y \leqslant \delta$ for a small $\delta > 0$. The most widely adopted method to solve the ill-posed operator equation (1.1) is through sparsity regularization

$$\min_x \frac{1}{q} \|Ax - y^\delta\|_Y^q + \alpha \|x\|_{w,p}^p, \tag{1.2}$$

where $1 \leqslant p < 2, 1 \leqslant q \leqslant 2, \alpha > 0, \frac{1}{q}\|Ax - y^\delta\|_Y^q$ is the fidelity term characterizing the misfit of the data $y^\delta$, $\|x\|_{w,p}^p = \sum w_i |\langle \varphi_i, x \rangle|^p$, $\{w_i > 0\}$ are the weights, and $\{\varphi_i\}$ is an orthonormal basis. In [1], Daubechies *et al* used a wavelet as an orthonormal basis and let $w = \mu w_0$, where $\mu > 0$ is a constant and $w_0$ is the sequence with all entries equal to 1. Over the past two decades, sparsity has become popular and great efforts have been devoted to investigating well-posedness issues and developing algorithms for solving the sparsity regularization problems—see [1–3] and the references therein.

Since the $\ell_p$-norm regularization with $1 \leqslant p < 2$ does not always provide the sparsest solution, the non-convex $\ell_p$-norm sparsity regularization with $0 \leqslant p < 1$ has been proposed as an alternative to (1.2) ([4, 5]). However, the non-convex regularized problem is generally more challenging to analyze and to solve due to the non-convexity and non-differentiability, especially if $p = 0$. In spite of the growing interest in non-convex sparsity regularization, limited work can be found on regularization properties, especially on convergence rates. Special regularization techniques are needed to analyze the $\ell_p$-norm sparsity regularization with $0 \leqslant p < 1$. In [4], a non-convex separable constrained sparsity regularization is investigated. Under an additional boundedness assumption on the chosen weights, a sparsity regularization with weighted regularization terms is analyzed and the well-posedness is proven in [5]. In [6], a generalized notion of Bregman distances is introduced that allows the derivation of convergence rate results for the Tikhonov regularization with non-convex regularization terms. In [7], $\ell_0$-norm regularization problems are investigated in finite-dimensional spaces. In [8], sparsity optimization is studied in infinite-dimensional sequence spaces $\ell_p$ with $p \in [0, 1]$. Recently, some new forms of regularization were proposed as alternatives to the non-convex $\ell_p$-norm. In [9], a Lipschitz continuous regularization term is proposed as the difference between $\ell_1$- and $\ell_2$-norms and the minimization of regularized functionals is studied in finite-dimensional spaces for solving compressed sensing problems. A new regularization term called sorted $\ell_1$ is proposed in [10]. For the non-convex sparsity regularization of nonlinear ill-posed problems, see [11–14] and references therein.

Although the $\ell_p$-norm regularization with $0 \leqslant p < 1$ provides a more sparse solution, the $\ell_1$ regularizer is preferred because it can easily be implemented. Hence a critical issue for non-convex sparsity regularization is the development of numerical algorithms. In [15], a reweighted iterative algorithm is proposed for the $\ell_{1/2}$ regularizer. An iterative algorithm based on the difference of convex functions algorithm is proposed in [9]. For ill-conditioned matrices, numerical examples show that the $\ell_{1-2}$ regularizer performs better than $\ell_1$ and $\ell_{1/2}$. A general framework for non-smooth and non-convex regularizations based on a generalized gradient projection method is analyzed in [16]. Discussions of splitting algorithms for solving

the non-convex sparsity regularization can be found in [17] and references therein. In [18], ADMM (alternating direction method of multipliers) is applied to a non-convex and non-smooth optimization problem and global convergence is obtained.

The aim of this paper is to study the following regularization method to solve the ill-posed linear equation (1.1):

$$\min \mathcal{J}_{\alpha,\beta}^{\delta}(x) = \frac{1}{q}\|Ax - y^{\delta}\|_{Y}^{q} + \mathcal{R}_{\alpha,\beta}(x) \tag{1.3}$$

in $\ell_2$ space with the standard $\ell_2$-norm $\|\cdot\|_{\ell_2}$, where $1 \leqslant q \leqslant 2$ and

$$\mathcal{R}_{\alpha,\beta}(x) := \alpha\|x\|_{\ell_1} - \beta\|x\|_{\ell_2}, \quad \alpha > \beta \geqslant 0. \tag{1.4}$$

Denoting $\eta = \beta/\alpha$, we can equivalently express the functional in (1.4) as

$$\mathcal{R}_{\alpha,\beta}(x) = \alpha\,\mathcal{R}_{\eta}(x),$$

where $\mathcal{R}_{\eta}(x) := \|x\|_{\ell_1} - \eta\|x\|_{\ell_2}$, $\alpha > 0$, $1 > \eta \geqslant 0$. The choice $\mathcal{R}_1(x) = \|x\|_{\ell_1} - \|x\|_{\ell_2}$ was first addressed in [19] for the nonnegative least squares problem. Then it was extended to compressive sensing problems in finite dimensional spaces ([9]). In figure 1, we illustrate contours of the $\ell_1$-norm, $\ell_{1/2}$-norm and $\mathcal{R}_{\alpha,\beta}(x)$ for several different ratios of parameter $\alpha$ and $\beta$.

We see that $\mathcal{R}_{\alpha,\beta}(x)$ behaves more and more like the $\ell_0$-norm as $\beta/\alpha \to 1$. Meanwhile, $\mathcal{R}_{\alpha,\beta}(x)$ converges to a constant multiple of the $\ell_1$-norm as $\beta/\alpha \to 0$. For the case $\beta/\alpha = 1$, $\mathcal{R}_{\alpha,\beta}(x)$ is a good approximation of a constant multiple of $\|x\|_{\ell_0}$. However, the contour of $\mathcal{R}_{\alpha,\beta}(x)$ does not intersect with the coordinate axes, i.e. it is not closed. The main motivation for investigating minimization using regularization (1.4) is that $\alpha\|x\|_{\ell_1} - \beta\|x\|_{\ell_2}$ can be viewed as an approximation of $\|x\|_{\ell_0}$. It has a simpler structure as compared to the regularization with the $\ell_0$- and $\ell_p$-norms for $p < 1$. Commonly used norms, such as $\ell_1$-, $\ell_2$-norm, and their derivatives can easily be evaluated and a numerical solution of problem (1.3) can be implemented by an iterative threshold algorithm. Furthermore, the numerical algorithm can easily be extended to solve nonlinear ill-posed equations—see section 3 for details. Moreover, $\alpha\|x\|_{\ell_1} - \beta\|x\|_{\ell_2}$ can be expressed as $\alpha(\|x\|_{\ell_1} - \|x\|_{\ell_2}) + (\alpha - \beta)\|x\|_{\ell_2}$. From the perspective of elastic-net regularization ([20, 21]), the additional term $(\alpha - \beta)\|x\|_{\ell_2}$ leads to more stable algorithms and allows improved error bounds.

In this paper, we investigate the regularizing properties and numerical algorithm of problem (1.3). Proofs of the existence, stability and convergence are along the lines of the classical regularization. However, some extra work is needed due to the presence of the non-convex regularization term $\mathcal{R}_{\alpha,\beta}(x)$. An inequality is derived under an additional source condition. The convergence rate $O(\delta)$ in the $\ell_2$-norm is proved by applying the inequality. As for a numerical method, we present an iterative soft thresholding (ST) algorithm for problem (1.3) which is based on the generalized conditional gradient method ([22, 23]) and the iterative shrinkage method ([1, 24]). In analogy to a technique presented in [23], we can rewrite the functional $\mathcal{J}_{\alpha,\beta}^{\delta}$ in (1.3) as

$$\mathcal{J}_{\alpha,\beta}^{\delta}(x) = F(x) + \Phi(x),$$

where $F(x) = \frac{1}{q}\|Ax - y^{\delta}\|_{Y}^{q} - \Theta(x)$, $\Phi(x) = \Theta(x) + \alpha\|x\|_{\ell_1} - \beta\|x\|_{\ell_2}$ and $\Theta(x) = \frac{\lambda}{2}\|x\|_{\ell_2}^{2} + \beta\|x\|_{\ell_2}$. We show that $F(x)$ and $\Phi(x)$ have the smoothness and convexity required for the application of the generalized conditional gradient method.

An outline of the rest of this paper is as follows. The next section provides the well-posedness and convergence rate results of the $\alpha\|\cdot\|_{\ell_1} - \beta\|\cdot\|_{\ell_2}$ regularization in $\ell_2$ space. In section 3, inspired by the generalized conditional gradient method, we propose a new iterative ST algorithm. Finally, numerical experiments are presented in section 4.
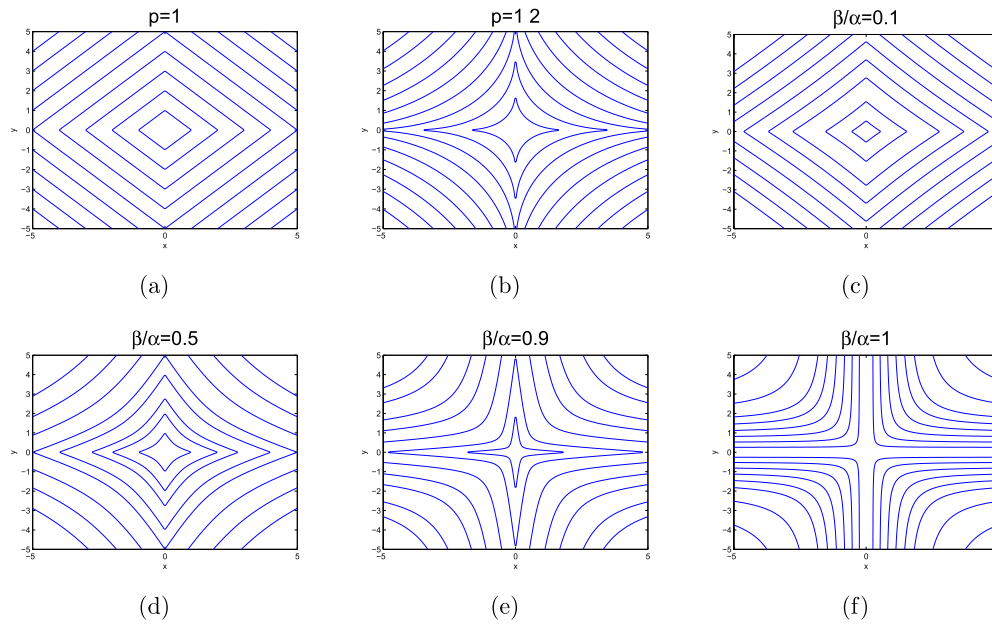
**Figure 1.** Contour plots of $\ell_1$, $\ell_{1/2}$ and $\mathcal{R}_{\alpha,\beta}(x)$ with different ratios between $\alpha$ and $\beta$.

## 2. Regularizing properties

### 2.1. Preliminaries

We denote by

$$x^\delta_{\alpha,\beta} = \arg\min_x \left\{ \frac{1}{q} \|Ax - y^\delta\|^q_Y + \mathcal{R}_{\alpha,\beta}(x) \right\} \tag{2.1}$$

a minimizer of the regularization functional $\mathcal{J}^\delta_{\alpha,\beta}(x)$ in (1.3) for every $\alpha > \beta \geqslant 0$. We use the following definition of the $\mathcal{R}_\eta$-minimum solution ([25]).

**Definition 2.1.** An element $x^\dagger \in \ell_2$ is called an $\mathcal{R}_\eta$-minimum solution of the linear problem $Ax = y$ if

$$Ax^\dagger = y \text{ and } \mathcal{R}_\eta(x^\dagger) = \min_x\{\mathcal{R}_\eta(x) \mid Ax = y\}.$$

We recall the definition of sparsity ([1]).

**Definition 2.2.** $x \in \ell_2$ is called sparse if $\text{supp}(x) := \{i \in \mathbb{N} \mid x_i \neq 0\}$ is finite, where $x_i$ is the $i$th component of $x$. $\|x\|_0 := \text{supp}(x)$ is the cardinality of $\text{supp}(x)$. If $\|x\|_0 = s$ for some $s \in \mathbb{N}$, then $x \in \ell_2$ is called $s$-sparse.

**Definition 2.3.** Define

$$I(x^\dagger) = \{i \in \mathbb{N} \mid x^\dagger_i \neq 0\},$$

where $x^\dagger_i$ is the $i$th component of $x^\dagger$ and $x^\dagger$ is an $\mathcal{R}_\eta$-minimum solution of the linear problem $Ax = y$.

**Remark 2.4.** If $x^\dagger$ is sparse, i.e. $I(x^\dagger)$ is finite, then there exists a number $m > 0$ such that

$$\min_{i \in I(x^\dagger)} |x_i^\dagger| = m.$$

**Lemma 2.5 (Coercivity).** *Assume $\alpha > \beta \geqslant 0$. The functional $\mathcal{R}_{\alpha,\beta} : \ell_2 \to [0, +\infty]$ is coercive, i.e. $\|x\|_{\ell_2} \to +\infty$ implies $\mathcal{R}_{\alpha,\beta}(x) \to +\infty$.*

**Proof.** Note that $\|x\|_{\ell_1} \geqslant \|x\|_{\ell_2}$. So

$$\mathcal{R}_{\alpha,\beta}(x) = \alpha(\|x\|_{\ell_1} - \|x\|_{\ell_2}) + (\alpha - \beta)\|x\|_{\ell_2} \geqslant (\alpha - \beta)\|x\|_{\ell_2},$$

from which it is obvious that $\|x\|_{\ell_2} \to +\infty$ implies $\mathcal{R}_{\alpha,\beta}(x) \to +\infty$. ∎

**Remark 2.6.** Write

$$\mathcal{R}_{\alpha,\beta}(x) = (\alpha - \beta)\|x\|_{\ell_1} + \beta(\|x\|_{\ell_1} - \|x\|_{\ell_2}) \geqslant (\alpha - \beta)\|x\|_{\ell_1}.$$

If $\|x\|_{\ell_1} \to +\infty$, then $\mathcal{R}_{\alpha,\beta}(x) \to +\infty$. So $\mathcal{R}_{\alpha,\beta}(x)$ is also coercive with respect to the $\ell_1$-norm.

Note that $\mathcal{R}_{\alpha,\beta}(x)$ is not coercive when $\alpha = \beta$. For example, let $x = (\underbrace{0, \cdots, 0, x_i}_{i}, 0, \cdots)$, then $\|x\|_{\ell_2} \to +\infty$ as $|x_i| \to +\infty$. However, $\mathcal{R}_{\alpha,\beta}(x) \equiv 0$ for any $x_i$.

Next we recall an extension of Fatou's lemma ([26, pp 321–2]).

**Lemma 2.7 (Extension of Fatou's lemma).** *Let $f_1, f_2, \ldots$ be a sequence of real-valued measurable functions defined on a measure space $(S, \Sigma, \mu)$. If there exists an integrable function $g$ on $S$ such that $f_n \geqslant -g$ for all $n$, then*

$$\int_S \liminf_n f_n \mathrm{d}\mu \leqslant \liminf_n \int_S f_n \mathrm{d}\mu.$$

Lemma 2.7 can be proven by applying Fatou's lemma to the non-negative sequence $\{f_n + g\}$. In lemma 2.7, $S$ is a (nonempty) set, $\Sigma$ is an $\sigma$-algebra on the set $S$, and $\mu$ is a measure on $(S, \Sigma)$. A $\sigma$-algebra (also $\sigma$-field) on a set $S$ is a collection $\Sigma$ of subsets of $S$ that includes $S$ itself, is closed under complement, and is closed under countable unions. Elements of the $\sigma$-algebra are called measurable sets. An ordered triad $(S, \Sigma, \mu)$ is called a measurable space.

**Lemma 2.8 (Weak lower semi-continuity).** *Let $M > 0$ be given. Then, for any $x_n \in \ell_2$ with $\mathcal{R}_{\alpha,\beta}(x_n) \leqslant M$, $\{x_n\}$ weakly converging to $x$ in $\ell_2$ implies $\liminf_n \mathcal{R}_{\alpha,\beta}(x_n) \geqslant \mathcal{R}_{\alpha,\beta}(x)$.*

**Proof.** By the definition of $\mathcal{R}_{\alpha,\beta}$ in (1.4), we obtain

$$
\begin{aligned}
\mathcal{R}_{\alpha,\beta}(x_n) - \mathcal{R}_{\alpha,\beta}(x) &= \alpha(\|x_n\|_{\ell_1} - \|x\|_{\ell_1}) - \beta(\|x_n\|_{\ell_2} - \|x\|_{\ell_2}) \\
&= \sum_i \alpha(|x_n^i| - |x^i|) - \beta \frac{\sum_i (|x_n^i| + |x^i|)(|x_n^i| - |x^i|)}{\|x_n\|_{\ell_2} + \|x\|_{\ell_2}} \quad (2.2) \\
&= \sum_i \left[\alpha - \frac{\beta(|x_n^i| + |x^i|)}{\|x_n\|_{\ell_2} + \|x\|_{\ell_2}}\right](|x_n^i| - |x^i|),
\end{aligned}
$$

where $x^i$ and $x_n^i$ are the $i$th components of $x$ and $x_n$, respectively. If $x_n \neq 0$ or $x \neq 0$, define $c_n^i := \alpha - \frac{\beta(|x_n^i| + |x^i|)}{\|x_n\|_{\ell_2} + \|x\|_{\ell_2}}$; then $0 < \alpha - \beta \leqslant c_n^i \leqslant \alpha$. If $x_n = 0$ and $x = 0$, then let $c_n^i = 0$. From (2.2),

$$\liminf_n (\mathcal{R}_{\alpha,\beta}(x_n) - \mathcal{R}_{\alpha,\beta}(x)) = \liminf_n \left[ \sum_i c_n^i (|x_n^i| - |x^i|) \right]. \tag{2.3}$$

By the definition of $c_n^i$, we have

$$c_n^i(|x_n^i| - |x^i|) \geqslant -c_n^i |x^i| \geqslant -\alpha |x^i|. \tag{2.4}$$

Meanwhile, $\mathcal{R}_{\alpha,\beta}(x_n) \leqslant M$ implies that $\{\|x_n\|_{\ell_1}\}$ is bounded. Then it follows from $\|x\|_{\ell_1} \leqslant \liminf_n \|x_n\|_{\ell_1}$ that $\|x\|_{\ell_1}$ is finite. Hence,

$$\sum_i \alpha |x^i| \neq \infty. \tag{2.5}$$

With (2.4) and (2.5) at our disposal, we apply lemma 2.7 to find

$$\liminf_n \left[ \sum_i c_n^i (|x_n^i| - |x^i|) \right] \geqslant \sum_i \liminf_n \left( c_n^i |x_n^i| - c_n^i |x^i| \right). \tag{2.6}$$

From the weak convergence of $x_n$ to $x$, we have $|x_n^i| \to |x^i|$ for all $i \in \mathbb{N}$. Since $0 < c_n^i \leqslant \alpha$, it is obvious that $c_n^i |x_n^i| - c_n^i |x^i| \to 0$. Then we have

$$\liminf_n (c_n^i |x_n^i| - c_n^i |x^i|) = 0. \tag{2.7}$$

Hence,

$$\sum_i \liminf_n \left( c_n^i |x_n^i| - c_n^i |x^i| \right) = 0. \tag{2.8}$$

A combination of (2.3), (2.6) and (2.8) implies that $\liminf_n (\mathcal{R}_{\alpha,\beta}(x_n) - \mathcal{R}_{\alpha,\beta}(x)) \geqslant 0$, which proves the lemma. ∎

**Remark 2.9.** Note that lemma 2.8 still holds when $\alpha = \beta$, i.e. $\|\cdot\|_{\ell_1} - \|\cdot\|_{\ell_2}$ is weakly lower semi-continuous. For the case $\alpha = \beta$, $0 \leqslant c_n^i \leqslant \alpha$, the above proof is still valid.

**Lemma 2.10 (Radon–Riesz property).** *Let $M > 0$ be given. Then, for any $x_n \in \ell_2$ with $\mathcal{R}_{\alpha,\beta}(x_n) \leqslant M$, if $x_n$ converges weakly to $x$ in $\ell_2$ and $\mathcal{R}_{\alpha,\beta}(x_n) \to \mathcal{R}_{\alpha,\beta}(x)$, then $x_n$ converges strongly to $x$ in $\ell_2$.*

**Proof.** By the assumption $\mathcal{R}_{\alpha,\beta}(x_n) \to \mathcal{R}_{\alpha,\beta}(x)$, we have

$$\alpha \|x_n\|_{\ell_1} - \beta \|x_n\|_{\ell_2} \to \alpha \|x\|_{\ell_1} - \beta \|x\|_{\ell_2},$$

i.e.

$$\alpha(\|x_n\|_{\ell_1} - \|x_n\|_{\ell_2}) + (\alpha - \beta)\|x_n\|_{\ell_2} \to \alpha(\|x\|_{\ell_1} - \|x\|_{\ell_2}) + (\alpha - \beta)\|x\|_{\ell_2}. \tag{2.9}$$

Next, we prove $\|x_n\|_{\ell_2} \to \|x\|_{\ell_2}$ and argue by contradiction. Suppose $\|x_n\|_{\ell_2} \not\to \|x\|_{\ell_2}$. Since $x_n \rightharpoonup x$ in $\ell_2$, we have $\|x\|_{\ell_2} \leqslant \liminf_n \|x_n\|_{\ell_2}$. Thus, there exists a constant $c > 0$ such that $c = \limsup_n \|x_n\|_{\ell_2} > \|x\|_{\ell_2}$. Consequently, there exists a subsequence $\{x_m\}$ of $\{x_n\}$ such that

$$\lim_m \|x_m\|_{\ell_2} = c > \|x\|_{\ell_2}.$$

Hence,

$$\lim_m (\alpha - \beta)\|x_m\|_{\ell_2} = c(\alpha - \beta) > (\alpha - \beta)\|x\|_{\ell_2}. \tag{2.10}$$

By (2.9), we have

$$\alpha(\|x_m\|_{\ell_1} - \|x_m\|_{\ell_2}) + (\alpha - \beta)\|x_m\|_{\ell_2} \to \alpha(\|x\|_{\ell_1} - \|x\|_{\ell_2}) + (\alpha - \beta)\|x\|_{\ell_2}. \tag{2.11}$$

A combination of (2.10) and (2.11) implies that

$$\lim_m \alpha(\|x_m\|_{\ell_1} - \|x_m\|_{\ell_2}) < \alpha(\|x\|_{\ell_1} - \|x\|_{\ell_2}).$$

Hence,

$$\liminf_n \alpha(\|x_n\|_{\ell_1} - \|x_n\|_{\ell_2}) \leqslant \liminf_m \alpha(\|x_m\|_{\ell_1} - \|x_m\|_{\ell_2}) < \alpha(\|x\|_{\ell_1} - \|x\|_{\ell_2}). \tag{2.12}$$

This contradicts the fact that $\|\cdot\|_{\ell_1} - \|\cdot\|_{\ell_2}$ is weakly lower semi-continuous— see remark 2.9. This argument shows that $\|x_n\|_{\ell_2} \to \|x\|_{\ell_2}$. Since $x_n \rightharpoonup x$ in $\ell_2$ by assumption, we conclude that $x_n \to x$ in $\ell_2$. ∎

### 2.2. Well-posedness of regularization

In this section, we consider the well-posedness of the regularization method. We prove the existence of the regularized solution $x_{\alpha,\beta}^\delta$ defined by (2.1), which continuously depends on the data $y^\delta$ and converges to an $\mathcal{R}_\eta$-minimum solution of the linear problem $Ax = y$. The proof is along the lines of the standard quadratic Tikhonov regularization ([25]) and sparsity regularization ([13, 21, 27, 28]). However, some extra work is needed due to the use of the non-convex regularization term $\mathcal{R}_{\alpha,\beta}(x)$.

**Theorem 2.11 (Existence).** *For all $\alpha > \beta \geqslant 0$ and $y^\delta \in Y$, problem (1.3) has a solution.*

**Proof.** Since $\mathcal{J}_{\alpha,\beta}^\delta(x)$ is nonnegative, there exists a minimizing sequence $\{x_n\}$ such that

$$\lim_{n \to +\infty} \mathcal{J}_{\alpha,\beta}^\delta(x_n) = \lim_{n \to +\infty} \left[\frac{1}{q}\|Ax_n - y^\delta\|_Y^q + \mathcal{R}_{\alpha,\beta}(x_n)\right] = c := \inf \mathcal{J}_{\alpha,\beta}^\delta(x) \geqslant 0.$$

We see that $\mathcal{R}_{\alpha,\beta}(x_n) = \alpha(\|x_n\|_{\ell_1} - \|x_n\|_{\ell_2}) + (\alpha - \beta)\|x_n\|_{\ell_2}$ is bounded with respect to $n$. Then $\{\|x_n\|_{\ell_1}\}$ and $\{\|x_n\|_{\ell_2}\}$ are bounded by lemma 2.5 and remark 2.6. Thus $\{x_n\}$ has a subsequence $\{x_{n_k}\}$ which is weakly convergent to an element $\bar{x}$ in $\ell_2$ space, i.e. $x_{n_k} \rightharpoonup \bar{x}$ in $\ell_2$. By lemma 2.8,

$$\mathcal{R}_{\alpha,\beta}(\bar{x}) \leqslant \liminf_{k \to +\infty} \mathcal{R}_{\alpha,\beta}(x_{n_k}). \tag{2.13}$$

On the other hand, since $A$ is bounded and linear, $A(x_{n_k}) - y^\delta \rightharpoonup A(\bar{x}) - y^\delta$ in $Y$, and it follows from the weak lower semi-continuity of the norm that

$$\frac{1}{q}\|A\bar{x} - y^\delta\|_Y^q \leqslant \liminf_{k \to +\infty} \frac{1}{q}\|Ax_{n_k} - y^\delta\|_Y^q. \tag{2.14}$$

A combination of (2.13) and (2.14) shows that

$$\frac{1}{q}\|A\bar{x} - y^\delta\|_Y^q + \mathcal{R}_{\alpha,\beta}(\bar{x}) \leqslant \liminf_{k \to +\infty} \frac{1}{q}\|Ax_{n_k} - y^\delta\|_Y^q + \liminf_{k \to +\infty} \mathcal{R}_{\alpha,\beta}(x_{n_k})$$

$$\leqslant \liminf_{k \to +\infty} \left[\frac{1}{q}\|Ax_{n_k} - y^\delta\|_Y^q + \mathcal{R}_{\alpha,\beta}(x_{n_k})\right].$$

Hence, $\bar{x}$ minimizes $\mathcal{J}_{\alpha,\beta}^\delta(x)$. ∎

**Theorem 2.12 (Stability).** *Let $\alpha > \beta \geqslant 0$ and $\{y_n\}$ and $\{x_n\}$ be sequences with* $\lim_{n \to +\infty} \|y_n - y^\delta\| = 0$, $x_n$ *being a minimizer of* $\mathcal{J}^{\delta_n}_{\alpha_n, \beta_n}(x)$, *where* $\alpha_n > \beta_n \geqslant 0$ *and* $\alpha_n \to \alpha$, $\beta_n \to \beta$ *as* $n \to +\infty$. *Then* $\{x_n\}$ *contains a convergent subsequence* $\{x_{n_k}\}$ *and the limit* $x^\delta_{\alpha, \beta}$ *of every convergent subsequence is a minimizer of* $\mathcal{J}^\delta_{\alpha, \beta}(x)$. *If the minimizer of* $\mathcal{J}^\delta_{\alpha, \beta}(x)$ *is unique, then* $\lim_{k \to +\infty} \|x_{n_k} - x^\delta_{\alpha, \beta}\|_{\ell_2} = 0$.

**Proof.** By definition of $x_n$, we have

$$\frac{1}{q} \|Ax_n - y_n\|_Y^q + \alpha_n \|x_n\|_{\ell_1} - \beta_n \|x_n\|_{\ell_2} \leqslant \frac{1}{q} \|Ax - y_n\|_Y^q + \alpha_n \|x\|_{\ell_1} - \beta_n \|x\|_{\ell_2}$$

$$(2.15)$$

for all $x \in \ell_1$. Then $\{\|x_n\|_{\ell_2}\}$ and $\{\|x_n\|_{\ell_1}\}$ are bounded. Hence, there exists a subsequence $\{x_{n_k}\}$ of $\{x_n\}$ and $x^\delta_{\alpha, \beta}$ such that

$$x_{n_k} \rightharpoonup x^\delta_{\alpha, \beta} \text{ in } \ell_2, \quad Ax_{n_k} \rightharpoonup Ax^\delta_{\alpha, \beta} \text{ in } Y.$$

By the weak lower semi-continuity of the norm, we obtain

$$\frac{1}{q} \|Ax^\delta_{\alpha, \beta} - y^\delta\|_Y^q \leqslant \liminf_{k \to +\infty} \frac{1}{q} \|Ax_{n_k} - y_{n_k}\|_Y^q.$$

$$(2.16)$$

By lemma 2.8, we have

$$\mathcal{R}_{\alpha, \beta}(x^\delta_{\alpha, \beta}) \leqslant \liminf_{k \to +\infty} \mathcal{R}_{\alpha, \beta}(x_{n_k}) = \liminf_{k \to +\infty} \left( \alpha_n \|x_{n_k}\|_{\ell_1} - \beta_n \|x_{n_k}\|_{\ell_2} \right). \quad (2.17)$$

A combination of (2.15), (2.16) and (2.17) implies that

$$\frac{1}{q} \|Ax^\delta_{\alpha, \beta} - y^\delta\|_Y^q + \mathcal{R}_{\alpha, \beta}(x^\delta_{\alpha, \beta}) \leqslant \liminf_{k \to +\infty} \left[ \frac{1}{q} \|Ax_{n_k} - y_{n_k}\|_Y^q + \alpha_n \|x_{n_k}\|_{\ell_1} - \beta_n \|x_{n_k}\|_{\ell_2} \right]$$

$$\leqslant \liminf_{k \to +\infty} \left[ \frac{1}{q} \|Ax - y_{n_k}\|_Y^q + \alpha_n \|x\|_{\ell_1} - \beta_n \|x\|_{\ell_2} \right]$$

$$= \frac{1}{q} \|Ax - y^\delta\|_Y^q + \mathcal{R}_{\alpha, \beta}(x)$$

for all $x \in \ell_2$. This implies that $x^\delta_{\alpha, \beta}$ is a minimizer of $\mathcal{J}^\delta_{\alpha, \beta}(x)$.

On the other hand, we note that

$$\limsup_{k \to +\infty} \left[ \frac{1}{q} \|Ax_{n_k} - y_{n_k}\|_Y^q + \mathcal{R}_{\alpha, \beta}(x_{n_k}) \right] = \limsup_{k \to +\infty} \left[ \frac{1}{q} \|Ax_{n_k} - y_{n_k}\|_Y^q + \alpha_n \|x_{n_k}\|_{\ell_1} - \beta_n \|x_{n_k}\|_{\ell_2} \right]$$

$$\leqslant \limsup_{k \to +\infty} \left[ \frac{1}{q} \|Ax^\delta_{\alpha, \beta} - y_{n_k}\|_Y^q + \alpha_n \|x^\delta_{\alpha, \beta}\|_{\ell_1} - \beta_n \|x^\delta_{\alpha, \beta}\|_{\ell_2} \right]$$

$$= \frac{1}{q} \|Ax^\delta_{\alpha, \beta} - y^\delta\|_Y^q + \mathcal{R}_{\alpha, \beta}(x^\delta_{\alpha, \beta}).$$

Hence,

$$\frac{1}{q} \|Ax_{n_k} - y_{n_k}\|_Y^q + \mathcal{R}_{\alpha, \beta}(x_{n_k}) \to \frac{1}{q} \|Ax^\delta_{\alpha, \beta} - y^\delta\|_Y^q + \mathcal{R}_{\alpha, \beta}(x^\delta_{\alpha, \beta}).$$

Since both $\|\cdot\|_Y$ and $\mathcal{R}_{\alpha,\beta}$ are weakly lower semi-continuous, this implies that $\mathcal{R}_{\alpha,\beta}(x_{n_k}) \to \mathcal{R}_{\alpha,\beta}(x_{\alpha,\beta}^\delta)$. If the minimizer of $\mathcal{J}_{\alpha,\beta}^\delta(x)$ is unique, then $\lim_{k\to+\infty}\|x_{n_k} - x_{\alpha,\beta}^\delta\|_{\ell_2} = 0$ through an application of lemma 2.10. ∎

**Theorem 2.13 (Convergence).**  *Let $x_{\alpha_n,\beta_n}^{\delta_n}$ be a minimizer of $\mathcal{J}_{\alpha_n,\beta_n}^{\delta_n}(x)$ defined by (2.1) with the data $y^{\delta_n}$ satisfying $\|y - y^{\delta_n}\| \leqslant \delta_n$, where $\delta_n \to 0$ if $n \to +\infty$ and $y^{\delta_n}$ belongs to the range of A. Assume $\alpha_n := \alpha(\delta_n)$, $\beta_n := \beta(\delta_n)$, $\alpha_n > \beta_n \geqslant 0$, are such that*

$$\lim_{n\to\infty}\alpha_n = 0, \quad \lim_{n\to\infty}\beta_n = 0 \quad and \quad \lim_{n\to\infty}\frac{\delta_n^q}{\alpha_n} = 0.$$

*Moreover, assume that $\eta = \lim_{n\to\infty}\eta_n \in [0,1)$ exists, where $\eta_n = \beta_n/\alpha_n$. Then there exists a subsequence of $\{x_{\alpha_n,\beta_n}^{\delta_n}\}$, still denoted by $\{x_{\alpha_n,\beta_n}^{\delta_n}\}$, such that $x_{\alpha_n,\beta_n}^{\delta_n}$ converges to an $\mathcal{R}_\eta$-minimizing solution $x^\dagger$ in $\ell_2$. If, in addition, the $\mathcal{R}_\eta$-minimizing solution $x^\dagger$ is unique, then*

$$\lim_{n\to+\infty}\|x_{\alpha_n,\beta_n}^{\delta_n} - x^\dagger\|_{\ell_2} = 0.$$

**Proof.**   Denote $y_n := y^{\delta_n}$, $x_n := x_{\alpha_n,\beta_n}^{\delta_n}$, $\eta_n := \eta^{\delta_n}$. By the definition of $x_n$, we obtain

$$\frac{1}{q}\|Ax_n - y_n\|_Y^q + \alpha_n\|x_n\|_{\ell_1} - \beta_n\|x_n\|_{\ell_2} \leqslant \frac{1}{q}\|Ax^\dagger - y_n\|_Y^q + \alpha_n\|x^\dagger\|_{\ell_1} - \beta_n\|x^\dagger\|_{\ell_2}$$

$$\leqslant \frac{1}{q}\delta_n^q + \alpha_n\|x^\dagger\|_{\ell_1} - \beta_n\|x^\dagger\|_{\ell_2}. \qquad (2.18)$$

By assumption, we have $\frac{1}{q}\delta_n^q + \alpha_n\|x^\dagger\|_{\ell_1} - \beta_n\|x^\dagger\|_{\ell_2} \to 0$ as $n \to +\infty$. Hence,

$$\|Ax_n - y_n\|_Y \to 0 \quad (n \to +\infty). \qquad (2.19)$$

Moreover, we have

$$\|Ax_n - y\|_Y \leqslant \|Ax_n - y_n\|_Y + \|y - y_n\|_Y \leqslant \|Ax_n - y_n\|_Y + \delta_n. \qquad (2.20)$$

A combination of (2.19) and (2.20) implies that

$$\lim_{n\to+\infty}Ax_n = y. \qquad (2.21)$$

On the other hand, it follows from (2.18) that

$$\limsup_{n\to+\infty}(\|x_n\|_{\ell_1} - \eta_n\|x_n\|_{\ell_2}) \leqslant \|x^\dagger\|_{\ell_1} - \eta\|x^\dagger\|_{\ell_2}. \qquad (2.22)$$

Since $\|x_n\|_{\ell_1} - \eta_n\|x_n\|_{\ell_2}$ is bounded, there exists an $x^* \in \ell_2$ and a subsequence of $\{x_n\}$, still denoted by $\{x_n\}$, such that $x_n \rightharpoonup x^*$ in $\ell_2$. Together with (2.21), it follows that

$$y = \lim_{n\to+\infty}Ax_n = A(x^*).$$

Meanwhile, by lemma 2.8, we have

$$\|x^*\|_{\ell_1} - \eta\|x^*\|_{\ell_2} \leqslant \liminf_n(\|x_n\|_{\ell_1} - \eta_n\|x_n\|_{\ell_2})$$
$$\leqslant \|x^\dagger\|_{\ell_1} - \eta\|x^\dagger\|_{\ell_2}. \qquad (2.23)$$

By the definition of $x^\dagger$, then $x^*$ is an $\mathcal{R}_\eta$-minimizing solution. If the $\mathcal{R}_\eta$-minimizing solution is unique, then $x^* = x^\dagger$. A combination of (2.22) and (2.23)

implies $\quad \|x_n\|_{\ell_1} - \eta_n\|x_n\|_{\ell_2} \to \|x^\dagger\|_{\ell_1} - \eta\|x^\dagger\|_{\ell_2}.\quad$ Thus, $\quad \mathcal{R}_{\alpha,\beta}(x_n) \to \mathcal{R}_{\alpha,\beta}(x^\dagger).\quad$ Then $\lim_{n\to+\infty}\|x_n - x^\dagger\|_{\ell_2} = 0$ by lemma 2.10. ∎

**Proposition 2.14 (Sparsity).** *Every minimizer $x$ of $\mathcal{J}^\delta_{\alpha,\beta}(x)$ is sparse.*

**Proof.** The proof is along the lines of the proof of proposition 4.5 in [5]. For simplicity, we only discuss the case $q = 2$. Define the sequence $\bar{x} := x - x_i e_i$ for $i \in \mathbb{N}$, where $e_i = (\underbrace{0, \cdots, 0, 1}_{i}, 0, \cdots)$, and $x_i$ is the $i$th component of $x$. By the definition of $x$, it follows that

$$\frac{1}{2}\|Ax - y^\delta\|^2_Y + \mathcal{R}_{\alpha,\beta}(x) \leqslant \frac{1}{2}\|A(x - x_i e_i) - y^\delta\|^2_Y + \mathcal{R}_{\alpha,\beta}(x - x_i e_i). \quad (2.24)$$

If $x = 0$, then $x$ is sparse. If $x \neq 0$, by (2.24), we see that

$$\alpha|x_i| - \beta\frac{|x_i|^2}{\|x\|_{\ell_2} + \|\bar{x}\|_{\ell_2}} = \mathcal{R}_{\alpha,\beta}(x) - \mathcal{R}_{\alpha,\beta}(\bar{x}) \leqslant \frac{1}{2}x_i^2\|Ae_i\|^2_Y - x_i\langle Ae_i, Ax - y^\delta\rangle$$

$$\leqslant \frac{1}{2}x_i^2\|A\|^2 - x_i\langle e_i, A^*(Ax - y^\delta)\rangle \quad (2.25)$$

for every $i \in \mathbb{N}$. Meanwhile, for any constant $0 < c \leqslant 1 - \frac{\beta}{\alpha}$,

$$\alpha c\frac{|x_i|}{1 + |x_i|} \leqslant \alpha|x_i| - \beta|x_i| \leqslant \alpha|x_i| - \beta\frac{|x_i|^2}{\|x\|_{\ell_2} + \|\bar{x}\|_{\ell_2}}. \quad (2.26)$$

Denote

$$K_i := \frac{(1 + \|x\|_{\ell_2})(\frac{1}{2}x_i\|A\|^2 - \langle e_i, A^*(Ax - y^\delta)\rangle)}{c\alpha}.$$

Then a combination of (2.25) and (2.26) implies that

$$K_i x_i \geqslant |x_i|, \quad i \in \mathbb{N}.$$

Since $x \in \ell_2$, $x_i \to 0$ as $i \to \infty$. Also, $\|A\|$ is finite since $A$ is linear and bounded. Moreover, $\langle e_i, A^*(Ax - y^\delta)\rangle = (A^*(Ax - y^\delta))_i$, where $(A^*(Ax - y^\delta))_i$ is the $i$th component of $A^*(Ax - y^\delta)$. Since $A^*(Ax - y^\delta) \in \ell_2$, $(A^*(Ax - y^\delta))_i \to 0$ as $i \to \infty$. Then we have $K_i \to 0$ as $i \to \infty$, which implies that $\Lambda := \{i \in \mathbb{N} \mid |K_i| \geqslant 1\}$ is finite. Obviously, $x_i = 0$ whenever $i \notin \Lambda$. This proves that $x$ is sparse. ∎

**Remark 2.15.** The regularization parameter $\alpha(\delta)$ depends on the noise level $\delta$; in particular, $\alpha(\delta) \to 0$ as $\delta \to 0$. In applications, the observed data $y^\delta$ contains noise and so the noise level $\delta > 0$. For each fixed $\delta$, there is a regularization parameter $\alpha > 0$. Then $K_i \to 0$ as $i \to \infty$.

If $\delta = 0$, then the regularization parameter $\alpha = 0$. The definition of $K_i$ is unreasonable. For this case, $\beta = \eta\alpha$ implies that $\beta = 0$. Then (1.3) becomes

$$\min \mathcal{J}_{\alpha,\beta}(x) = \frac{1}{q}\|Ax - y\|^q_Y.$$

Since this paper is concerned with solving an ill-posed operator equation of the form $Ax = y$ with a sparse solution, the minimizer $x$ of (1.3) is sparse. So the minimizer $x$ of (1.3) is sparse whenever $\alpha \geqslant 0$.

Note that if the ill-posed operator equation $Ax = y$ does not have a sparse solution, then for the case $\alpha = 0$, the minimizer of (1.3) is non-sparse. A natural question is whether one should use sparsity regularization when a linear ill-posed equation does not have a sparse solution. We refer the reader to [29] which provides a discussion regarding solutions that are not completely sparse but have a fast decaying nonzero part.

### 2.3. Convergence rate of the regularized solution

In this section, we present the convergence rate results of *a priori* and *a posteriori* parameter choice rules. An inequality is derived under a source condition, and we obtain the convergence rate $O(\delta)$ in the $\ell_2$-norm based on the inequality. The source condition is stated next.

**Assumption 2.16.** *Let $x^\dagger \neq 0$ be a sparse $\mathcal{R}_\eta$-minimizing solution of the problem $Ax = y$. Assume that*

$$e_i \in R(A^*) \quad \forall\, i \in I(x^\dagger),$$

*where $e_i = (\underbrace{0, \cdots, 0, 1}_{i}, 0, \cdots)$ and $I(x^\dagger)$ is defined in definition 2.3. In other words, for each $i \in I(x^\dagger)$ there exists an element $\omega_i \in D(A^*)$ such that $e_i = A^*\omega_i$.*

Assumption 2.16 and its modified form were introduced in [4, 11]. This assumption can be viewed as a source condition and it implies that the operator $A$ fulfills some kind of 'finite basis injectivity condition' which is commonly used in sparsity regularization.

Next, we present an inequality under the source condition. The linear convergence rate $O(\delta)$ can be derived from this inequality.

**Lemma 2.17.** *Let assumption 2.16 hold and $\mathcal{R}_{\alpha,\beta}(x) \leqslant M$ for a given $M > 0$. Then there exist constants $c_1 > c_2$ with $c_1 > 0$ such that*

$$(\alpha - \beta)\|x - x^\dagger\|_{\ell_1} \leqslant \mathcal{R}_{\alpha,\beta}(x) - \mathcal{R}_{\alpha,\beta}(x^\dagger) + (c_1\alpha - c_2\beta)\|Ax - Ax^\dagger\|_Y. \tag{2.27}$$

**Proof.** From the definition of index set $I(x^\dagger)$, we have

$$(\alpha - \beta)\|x - x^\dagger\|_{\ell_1} = (\alpha - \beta)\left(\sum_{i \in I(x^\dagger)} |x_i - x_i^\dagger| + \sum_{i \notin I(x^\dagger)} |x_i|\right).$$

Then,

$$
\begin{aligned}
&(\alpha - \beta)\|x - x^\dagger\|_{\ell_1} - (\mathcal{R}_{\alpha,\beta}(x) - \mathcal{R}_{\alpha,\beta}(x^\dagger)) \\
&= -\alpha \sum_{i \in I(x^\dagger)} \left(|x_i| - |x_i^\dagger|\right) + (\alpha - \beta) \sum_{i \in I(x^\dagger)} |x_i - x_i^\dagger| + \beta(T_1 - T_2),
\end{aligned} \tag{2.28}
$$

where

$$T_1 = \left(\sum_i |x_i|^2\right)^{\frac{1}{2}} - \left(\sum_{i \notin I(x^\dagger)} |x_i|^2\right)^{\frac{1}{2}} - \left(\sum_{i \in I(x^\dagger)} |x_i^\dagger|^2\right)^{\frac{1}{2}},$$

$$T_2 = \sum_{i \notin I(x^\dagger)} |x_i| - \left(\sum_{i \notin I(x^\dagger)} |x_i|^2\right)^{\frac{1}{2}}.$$

Observe that $T_2 \geqslant 0$. Since

$$\left( \sum_i |x_i|^2 \right)^{\frac{1}{2}} \leqslant \left( \sum_{i \in I(x^\dagger)} |x_i|^2 \right)^{\frac{1}{2}} + \left( \sum_{i \notin I(x^\dagger)} |x_i|^2 \right)^{\frac{1}{2}},$$

we see that

$$T_1 \leqslant T_3 := \left( \sum_{i \in I(x^\dagger)} |x_i|^2 \right)^{\frac{1}{2}} - \left( \sum_{i \in I(x^\dagger)} |x_i^\dagger|^2 \right)^{\frac{1}{2}}. \tag{2.29}$$

Thus, from (2.28),

$$(\alpha - \beta)\|x - x^\dagger\|_{\ell_1} \leqslant \mathcal{R}_{\alpha,\beta}(x) - \mathcal{R}_{\alpha,\beta}(x^\dagger) + \alpha \sum_{i \in I(x^\dagger)} |x_i - x_i^\dagger|$$
$$+ (\alpha - \beta) \sum_{i \in I(x^\dagger)} |x_i - x_i^\dagger| + \beta T_3. \tag{2.30}$$

Let $m_1$ be a constant upper bound of the terms of the form $|x_i| + |x_i^\dagger|$ and let $m_2 = \left( \sum_{i \in I(x^\dagger)} |x_i^\dagger|^2 \right)^{\frac{1}{2}}$. Then $0 < m_2 \leqslant \left( \sum_{i \in I(x^\dagger)} |x_i|^2 \right)^{\frac{1}{2}} + \left( \sum_{i \in I(x^\dagger)} |x_i^\dagger|^2 \right)^{\frac{1}{2}}$, and

$$T_3 = \frac{\sum_{i \in I(x^\dagger)} (|x_i| - |x_i^\dagger|)(|x_i| + |x_i^\dagger|)}{\left( \sum_{i \in I(x^\dagger)} |x_i|^2 \right)^{\frac{1}{2}} + \left( \sum_{i \in I(x^\dagger)} |x_i^\dagger|^2 \right)^{\frac{1}{2}}} \leqslant \frac{m_1}{m_2} \sum_{i \in I(x^\dagger)} |x_i - x_i^\dagger|. \tag{2.31}$$

A combination of (2.30) and (2.31) shows that

$$(\alpha - \beta)\|x - x^\dagger\|_{\ell_1} \leqslant \mathcal{R}_{\alpha,\beta}(x) - \mathcal{R}_{\alpha,\beta}(x^\dagger) + \left[ 2\alpha - \left( 1 - \frac{m_1}{m_2} \right)\beta \right] \sum_{i \in I(x^\dagger)} |x_i - x_i^\dagger|. \tag{2.32}$$

Furthermore, by assumption 2.16,

$$|x_i - x_i^\dagger| = |\langle e_i, x - x^\dagger \rangle| = |\langle \omega_i, Ax - Ax^\dagger \rangle| \leqslant \max_{i \in I(x^\dagger)} \|\omega_i\|_Y \|Ax - Ax^\dagger\|_Y$$

for all $i \in I(x^\dagger)$. Hence,

$$\sum_{i \in I(x^\dagger)} |x_i - x_i^\dagger| \leqslant |I(x^\dagger)| \max_{i \in I(x^\dagger)} \|\omega_i\|_Y \|Ax - Ax^\dagger\|_Y, \tag{2.33}$$

where $|I(x^\dagger)|$ denotes the size of the index set $I(x^\dagger)$. A combination of (2.32) and (2.33) implies that

$$(\alpha - \beta)\|x - x^\dagger\|_{\ell_1} \leqslant \mathcal{R}_{\alpha,\beta}(x) - \mathcal{R}_{\alpha,\beta}(x^\dagger)$$
$$+ \left[ 2\alpha - \left( 1 - \frac{m_1}{m_2} \right)\beta \right] |I(x^\dagger)| \max_{i \in I(x^\dagger)} \|\omega_i\|_Y \|Ax - Ax^\dagger\|_Y,$$

i.e.

$$(\alpha - \beta)\|x - x^\dagger\|_{\ell_1} \leqslant \mathcal{R}_{\alpha,\beta}(x) - \mathcal{R}_{\alpha,\beta}(x^\dagger) + (c_1\alpha - c_2\beta)\|Ax - Ax^\dagger\|_Y,$$

where $\quad c_1 = 2|I(x^\dagger)|\max_{i\in I(x^\dagger)}\|\omega_i\|_Y, \quad c_2 = \left(1 - \frac{m_1}{m_2}\right)|I(x^\dagger)|\max_{i\in I(x^\dagger)}\|\omega_i\|_Y \quad$ and $c_1\alpha - c_2\beta > 0$. ∎

We comment that the condition $\mathcal{R}_{\alpha,\beta}(x) \leqslant M$ in lemma 2.17 is reasonable in the study of problem (2.1).

**Theorem 2.18 (Convergence rate $O(\delta)$).** *Keep assumption 2.16, let $x_{\alpha,\beta}^\delta$ be defined by (2.1), and let the constants $c_1 > c_2$ be as in lemma 2.17.*

*Case 1. If $q = 1$ and $1 - (c_1\alpha - c_2\beta) > 0$, then*

$$\|x_{\alpha,\beta}^\delta - x^\dagger\|_{\ell_1} \leqslant \frac{1 + (c_1\alpha - c_2\beta)}{(\alpha - \beta)}\delta, \qquad \|Ax_{\alpha,\beta}^\delta - y^\delta\|_Y \leqslant \frac{1 + (c_1\alpha - c_2\beta)}{1 - (c_1\alpha - c_2\beta)}\delta.$$
$$(2.34a)$$

*Case 2. If $q > 1$, then*

$$\|x_{\alpha,\beta}^\delta - x^\dagger\|_{\ell_1} \leqslant \frac{1}{\alpha - \beta}\left[\frac{\delta^q}{q} + (c_1\alpha - c_2\beta)\delta + \frac{(q-1)2^{\frac{1}{q-1}}(c_1\alpha - c_2\beta)^{\frac{q}{q-1}}}{q}\right],$$

$$\|Ax_{\alpha,\beta}^\delta - y^\delta\|_Y^q \leqslant q\left[\frac{\delta^q}{q} + (c_1\alpha - c_2\beta)\delta + \frac{(q-1)2^{\frac{1}{q-1}}(c_1\alpha - c_2\beta)^{\frac{q}{q-1}}}{q}\right].$$
$$(2.34b)$$

**Proof.** Due to the minimization property of $x_{\alpha,\beta}^\delta$, it is clear that

$$\frac{1}{q}\|Ax_{\alpha,\beta}^\delta - y^\delta\|_Y^q + \mathcal{R}_{\alpha,\beta}(x_{\alpha,\beta}^\delta) \leqslant \frac{\delta^q}{q} + \mathcal{R}_{\alpha,\beta}(x^\dagger).$$

Then $\mathcal{R}_{\alpha,\beta}(x_{\alpha,\beta}^\delta)$ is bounded. From lemma 2.17 we see that

$$\frac{\delta^q}{q} \geqslant \mathcal{R}_{\alpha,\beta}(x_{\alpha,\beta}^\delta) - \mathcal{R}_{\alpha,\beta}(x^\dagger) + \frac{1}{q}\|Ax_{\alpha,\beta}^\delta - y^\delta\|_Y^q$$

$$\geqslant (\alpha - \beta)\|x_{\alpha,\beta}^\delta - x^\dagger\|_{\ell_1} - (c_1\alpha - c_2\beta)\|Ax_{\alpha,\beta}^\delta - Ax^\dagger\|_Y + \frac{1}{q}\|Ax_{\alpha,\beta}^\delta - y^\delta\|_Y^q$$

$$\geqslant (\alpha - \beta)\|x_{\alpha,\beta}^\delta - x^\dagger\|_{\ell_1} - (c_1\alpha - c_2\beta)\|Ax_{\alpha,\beta}^\delta - y^\delta\|_Y - (c_1\alpha - c_2\beta)\delta + \frac{1}{q}\|Ax_{\alpha,\beta}^\delta - y^\delta\|_Y^q.$$
$$(2.35)$$

So if $q = 1$ and $1 - (c_1\alpha - c_2\beta) > 0$, then (2.34a) holds. For the case $q > 1$, we apply Young's inequality $ab \leqslant \frac{a^q}{q} + \frac{b^{q^*}}{q^*}$. We have

$$(c_1\alpha - c_2\beta)\|Ax_{\alpha,\beta}^\delta - y^\delta\|_Y = 2^{\frac{1}{q}}(c_1\alpha - c_2\beta)2^{-\frac{1}{q}}\|Ax_{\alpha,\beta}^\delta - y^\delta\|_Y$$

$$\leqslant \frac{1}{2q}\|Ax_{\alpha,\beta}^\delta - y^\delta\|_Y^q + \frac{(q-1)2^{\frac{1}{q-1}}(c_1\alpha - c_2\beta)^{\frac{q}{q-1}}}{q}.$$
$$(2.36)$$

A combination of (2.35) and (2.36) implies (2.34b). ∎

**Remark 2.19 (*A priori* estimation).** Assume $\beta = \eta\alpha$ for a constant $\eta > 0$. If $\alpha \sim \delta^{q-1}(q > 1)$, then $\|x_{\alpha,\beta}^\delta - x^\dagger\|_{\ell_1} \leqslant c\delta$ for some constant $c > 0$. It also follows that $\|x_{\alpha,\beta}^\delta - x^\dagger\|_{\ell_2} \leqslant c\delta$.

Next we provide a convergence rate result by the discrepancy principle.

**Theorem 2.20 (Discrepancy principle).** *Keep the assumptions of lemma 2.17 and let $x_{\alpha,\beta}^\delta$ be defined by (2.1), where the parameters $\alpha$ and $\beta$ ($\beta = \eta\alpha$) are defined via the discrepancy principle*

$$\delta \leqslant \|Ax_{\alpha,\beta}^\delta - y^\delta\|_Y \leqslant \tau\delta \ (\tau \geqslant 1).$$

*Then*

$$\|x_{\alpha,\beta}^\delta - x^\dagger\|_{\ell_2} \leqslant \frac{(c_1 - c_2\eta)(\tau + 1)\delta}{1 - \eta}.$$

**Proof.** By the definition of $x_{\alpha,\beta}^\delta$, $\alpha$ and $\beta$, we see that

$$\frac{1}{q}\delta^q + \mathcal{R}_{\alpha,\beta}(x_{\alpha,\beta}^\delta) \leqslant \frac{1}{q}\|Ax_{\alpha,\beta}^\delta - y^\delta\|_Y^q + \mathcal{R}_{\alpha,\beta}(x_{\alpha,\beta}^\delta) \leqslant \frac{1}{q}\|Ax^\dagger - y^\delta\|_Y^q + \mathcal{R}_{\alpha,\beta}(x^\dagger). \tag{2.37}$$

Hence $\mathcal{R}_{\alpha,\beta}(x_{\alpha,\beta}^\delta) \leqslant \mathcal{R}_{\alpha,\beta}(x^\dagger)$. It follows from lemma 2.17 that

$$\begin{aligned}
0 \geqslant \mathcal{R}_{\alpha,\beta}(x_{\alpha,\beta}^\delta) - \mathcal{R}_{\alpha,\beta}(x^\dagger) &\geqslant (\alpha - \beta)\|x_{\alpha,\beta}^\delta - x^\dagger\|_{\ell_1} - (c_1\alpha - c_2\beta)\|Ax_{\alpha,\beta}^\delta - Ax^\dagger\|_Y \\
&\geqslant (\alpha - \beta)\|x_{\alpha,\beta}^\delta - x^\dagger\|_{\ell_1} - (c_1\alpha - c_2\beta)(\tau + 1)\delta.
\end{aligned} \tag{2.38}$$

Then

$$\|x_{\alpha,\beta}^\delta - x^\dagger\|_{\ell_2} \leqslant \|x_{\alpha,\beta}^\delta - x^\dagger\|_{\ell_1} \leqslant \frac{(c_1\alpha - c_2\beta)(\tau + 1)\delta}{\alpha - \beta}.$$

The theorem is proven with $\beta = \eta\alpha$. ∎

## 3. Computational approach

In this section, we introduce an algorithm to solve problem (1.3) and study its convergence property. We will adapt the generalized conditional gradient method ([22, 23]) and show that this algorithm can be applied to minimize the functional with the non-convex and non-smooth regularization term $\alpha\|x\|_{\ell_1} - \beta\|x\|_{\ell_2}$, $\alpha > \beta \geqslant 0$.

### 3.1. Generalized conditional gradient method

For the sake of completeness, we start with a short description of the generalized conditional gradient method. The starting point for this part is [23], where a generalized conditional gradient method is proposed for solving minimization problems of the form

$$\min_{x \in X} F(x) + \Phi(x)$$

where $X$ is a Hilbert space. Assume that the functional $\Phi(x) : X \to \mathbb{R} \cup \{+\infty\}$ is proper, convex, lower semi-continuous and coercive:

**Condition 3.1.**

1. $\Phi(x) < +\infty$ for some $x \in X$.
2. $\Phi(sx + (1 - s)y) \leqslant s\Phi(x) + (1 - s)\Phi(y)$ for all $x, y \in X$ and $s \in [0, 1]$.
3. $\Phi(x) \leqslant \liminf \Phi(x^k)$ whenever $\lim x^k = x$ in $X$.
4. $\Phi(x)/\|x\| \overset{k \to +\infty}{\to} +\infty$ whenever $\|x\|^k \overset{k \to +\infty}{\to} +\infty$.

The generalized conditional gradient method from [23] is stated in the form of algorithm 1.

---

**Algorithm 1.** Generalized conditional gradient method.

---

1: Set $k = 0$, $x_0 \in X$ such that $\Phi(x_0) < +\infty$.

2: Determine a descent direction $z^k$ as a solution of

$$\min_{z \in X} \langle F'(x^k), z \rangle + \Phi(z).$$

3: Determine a step size $s^k$ as a solution of

$$\min_{s \in [0,1]} F(x^k + s(z^k - x^k)) + \Phi(x^k + s(z^k - x^k)).$$

4: $x^{k+1} = x^k + s_k(z^k - x^k)$, and $k = k + 1$, return to step 2.

---

We recall a convergence result on algorithm 1 proved in [23].

**Theorem 3.2.** *Let $\Phi$ satisfy condition 3.1 and assume $E_t = \{x \in X : \Phi(u) \leqslant t\}$ is compact for every $t \in \mathbb{R}$. Furthermore, let F be a continuously Fréchet differentiable functional, which is bounded on bounded sets with $F + \Phi$ weakly coercive, i.e. $\|x\| \to +\infty \Rightarrow F(x) + \Phi(x) \to +\infty$, and assume $x_0 \in X$ with $\Phi(x_0) < +\infty$. Let $\{x_n\}$ be the sequence generated by the generalized conditional gradient method. Then $\{x_n\}$ contains a convergent subsequence, and every convergent subsequence of $\{x_n\}$ converges to a stationary point of the functional $F + \Phi$.*

### 3.2. Generalized conditional gradient method for a non-convex sparsity regularization

For the sake of convenience, we only consider the case $q = 2$ in (1.3). Since $\mathcal{R}_{\alpha,\beta}(x) := \alpha\|x\|_{\ell_1} - \beta\|x\|_{\ell_2}$, $\alpha > \beta \geqslant 0$ is non-convex, the generalized conditional gradient method cannot be applied to problem (1.3) directly. We rewrite $\mathcal{J}_{\alpha,\beta}^{\delta}(x)$ in (1.3) as

$$\mathcal{J}_{\alpha,\beta}^{\delta}(x) = F(x) + \Phi(x), \tag{3.1}$$

where $F(x) = \frac{1}{2}\|Ax - y^{\delta}\|_Y^2 - \Theta(x)$, $\Phi(x) = \Theta(x) + \alpha\|x\|_{\ell_1} - \beta\|x\|_{\ell_2}$, $\Theta(x) = \frac{\lambda}{2}\|x\|_{\ell_2}^2 + \beta\|x\|_{\ell_2}$ and $\lambda > 0$. There are two reasons why we propose $\Theta(x) = \frac{\lambda}{2}\|x\|_{\ell_2}^2 + \beta\|x\|_{\ell_2}$. First, $\Phi(x) = \frac{\lambda}{2}\|x\|_{\ell_2}^2 + \alpha\|x\|_{\ell_1}$ has certain desirable properties, for example, it is proper, convex, lower semi-continuous and coercive. Another reason is that the iterative ST algorithm can be applied to the minimization of (3.1) directly.

Now we examine the minimization problem in the second step of algorithm 1. The Fréchet derivative of $F(x)$ is given by

$$F'(x) = A^*(Ax - y^{\delta}) - \lambda x - \frac{\beta x}{\|x\|_{\ell_2}}.$$

The minimization problem for determining a descent direction $z^k$ is given by

$$\min_z \langle A^*(Ax^k - y^{\delta}) - \lambda x^k - \frac{\beta x^k}{\|x^k\|_{\ell_2}}, z \rangle + \frac{\lambda}{2}\|z\|_{\ell_2}^2 + \alpha\|z\|_{\ell_1}. \tag{3.2}$$

The minimizer of (3.2) can be calculated explicitly componentwise. The component $z_i$ has to satisfy

$$z_i + \frac{\alpha}{\lambda}\text{sign}(z_i) = \left(x^k + \frac{\beta x^k}{\lambda\|x^k\|_{\ell_2}} - \lambda^{-1}A^*(A(x^k) - y^\delta)\right)_i. \tag{3.3}$$

The solution of (3.3) can be expressed by the ST function $\mathbb{S}_{\alpha/\lambda}$ and $S_{\alpha/\lambda}$, where $\mathbb{S}_{\alpha/\lambda}(x)$ is defined by

$$\mathbb{S}_{\alpha/\lambda}(x) = \sum_i S_{\frac{\alpha}{\lambda}}(x_i)e_i \tag{3.4}$$

and $S_{\alpha/\lambda}(t)$, $t \in \mathbb{R}$ is defined by

$$S_{\alpha/\lambda}(t) = \begin{cases} t - \frac{\alpha}{\lambda} & \text{if } t \geqslant \frac{\alpha}{\lambda}, \\ 0 & \text{if } |t| < \frac{\alpha}{\lambda}, \\ t + \frac{\alpha}{\lambda} & \text{if } t \leqslant -\frac{\alpha}{\lambda}. \end{cases} \tag{3.5}$$

**Lemma 3.3.** *If $x^k \neq 0$, then the minimizer of problem (3.2) is given by*

$$z^k = \mathbb{S}_{\alpha/\lambda}\left(\left(\frac{\beta}{\lambda\|x^k\|_{\ell_2}} + 1\right)x^k - \frac{1}{\lambda}A^*(Ax^k - y^\delta)\right). \tag{3.6}$$

**Proof.**　The proof is similar to that of lemma 2.3 in [22]. Problem (3.2) is equivalent to the problem

$$\min_z \sum_i \frac{\lambda}{2}\left|z_i - \left(x^k + \frac{\beta x^k}{\lambda\|x^k\|_{\ell_2}} - \lambda^{-1}A^*(A(x^k) - y^\delta)\right)_i\right|^2 + \alpha|z_i|. \tag{3.7}$$

From a result in [30, chapter 10], for every proper convex $g : \mathbb{R} \to \mathbb{R}$ and every $\lambda > 0$,

$$(I + \frac{1}{\lambda}\partial(\alpha\|\cdot\|_{\ell_1}))^{-1}(x) = \arg\min_\omega\left\{\frac{\lambda}{2}|\omega - x|^2 + g(\omega)\right\}.$$

Then the minimizer $z^k$ is given by

$$z^k = \sum_i \left[(I + \frac{1}{\lambda}\partial(\alpha\|\cdot\|_{\ell_1}))^{-1}\left(x^k + \frac{\beta}{\lambda\|x^k\|_{\ell_2}} - \lambda^{-1}A^*(A(x^k) - y^\delta)\right)_i\right] \cdot e_i. \tag{3.8}$$

Using the definition (3.4) and (3.5), we can rewrite (3.8) in the form of (3.6). ∎

If $\beta = 0$, (3.6) reduces to the standard ST iteration. We note that the functional $\|x\|_{\ell_2}$ is differentiable at $x \neq 0$ with gradient $x/\|x\|_{\ell_2}$, and is not differentiable at $x = 0$ where the sub-differential contains the element 0. We see that $F(x)$ fails to satisfy the smoothness condition required in the generalized conditional gradient method. Thus we formulate a strategy where the iteration is divided into two steps. We summarize the strategy (ST-$(\alpha\ell_1 - \beta\ell_2)$ algorithm) in algorithm 2.

---

**Algorithm 2.** ST-$(\alpha\ell_1 - \beta\ell_2)$ algorithm for problem (1.3).

---

Set $k = 0$, $x_0 \in X$ such that $\Phi(x_0) < +\infty$,

for $k = 0, 1, 2, \cdots$, do

   If $x^k = 0$ then

   $x^{k+1} = \arg\min \frac{1}{2}\|F(x) - y^\delta\|_Y^2 + \alpha\|x\|_{\ell_1}$

   else

     Determine a descent direction $z^k$ by

$$z^k = \mathbb{S}_{\alpha/\lambda}\left(\left(\frac{\beta}{\lambda\|x^k\|_{\ell_2}} + 1\right)x^k - \frac{1}{\lambda}A^*(Ax^k - y^\delta)\right)$$

     Determine a step size $s^k$ as a solution of

$$\min_{s\in[0,1]} F(x^k + s(z^k - x^k)) + \Phi(x^k + s(z^k - x^k))$$

   $x^{k+1} = x^k + s^k(z^k - x^k)$

   end if

   $k = k + 1$

end for

---

We now turn to the convergence properties of the two-step generalized conditional gradient algorithms. The first order necessary condition of problem (3.1) is (see [23, lemma 1])

$$x \in \ell_2: \quad \langle F'(x), y - x \rangle \geqslant \Phi(x) - \Phi(y) \quad \forall y \in \ell_2. \tag{3.9}$$

**Lemma 3.4.** *Suppose $x^k$ does not fulfill the first order optimality conditions (3.9). Then algorithm 2 determines an $x^{k+1}$ such that*

$$\mathcal{J}_{\alpha,\beta}^\delta(x^{k+1}) = F(x^{k+1}) + \Phi(x^{k+1}) \leqslant F(x^k) + \Phi(x^k) = \mathcal{J}_{\alpha,\beta}^\delta(x^k).$$

**Proof.** If $x^k = 0$, from algorithm 2 we see that

$$\begin{aligned}
\mathcal{J}_{\alpha,\beta}^\delta(x^{k+1}) &= F(x^{k+1}) + \Phi(x^{k+1}) \\
&= \frac{1}{2}\|Ax^{k+1} - y^\delta\|_Y^2 + \alpha\|x^{k+1}\|_{\ell_1} - \beta\|x^{k+1}\|_{\ell_2} \\
&\leqslant \frac{1}{2}\|A0 - y^\delta\|_Y^2 + \alpha\|0\|_{\ell_1} - \beta\|x^{k+1}\|_{\ell_2} \\
&\leqslant \frac{1}{2}\|A0 - y^\delta\|_Y^2 + \alpha\|0\|_{\ell_1} - \beta\|0\|_{\ell_2} \\
&= \mathcal{J}_{\alpha,\beta}^\delta(x^k).
\end{aligned}$$

If $x^k \neq 0$, then $F(x)$ is Fréchet differentiable and $\Phi(x) = \alpha\|x\|_{\ell_1} + \frac{\lambda}{2}\|x\|_{\ell_2}^2$ is proper, convex, lower semi-continuous and coercive. The rest of the proof is similar to that of lemma 2 in [23]. ∎

In order to prove the convergence, we need to analyze the relation between $x^k$ and 0. If $0 = x^0 = x^1$, then we stop the iteration and 0 is the iterative solution. Otherwise, we can see from lemma 3.4 that

$$\mathcal{J}_{\alpha,\beta}^\delta(x^1) - \mathcal{J}_{\alpha,\beta}^\delta(x^0) \leqslant -\beta\|x^1\|_{\ell_2} < 0.$$

Since $\mathcal{J}_{\alpha,\beta}^\delta(x^k)$ decreases, $x^k \neq 0$ for $k \geqslant 1$. So in the following we let $x^k \neq 0$ whenever $k \geqslant 1$.

**Theorem 3.5.**   *Let $\{x^k\}$ denote the sequence generated by algorithm 2. Then $\{x^k\}$ contains a convergent subsequence and every convergent subsequence of $\{x^k\}$ converges to a stationary point of the functional $\mathcal{J}^{\delta}_{\alpha,\beta}(x)$.*

**Proof.**   We apply theorem 3.2 to prove this result. The function $\Phi(x) = \alpha\|x\|_{\ell_1} + \frac{\lambda}{2}\|x\|^2_{\ell_2}$ is weakly lower semi-continuous, and the set $E_t = \{x \in \ell_2 \mid \Phi(x) \leqslant t\}$ is compact for every $t \in \mathbb{R}$. Since $\Phi(0) = \alpha\|0\|_{\ell_1} + \frac{\lambda}{2}\|0\|^2_{\ell_2} = 0 < +\infty$, $\Phi(x)$ is proper. The convexity and coercivity of $\Phi(x)$ follow from the convexity and coercivity of $\ell_1$- and $\ell_2$-norm. We see that $F$ is Fréchet differentiable and

$$F'(x)h = \langle A^*(Ax - y^{\delta}), h \rangle - \Theta'(x)h.$$

Then,

$$\|F'(x) - F'(y)\| \leqslant \|Ax - Ay\|\|A\| + \|\Theta'(x) - \Theta'(y)\| \quad \forall x \neq 0, y \neq 0 \in \ell_2.$$

The continuity of $F'$ follows from the continuity of $A$ and $\Theta'$. It is clear that

$$|F(x)| \leqslant \frac{1}{2}\|Ax\|^2 + \langle Ax, y^{\delta} \rangle + \frac{1}{2}\|y^{\delta}\|^2 + |\Theta(x)|.$$

Since $A$ and $\Theta$ are bounded, $F(x)$ is bounded. By lemma 2.5, $F + \Phi$ is weakly coercive. Then the assumption on $F$ and $\Phi$ in theorem 3.2 is valid and we can apply theorem 3.2. ∎

We comment that if $\beta = 0$, (1.3) reduces to the classical $\ell_1$ sparsity regularization, which implies that the proposed algorithm is a generalization of the classical sparsity regularization. Furthermore, we can extend the above discussion for the solution of a nonlinear ill-posed equation. Meanwhile, if we choose a suitable $\Theta(x)$, the proposed algorithm can be utilized to solve the elastic-net sparsity regularization. This will be implemented in forthcoming papers.

## 4. Numerical experiments

In this section, we present the results from two numerical experiments to demonstrate the efficiency of the proposed method. We analyze the influence of the parameter $\eta$ on the reconstruction of $x^*$ and compare the iterative solutions with that of the classical $\ell_1$ regularization. The classical $\ell_1$ sparsity regularization is as follows

$$\min_x \frac{1}{2}\|Ax - y^{\delta}\|^2_Y + \alpha\|x\|_{\ell_1}.$$

If $\eta = 0$, i.e. $\beta = 0$, (1.3) reduces to the classical $\ell_1$ sparsity regularization, then (3.6) reduces to the classical soft threshold iteration

$$z^k = \mathbb{S}_{\alpha/\lambda}\left(x^k - \frac{1}{\lambda}A^*(Ax^k - y^{\delta})\right).$$

The first example deals with a well-conditioned compressive sensing problem. The second example deals with an ill-conditioned image deblurring problem.

### 4.1. Well-conditioned compressive sensing with random Gaussian matrix

In the first example, we test the commonly used random Gaussian matrix. The compressive sensing problem is defined as $A_{m \times n}x_n = y_m$, where $A_{m \times n}$ is a well-conditioned random

**Table 1.** SNR of reconstruction $x^*$ for different values of $\eta$ and $\alpha$: example 1.

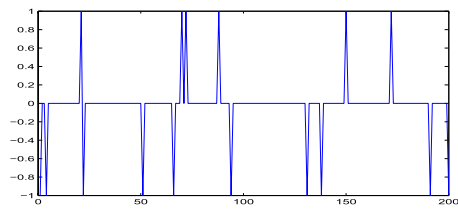| $\eta$ | $\alpha = 4.0 \times 10^{-2}$ | $\alpha = 4.4 \times 10^{-2}$ | $\alpha = 4.8 \times 10^{-2}$ | $\alpha = 5.2 \times 10^{-2}$ | $\alpha = 5.6 \times 10^{-2}$ | $\alpha = 6.0 \times 10^{-2}$ |
|---|---|---|---|---|---|---|
| 0.0 | 21.3059 | 20.6180 | 19.9593 | 19.3292 | 18.7275 | 18.1511 |
| 0.1 | 22.9104 | 22.2370 | 21.5823 | 20.9350 | 20.2948 | 19.6866 |
| 0.2 | 24.3090 | 23.6644 | 23.0212 | 22.3882 | 21.7378 | 21.1106 |
| 0.3 | 25.7979 | 25.1602 | 24.4664 | 23.7725 | 23.1027 | 22.4468 |
| 0.4 | 27.3444 | 26.6692 | 25.9810 | 25.2722 | 24.5976 | 23.9564 |
| 0.5 | 28.9246 | 28.3115 | 27.6261 | 26.9644 | 26.3282 | 25.7171 |
| 0.6 | 30.5250 | 30.0179 | 29.4418 | 28.8749 | 28.3165 | 27.7694 |
| 0.7 | 31.8789 | 31.6224 | 31.2449 | 30.8514 | 30.4465 | 30.0342 |
| 0.8 | 32.5463 | 32.6080 | 32.4561 | 32.2891 | 32.1079 | 31.9132 |
| 0.9 | 32.2278 | 32.3678 | 32.2995 | 32.2128 | 32.1232 | 32.0307 |
| 1.0 | 4.5223 | 31.0753 | 30.8970 | 30.6874 | 30.4755 | 0.9429 |

Gaussian matrix by calling $A = \text{randn}(m, n)$ in MATLAB. The exact data $y^\dagger$ is generated by $y^\dagger = Ax^\dagger$. The exact solution $x^\dagger$ is an $s$-sparse signal supported on a random index set. White Gaussian noise is added to the exact data $y^\dagger$ by calling $y^\delta = \text{awgn}(Ax^\dagger, \delta)$ in MATLAB, where $\delta$ is the noise level, measured in dB, which measures the ratio between the true (noise free) data $y^\dagger$ or $Ax^\dagger$ and Gaussian noise. $x^*$ denotes the reconstruction computed by the proposed algorithm. We use a signal-to-noise ratio (SNR) to evaluate the performance of reconstruction $x^*$, where SNR is defined by

$$\text{SNR} := -10 \log_{10} \frac{\|x^* - x^\dagger\|_{\ell_2}^2}{\|x^\dagger\|_{\ell_2}^2}.$$

We choose $n = 200$, $m = 0.4n$, $s = 0.2m$. The value of $\|A_{m \times n}\|_2$ is around 22 and the condition number of $A_{m \times n}$ is around 4. We rescale the matrix $A_{m \times n}$ by $A_{m \times n} \to 0.05A_{m \times n}$. The 2-norm of the rescaled matrix is around 0.8. Note that the condition number does not change under the matrix rescaling. We let $\lambda = 0.2$, step size $s^k = 1$ and the maximum number of iterations maxiter = 1000. The initial value $x^0$ is generated by calling $x^0 = \text{ones}(n, 1)$.

In section 2.3, we use an *a priori* rule or discrepancy principle to choose the regularization parameter $\alpha$. The *a priori* rule requires $\alpha = O(\delta)$. However, for the numerical implementation, it is difficult to find a good estimate for the optimal value of $\alpha$. In this section, we utilize the discrepancy principle to determine the regularization parameter $\alpha$ such that the residual norm for the regularized solution satisfies $\|Ax^* - y\|_Y = \delta$. When a good estimate for the noise level $\delta$ is known, this method yields a good regularization parameter.

The regularization parameter $\alpha$ is determined by calling $\alpha = \text{discrep}(U, d, V, y, \delta, x^0)$ in MATLAB regularization tools ([31]). Here, $x^0$ is an initial estimate of the solution, $\delta$ is the noise level, $y$ is the observed data, whereas $U$, $d$ and $V$ are the results of the singular value decomposition of $A$ by calling $[U, d, V] = \text{csvd}(A)$ in MATLAB regularization tools. Note that we choose the regularization parameter with the help of Hansen's MATLAB tools. If a regularization parameter $\alpha$ determined by Hansen's MATLAB tools does not satisfy the discrepancy principle, we try $\alpha_j = \frac{\alpha}{2^j}$, $j = 1, 2, \cdots$. As $j$ increases, we calculate $x_{\alpha,\beta}^\delta$ until the regularization parameter satisfies the discrepancy principle. In the numerical experiments of this paper, we determine the regularization parameter using Hansen MATLAB tools through the above procedure—see [32, 33]. The parameter $\alpha$ determined by the discrepancy principle is only an estimate of the optimal regularization parameter. To test the sensitivity of algorithm 2 with respect to $\alpha$, we choose several different regularization parameters in tables 1 and 5.
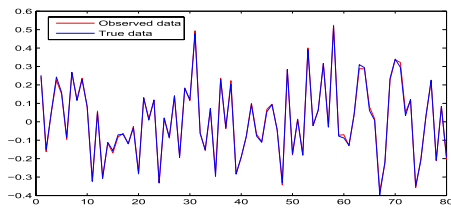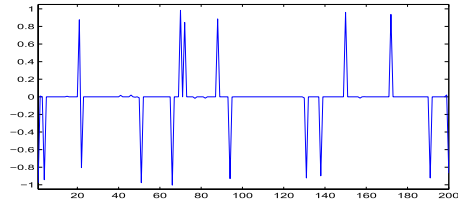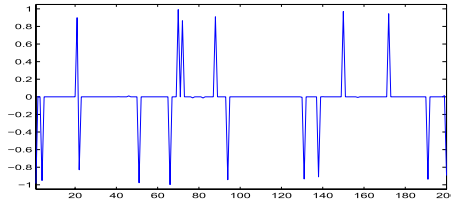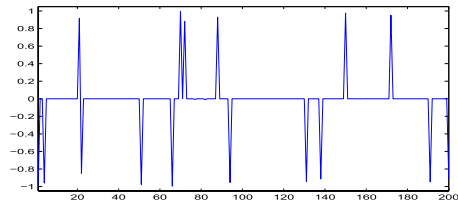
(a) True signal



(b) Observed data ($\delta$=40dB)



(c) $\eta = 0.0$, SNR=19.6401



(d) $\eta = 0.1$, SNR=21.2633
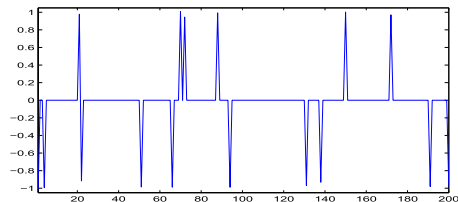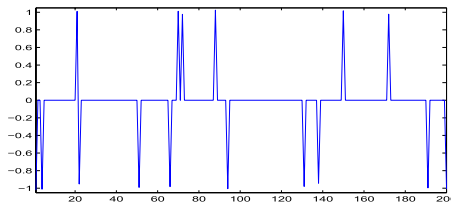


(e) $\eta = 0.2$, SNR=22.7077



(f) $\eta = 0.4$, SNR=25.6228



(g) $\eta = 0.6$, SNR=29.1575



(h) $\eta = 0.8$, SNR=32.3745



(i) $\eta = 0.9$, SNR=32.2565



(j) $\eta = 1.0$, SNR=30.7925

**Figure 2.** (a) True signal. (b) Observed data. (c)–(j) The recovered signal with different $\eta$ at a fixed regularization parameter $\alpha = 5.0 \times 10^{-2}$.

**Table 2.** SNR of reconstruction $x^*$ with different fixed step sizes.

| $s^k$ | 0.3 | 0.6 | 0.9 | 1.2 | 1.5 | 1.8 | 2.1 | 2.4 | 2.7 | 3 |
|---|---|---|---|---|---|---|---|---|---|---|
| SNR | 32.1784 | 32.1784 | 32.1784 | 32.1784 | 32.1784 | 32.1784 | NaN | NaN | NaN | NaN |

Note that the discrepancy principle requires a good estimate of the noise level $\delta$. In the numerical experiments, the added noise is an artificial Gaussian noise, so we do have a good estimate for noise level $\delta$. However, in practical applications, the noise level $\delta$ is not available exactly and only the observed data $y^\delta$ is known. One cannot obtain a good estimate for the noise level $\delta$. In this situation, a reasonable regularization parameter choice rule is based on a heuristic method, for example, L-curve, generalized cross-validation and quasi-optimality criterion. For details and related Matlab codes, see [31] and references therein.

In the first test, a noise $\delta$ is added to exact data $y^\dagger$ by calling $y^\delta = \mathrm{awg}n(Ax^\dagger, \delta)$, with the noise level $\delta = 40$ dB. In order to analyze the influence of $\eta$, we choose different values for the parameters $\eta$ and $\alpha$. Table 1 shows that the proposed algorithm performs well with the appropriate regularization parameters. From each column in table 1, we see that, for a fixed regularization parameter $\alpha$, the results of the reconstruction improve as $\eta$ increases, which implies that the non-convex regularization (for $\eta > 0$) performs better than classical $\ell_1$ regularization (for $\eta = 0$). However, with a large $\eta$ close to 1, the accuracy of recovery decreases and $\eta = 0.8$ is optimal. Meanwhile, too large or small $\alpha$ will lead to divergence when $\eta = 1$. It shows that the case $\eta = 1$ is not stable corresponding to regularization parameter $\alpha$. When $\alpha = 5.6 \times 10^{-2}$ or $6.0 \times 10^{-2}$, the optimal $\eta$ is 0.9. However, the accuracy is worse than the optimal case $\alpha = 4.8 \times 10^{-2}$. Figure 2 shows the graphs of the reconstruction $x^*$ when regularization parameter $\alpha = 5 \times 10^{-2}$.

Note that in algorithm 2, $s^k$ is determined by an optimization problem. However, we let $s^k = 1$ in the numerical experiments. There are two reasons why we chose the step size $s^k$ to be a fixed constant. First, the algorithm is easy to implement for a fixed step size $s^k$. Moreover, it is proved in [22] that under some additional assumptions, the generalized conditional gradient method is convergent with a fixed step size $s^k = 1$. In table 2, we set $\alpha = 5 \times 10^{-2}$, $\eta = 0.8$ and check the convergence of algorithm 2 with different fixed step sizes. Table 2 shows that there exists a threshold $s > 0$, and algorithm 2 does not converge for any fixed step size $s^k > s$. Algorithm 2 converges when $s^k < s$ and it provides the same inversion results. In algorithm 2, we require step size $s^k \in [0, 1]$. However, table 2 shows that $s^k$ can be chosen larger than 1, which implies that one can propose an accelerating version for algorithm 2.

Next we examine the effect of the parameter $\lambda$. Theoretically, (1.3) is the same as (3.1) for any $\lambda$, implying that the inversion results do not change with respect to $\lambda$. However, from the perspective of computation, a small value of $\lambda$ admits a larger $\frac{\beta x^k}{\lambda \|x^k\|_{\ell_2}}$ in (3.6). Indeed, $\|\frac{\beta x^k}{\lambda \|x^k\|_{\ell_2}}\|_{\ell_2} = \frac{\beta}{\lambda} \to \infty$ as $\lambda \to 0$, which leads to divergence. On the other hand, a larger value of $\lambda$ admits a smaller value of the threshold $\frac{\alpha}{\lambda}$. The value of the threshold is a crucial factor for the iterative ST algorithm. A small threshold value leads to divergence. In [34], Daubechies *et al* provided a choice of the threshold. However, for the present paper, we do not have a formula to determine the optimal $\lambda$. So, we provide table 3, which provides a clue as to how to choose a reasonable $\lambda$. In table 3, we set $\alpha = 5 \times 10^{-2}$, $\eta = 0.8$, $s^k = 1$ and provide the reconstruction results of algorithm 2 with different $\lambda$. It is shown that $0.15 \leqslant \lambda \leqslant 0.40$ is a a good choice.

In the second test, we test the stability of the proposed algorithm. Various noise levels $\delta$ are added to the exact data $y^\dagger$. We choose the optimal regularization parameters by the discrepancy principle. The numerical results are shown in table 4. One can clearly see that the SNR of reconstruction $x^*$ decreases as the noise level increases. In the noise-free case, we can obtain a

**Table 3.** SNR of reconstruction $x^*$ with different $\lambda$.

| $\lambda$ | 0.05 | 0.10 | 0.15 | 0.20 | 0.25 | 0.30 | 0.35 | 0.40 | 0.45 |
|---|---|---|---|---|---|---|---|---|---|
| SNR | NaN | 4.54e-12 | 31.0688 | 31.0688 | 31.0688 | 31.0687 | 31.0666 | 31.0500 | 29.9480 |
| $\lambda$ | 0.50 | 0.55 | 0.60 | 0.65 | 0.70 | 0.75 | 0.80 | 0.85 | 0.90 |
| SNR | 28.4993 | 26.2270 | 22.1077 | 16.6312 | 12.0806 | 6.7939 | 4.8894 | 3.2618 | 1.8777 |

**Table 4.** SNR of reconstruction $x^*$ with various noise levels.

| $\delta$ | $\eta = 0$ | $\eta = 0.2$ | $\eta = 0.4$ | $\eta = 0.7$ | $\eta = 0.8$ | $\eta = 0.9$ | $\eta = 1.0$ |
|---|---|---|---|---|---|---|---|
| Noise free, $\alpha = 0.007$ | 38.1346 | 40.6385 | 43.1617 | 45.1865 | 46.6969 | 50.2843 | 52.7298 |
| 50 dB, $\alpha = 0.022$ | 24.9646 | 28.3309 | 31.7742 | 38.5733 | 41.1439 | 42.1468 | 40.4161 |
| 40 dB, $\alpha = 0.046$ | 19.8962 | 23.6644 | 26.6692 | 31.6224 | 32.6080 | 32.3678 | 31.0753 |
| 30 dB, $\alpha = 0.088$ | 12.5288 | 15.3119 | 17.8542 | 19.9705 | 19.7516 | 18.9072 | 0.4470 |
| 20 dB, $\alpha = 1.760$ | NaN | 7.6190 | 3.1342 | 0.8195 | 0.8618 | 0.5943 | 0.4131 |

better performance as $\eta$ increases and $\eta = 1$ is optimal. The results coincide with the theory of the proposed non-convex regularization. Theoretically, in the noiseless case, the fidelity term $\frac{1}{2}\|Ax - y\|^2$ is 0. The regularization term $\mathcal{R}(x) = \alpha(\|x\|_{\ell_1} - \|x\|_{\ell_2}) + (\alpha - \beta)\|x\|_{\ell_2}$ will be minimum if $\alpha = \beta$ ($\eta = 1$). When the noise levels are lower, $\eta$ should be chosen as 0.9 or 0.8. When the noise level is 30 dB, the optimal $\eta$ is 0.7. As the noise level increases, we see that the value of the optimal $\eta$ decreases. When the noise level $\delta$ is 20 dB, the proposed algorithm does not converge except $\eta = 0.2$ and 0.4. Meanwhile, table 4 shows that the method with $\eta = 1$ has poor stability corresponding to the noise level. The algorithm does not converge for the case $\eta = 1$ when the noise level is 30 dB and 20 dB.

Next we discuss how to choose $\eta$ or $\beta$. From figure 1, it can be seen that the effect of $\eta$ is similar to that of the exponent $p$ in $\ell_p$-norm ($0 < p < 1$). Theoretically, $\mathcal{R}_{\alpha,\beta}(x)$ behaves more and more like the $\ell_0$-norm as $\beta/\alpha \to 1$ and we obtain the best recovery results for the noise-free case. In the case of the presence of noise, the situation is more complicated. In [35], a more flexible way of sparse regularization is introduced by varying exponents and it is observed that for $\ell_p$ ($0 < p < 1$) regularization, it is challenging to identify the optimal exponent $p$. The question of how to choose a suitable exponent $p$ is worth further investigation.

From tables 4 and 6, it can be seen that the optimal $\eta$ decreases as the noise level increases. If the noise level is low, a larger value of $\eta$ towards 1 is a reasonable choice. Otherwise, a smaller value of $\eta$ towards 0 is more appropriate.

The two tests show that the proposed non-convex sparsity regularization performs better compared with the $\ell_1$ regularization. Though $\ell_1 - \ell_2$, i.e. the case $\eta = 1$ is a good approximation of $\ell_0$, it is not optimal in the presence of noise—see similar statements in [36, chapter I]. For example, the choice of $\ell_{1.1}$ regularization gives better results than that of $\ell_1$ regularization.

### 4.2. Ill-conditioned image deblurring problem

In the second example, we test the ill-conditioned image deblurring problem which is the process of removing blurring artifacts from images, such as blur caused by defocus aberration or motion blur. The blur is typically modeled by a Fredholm integral equation of the first kind

$$\int_a^b K(s,t)f(t)\mathrm{d}t = g(s),$$

**Table 5.** SNR of reconstruction $x^*$ with different values of $\eta$ and $\alpha$.

| $\eta$ | $\alpha = 3.4 \times 10^{-2}$ | $\alpha = 3.6 \times 10^{-2}$ | $\alpha = 3.8 \times 10^{-2}$ | $\alpha = 4.0 \times 10^{-2}$ | $\alpha = 4.2 \times 10^{-2}$ | $\alpha = 4.4 \times 10^{-2}$ |
|---|---|---|---|---|---|---|
| 0.0 | 32.0565 | 32.0823 | 32.1029 | 32.1227 | 32.1245 | 32.1215 |
| 0.1 | 32.3823 | 32.4225 | 32.4611 | 32.4977 | 32.5139 | 32.5282 |
| 0.2 | 32.6610 | 32.7180 | 32.7739 | 32.8199 | 32.8488 | 32.8757 |
| 0.3 | 32.9108 | 32.9823 | 33.0528 | 33.1069 | 33.1411 | 33.1708 |
| 0.4 | 33.1162 | 33.1898 | 33.2615 | 33.3100 | 33.3494 | 33.3843 |
| 0.5 | 33.2733 | 33.3534 | 33.4295 | 33.4760 | 33.5161 | 33.5442 |
| 0.6 | 33.4032 | 33.4858 | 33.5457 | 33.5826 | 33.6153 | 33.6397 |
| 0.7 | 33.4999 | 33.5609 | 33.6123 | 33.6433 | 33.6649 | 33.6642 |
| 0.8 | 33.5649 | 33.6174 | 33.6390 | 33.6340 | 33.6197 | 33.5987 |
| 0.9 | 33.5624 | 33.5945 | 33.5973 | 33.5703 | 33.5314 | 33.4840 |
| 1.0 | 33.5067 | 33.5154 | 33.4969 | 33.4584 | 33.3987 | 33.3262 |

where $K(s,t)$ is the kernel function, $g(s)$ is the observed image and $f(t)$ is the true image. We utilize the blur problem from MATLAB regularization tools ([31]) by calling $[A, b, x^\dagger] = \mathrm{blur}(n, \mathrm{band}, \sigma)$, where the Gaussian point-spread function is used as the kernel function

$$K(s,t) = \frac{1}{\pi\sigma^2}\exp\left(-\frac{s^2 + t^2}{2\sigma^2}\right).$$

The matrix $A$ is a symmetric $n^2 \times n^2$ Toeplitz matrix and is given by $A = (2\pi\sigma^2)^{-1}T \otimes T$, where $T$ is an $n \times n$ symmetric banded Toeplitz matrix whose first row is given by calling

$$z = [\exp(-([0 : \mathrm{band} - 1].^2)/(2\sigma^2)); \mathrm{zeros}(1, N - \mathrm{band})].$$

We note that the parameter $\sigma$ controls the shape of the Gaussian point spread function and thus the amount of smoothing (the larger value of $\sigma$, the wider the function, and the less ill-posed the problem). We choose $n = 16$, band $= 3$, $\sigma = 0.7$. The value of $\|A\|_2$ is around 1 and the condition number is around 30. We use a similar setting as in section 4.1. We let $\lambda = 0.2$, step size $s^k = 1$ and the maximum number of iterations maxiter $= 500$. The initial value $x^0$ is generated by calling $x^0 = \mathrm{ones}(n, 1)$.

Table 5 shows the performance of reconstruction with the different regularization parameters and $\eta$. As expected, similar results are obtained in the test. The results of reconstruction improve as $\eta$ increases, but the choice $\eta = 1$ is not optimal. However, compared with the results from section 4.1, the performance for the case $\eta = 1$ is more stable corresponding to the regularization parameter $\alpha$. One can also obtain good results even with larger or smaller $\alpha$. Figure 3 shows graphs of the reconstruction $x^*$ when the regularization parameter $\alpha = 4.0 \times 10^{-2}$.

The stability of reconstruction corresponding to various noise levels $\delta$ is illustrated in table 6. We see that the accuracy decreases as the noise level increases. Its stability is better than that in section 4.1. One can obtain stable reconstruction even with a noise level $\delta = 20$ dB. If the noise level is lower, the optimal choice of $\eta$ is close to 1. As the noise level increases, we see that the value of the optimal $\eta$ decreases. When the noise levels are higher, for example, 20 dB or 10 dB, the non-convex regularization does not display an advantage over the classical $\ell_1$ regularization.
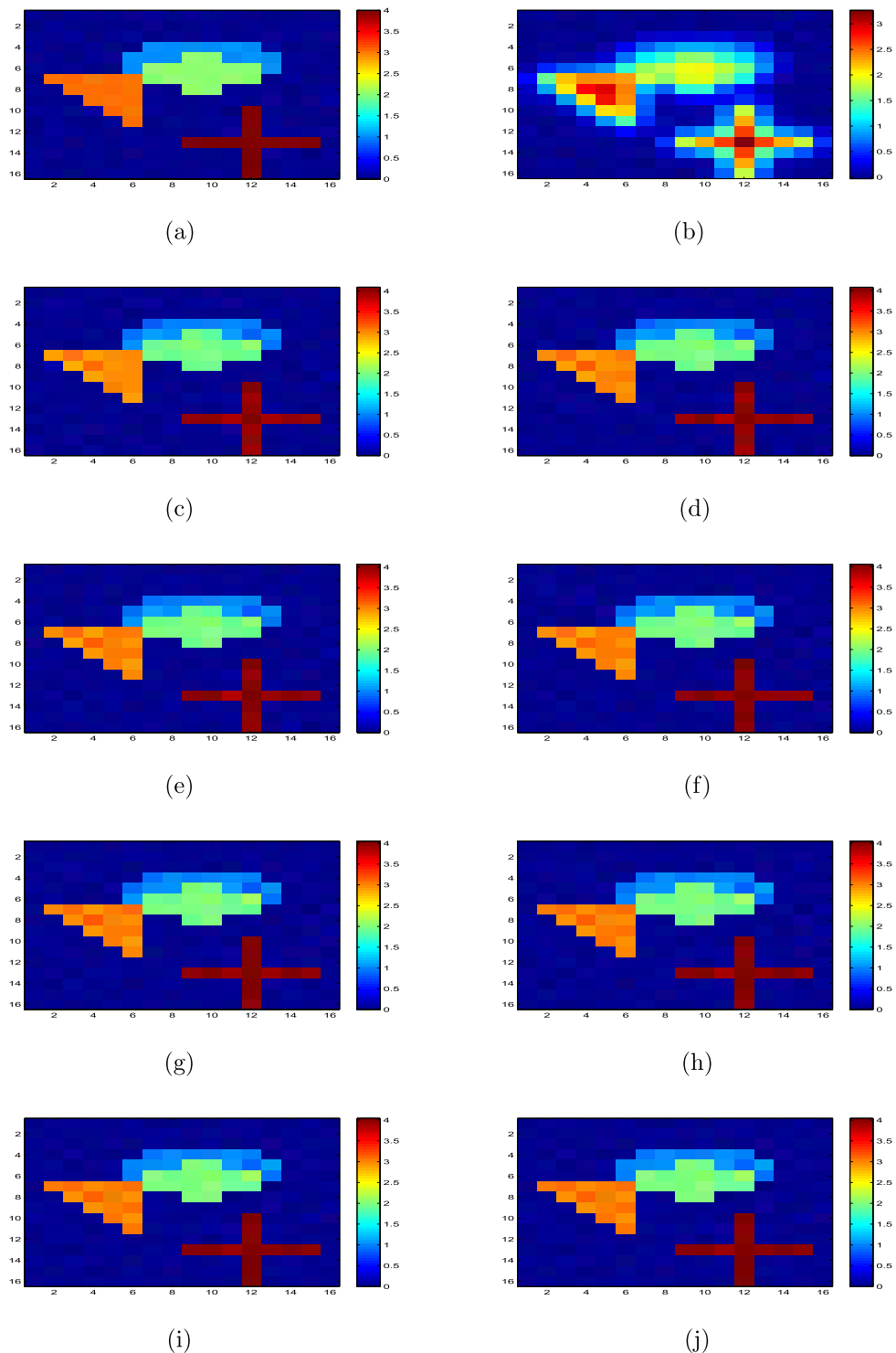
**Figure 3.** True image and its blurred and noisy observation together with reconstructions $x^*$ for $\alpha = 4.0 \times 10^{-2}$ with different $\eta$. (a) True image. (b) Blurred image ($\delta = 40\,\mathrm{dB}$). (c) $\eta = 0.0$, SNR $=32.0232$. (d) $\eta = 0.1$, SNR $=32.3958$. (e) $\eta = 0.2$, SNR $=32.7183$. (f) $\eta = 0.4$, SNR $=33.2029$. (g) $\eta = 0.6$, SNR $=33.4887$. (h) $\eta = 0.8$, SNR $=33.6092$. (i) $\eta = 0.9$, SNR $=33.5664$. (j) $\eta = 1.0$, SNR $=33.4856$.

**Table 6.** SNR of reconstruction $x^*$ with different noise levels and parameter $\eta$.

| $\delta$ | $\eta = 0$ | $\eta = 0.2$ | $\eta = 0.4$ | $\eta = 0.7$ | $\eta = 0.8$ | $\eta = 0.9$ | $\eta = 1.0$ |
|---|---|---|---|---|---|---|---|
| Noise free, $\alpha = 0.004$ | 53.4600 | 54.6527 | 55.6761 | 56.8733 | 57.1769 | 57.4255 | 57.6158 |
| 50 dB, $\alpha = 0.010$ | 43.0779 | 43.0779 | 43.5102 | 43.7562 | 43.7734 | 43.7723 | 43.7524 |
| 40 dB, $\alpha = 0.038$ | 32.1029 | 32.7739 | 32.2615 | 33.6483 | 33.6390 | 33.5973 | 33.4969 |
| 30 dB, $\alpha = 0.060$ | 19.7652 | 20.0685 | 19.4456 | 18.3547 | 17.9399 | 17.5027 | 17.0431 |
| 20 dB, $\alpha = 0.154$ | 14.5627 | 14.1465 | 13.0871 | 9.6716 | 8.5191 | 7.9043 | 7.3759 |
| 10 dB, $\alpha = 0.416$ | 4.7682 | 4.5543 | 3.9252 | 3.3839 | 3.3211 | 3.2126 | 3.0283 |

## 5. Conclusion

We proposed and analyzed a new non-convex $\alpha\ell_1 - \beta\ell_2$ $(\alpha > \beta \geqslant 0)$ regularization method for sparse recovery. The convergence rate $O(\delta)$ was derived under a source condition for both *a priori* and *a posteriori* parameter choice rules. An ST-$(\alpha\ell_1 - \beta\ell_2)$ algorithm was presented based on the generalized conditional gradient method. The critical parameter $\eta$ decreases as the noise level increases. Numerical experiments indicate that with a lower noise level, the proposed algorithm performs better compared with that of the $\ell_1$ regularization whether the operator $A$ is well- or ill-conditioned. If noise levels are higher, for example, 10 dB, the proposed non-convex method is not more advantageous; however, in this case, it is questionable whether any method for the reconstruction problem can yield a practically useful solution due to the high level of noise.

## Acknowledgments

## ORCID iDs

Liang Ding ⓘ https://orcid.org/0000-0003-1543-8614

## References

[1] Daubechies I, Defrise M and De Mol C 2004 An iterative thresholding algorithm for linear inverse problems with a sparsity constraint *Commun. Pure Appl. Math.* **57** 1413–57
[2] Fornasier M (ed) 2010 *Theoretical Foundations and Numerical Methods for Sparse Recovery* (Berlin: de Gruyter & Co)
[3] Jin B and Maass P 2012 Sparsity regularization for parameter identification problems *Inverse Problems* **28** 123001
[4] Bredies K and Lorenz D A 2009 Regularization with non-convex separable constraints *Inverse Problems* **25** 085011
[5] Grasmair M 2010 Non-convex sparse regularisation *J. Math. Anal. Appl.* **365** 19–28
[6] Grasmair M 2010 Generalized Bregman distances and convergence rates for non-convex regularization methods *Inverse Problems* **26** 115014
[7] Nikolova M 2013 Description of the minimizers of least squares regularized with $\ell_0$-norm. Uniqueness of the global minimizer *SIAM J. Imaging Sci.* **6** 904–37

[8] Ito K and Kunisch K 2014 A variational approach to sparsity optimization based on Lagrange multiplier theory *Inverse Problems* **30** 015001

[9] Yin P, Lou Y, He Q and Xin J 2015 Minimization of $\ell_{1-2}$ for compressed sensing *SIAM J. Sci. Comput.* **37** A536–63

[10] Huang X, Shi L and Yan M 2015 Nonconvex sorted $\ell_1$ minimization for sparse approximation *J. Oper. Res. Soc. China* **3** 207–29

[11] Grasmair M 2009 Well-posedness and convergence rates for sparse regularization with sublinear $\ell^q$ penalty term *Inverse Problems Imaging* **3** 383–7

[12] Ramlau R and Zarzer C A 2012 On the minimization of a Tikhonov functional with a non-convex sparsity constraint *Electron. Trans. Numer. Anal.* **39** 476–507

[13] Wang W, Lu S, Mao H and Cheng J 2013 Multi-parameter Tikhonov regularization with the $\ell_0$ sparsity constrain *Inverse Problems* **29** 065018

[14] Zarzer C A 2009 On Tikhonov regularization with non-convex sparsity constraints *Inverse Problems* **25** 025006

[15] Xu Z, Zhang H, Wang Y, Chang X and Liang Y 2010 $L$ regularization *Sci. China: Inf. Sci.* **53** 1159–69

[16] Bredies K, Lorenz D and Reiterer S 2015 Minimization of non-smooth, non-convex functionals by iterative thresholding *J. Optim. Theory Appl.* **165** 78–112

[17] Glowinski R, Osher S and Yin W (ed) 2016 *Splitting Methods in Communication, Imaging, Science, and Engineering* (Berlin: Springer)

[18] Wang Y, Yin W and Zeng J 2019 Global convergence of ADMM in nonconvex nonsmooth optimization *J. Sci. Comput.* **78** 29–63

[19] Esser E, Lou Y and Xin J 2013 A method for finding structured sparse solutions to non-negative least squares problems with applications *SIAM J. Imaging Sci.* **6** 2010–46

[20] Chen D, Hofmann B and Zou J 2017 Elastic-net regularization versus $\ell_1$-regularization for linear inverse problems with quasi-sparse solutions *Inverse Problems* **33** 015004

[21] Jin B, Lorenz D A and Schiffer S 2009 Elastic-net regularization: error estimates and active set methods *Inverse Problems* **25** 115022

[22] Bonesky T, Bredies K, Lorenz D A and Maass P 2007 A generalized conditional gradient method for nonlinear operator equations with sparsity constraints *Inverse Problems* **23** 2041–58

[23] Bredies K, Lorenz D A and Maass P 2009 A generalized conditional gradient method and its connection to an iterative shrinkage method *Comput. Optim. Appl.* **42** 173–93

[24] Amir B and Marc T 2009 A fast iterative shrinkage-thresholding algorithm for linear inverse problems *SIAM J. Imaging Sci.* **2** 183–202

[25] Engl H W, Hanke M and Neubauer A 1996 *Regularization of Inverse Problems* (*Mathematics and its Applications* vol 375) (Kluwer: Dordrecht)

[26] Carothers N L 2000 *Real Analysis* (New York: Cambridge University Press)

[27] Grasmair M, Haltmeier M and Scherzer O 2008 Sparse regularization with $\ell^q$ penalty term *Inverse Problems* **24** 055020

[28] Scherzer O, Grasmair M, Grossauer H, Haltmeier M and Lenzen F 2009 *Variational Methods in Imaging* (*Applied Mathematical Sciences* vol 167) (New York: Springer)

[29] Burger M, Flemming J and Hofmann B 2013 Convergence rates in $\ell_1$-regularization if the sparsity assumption fails *Inverse Problems* **29** 025013

[30] Rockafellar R T and Wets R J-B 1998 *Variational Analysis* (Berlin: Springer)

[31] Hansen P C 2007 Regularization tools version 4.0 for Matlab 7.3 *Numer. Algorithms* **46** 189–94

[32] Anzengruber S T and Ramlau R 2009 Morozov's discrepancy principle for Tikhonov type functionals with nonlinear operators *Inverse Problems* **26** 025001

[33] Wang W, Lu S, Hofmann B and Cheng J 2019 Tikhonov regularization with $\ell_0$-term complementing a convex penalty: $\ell_1$-convergence under sparsity constraints *J. Inverse Ill-Posed Problems* **27** 575–90

[34] Daubechies I, Defrise M and De Mol C 2003 An iterative thresholding algorithm for linear inverse problems with a sparsity constraint in preparation (arXiv:math/0307152)

[35] Lorenz D A and Resmerita E 2017 Flexible sparse regularization *Inverse Problems* **33** 014002

[36] Schuster T, Kaltenbacher B, Hofmann B and Kazimierski K 2012 *Regularization Methods in Banach Spaces* (Berlin: de Gruyter & Co)