

## ITERATIVE $\ell_1$ MINIMIZATION FOR NON-CONVEX COMPRESSED SENSING\*

Penghang Yin

*Department of Mathematics, University of California, Los Angeles, CA 90095, USA*

*Email: yph@ucla.edu*

Jack Xin

*Department of Mathematics, University of California, Irvine, CA 92697, USA*

*Email: jxin@math.uci.edu*

### Abstract

An algorithmic framework, based on the difference of convex functions algorithm (DCA), is proposed for minimizing a class of concave sparse metrics for compressed sensing problems. The resulting algorithm iterates a sequence of  $\ell_1$  minimization problems. An exact sparse recovery theory is established to show that the proposed framework always improves on the basis pursuit ( $\ell_1$  minimization) and inherits robustness from it. Numerical examples on success rates of sparse solution recovery illustrate further that, unlike most existing non-convex compressed sensing solvers in the literature, our method always outperforms basis pursuit, no matter how ill-conditioned the measurement matrix is. Moreover, the iterative  $\ell_1$  (IL<sub>1</sub>) algorithm lead by a wide margin the state-of-the-art algorithms on  $\ell_{1/2}$  and logarithmic minimizations in the strongly coherent (highly ill-conditioned) regime, despite the same objective functions. Last but not least, in the application of magnetic resonance imaging (MRI), IL<sub>1</sub> algorithm easily recovers the phantom image with just 7 line projections.

*Mathematics subject classification:* 90C26, 65K10, 49M29.

*Key words:* Compressed sensing, Non-convexity, Difference of convex functions algorithm, Iterative  $\ell_1$  minimization.

## 1. Introduction

Compressed sensing (CS) techniques [5, 6, 8, 17] enable efficient reconstruction of a sparse signal under linear measurements far less than its physical dimension. Mathematically, CS aims to recover an  $n$ -dimensional vector  $\bar{x} \in \mathbb{R}^n$  with few non-zero components from an under-determined linear system  $Ax = A\bar{x}$  of just  $m \ll n$  equations, where  $A \in \mathbb{R}^{m \times n}$  is a known measurement matrix. The first CS technique is the convex  $\ell_1$  minimization or the so-called basis pursuit [15]:

$$\min_{x \in \mathbb{R}^n} \|x\|_1 \quad \text{s.t.} \quad Ax = A\bar{x}. \quad (1.1)$$

Breakthrough results [8] have established that when matrix  $A$  satisfies certain restricted isometry property (RIP), the solution to (1.1) is exactly  $\bar{x}$ . It was shown that with overwhelming probability, several random ensembles such as random Gaussian, random Bernoulli, and random partial Fourier matrices, are of RIP type [8, 13, 32]. Note that (1.1) is just a minimization principle rather than an algorithm for retrieving  $\bar{x}$ . Algorithms for solving (1.1) and its associated

---

\* Received April 26, 2016 / Revised version received September 23, 2016 / Accepted October 14, 2016 /  
Published online June 1, 2017 /

$\ell_1$  regularization problem [36]:

$$\min_{x \in \mathbb{R}^n} \frac{1}{2} \|Ax - b\|^2 + \lambda \|x\|_1 \quad (1.2)$$

include Bregman methods [24, 43], alternating direction algorithms [3, 18, 40], iterative thresholding methods [1, 14] among others [25].

Inspired by the success of basis pursuit, researchers then began to investigate various non-convex CS models and algorithms. More and more empirical studies have shown that non-convex CS methods usually outperform basis pursuit when matrix  $A$  is RIP-like, in the sense that they require fewer linear measurements to reconstruct signals of interest. Instead of minimizing  $\ell_1$  norm, it is natural to consider minimization of non-convex (concave) sparse metrics, for instance,  $\ell_q$  (quasi-)norm ( $0 < q < 1$ ) [11, 12, 27], capped- $\ell_1$  [30, 45], and transformed- $\ell_1$  [28, 44]. Another category of CS methods in spirit rely on support detection of  $\bar{x}$ . To name a few, there are orthogonal matching pursuit (OMP) [37], iterative hard thresholding (IHT) [2], (re)weighted- $\ell_1$  scheme [7], iterative support detection (ISD) [38], and their variations [26, 31, 46].

On the other hand, it has been proved that even if  $A$  is not RIP-like and contains highly correlated columns, basis pursuit still enables sparse recovery under certain conditions of  $\bar{x}$  involving its support [4]. In this scenario, most of the existing non-convex CS methods, however, are not that robust to the conditioning of  $A$ , as suggested by [41]. Their success rates will drop as columns of  $A$  become more and more correlated. In [41], based on the difference of convex functions algorithm (DCA) [34, 35], the authors propose DCA- $\ell_{1-2}$  for minimizing the difference of  $\ell_1$  and  $\ell_2$  norms [19, 42]. Extensive numerical experiments [29, 30, 41] imply that DCA- $\ell_{1-2}$  algorithm consistently outperforms  $\ell_1$  minimization, irrespective of the conditioning of  $A$ .

Stimulated by the empirical evidence found in [29, 30, 41], we propose a general DCA-based CS framework for the minimization of a class of concave sparse metrics. More precisely, we consider the reconstruction of a sparse vector  $\bar{x} \in \mathbb{R}^n$  by minimizing sparsity-promoting metrics:

$$\min_{x \in \mathbb{R}^n} P(|x|) \quad \text{s.t.} \quad Ax = A\bar{x}. \quad (1.3)$$

Throughout the paper, we assume that  $P(x)$  always takes the form  $\sum_{i=1}^n p(x_i)$  unless otherwise stated, where  $p$  defined on  $[0, +\infty)$  satisfies:

- $p$  is concave and increasing.
- $p$  is continuous with the right derivative  $p'(0+) > 0$ .

The first condition encourages zeros in  $|x|$  rather than small entries, since  $p$  changes rapidly around the origin; the second one is imposed for the good of the proposed algorithm, as will be seen later. A number of sparse metrics in the literature enjoy the above properties, including smoothly clipped absolute deviation (SCAD) [20], capped- $\ell_1$ , transformed- $\ell_1$ , and of course  $\ell_1$  itself. Although  $\ell_q$  ( $q \in (0, 1)$ ) and logarithm functional do not meet the second condition, their smoothed versions  $p(t) = (t + \varepsilon)^q$  and  $p(t) = \log(t + \varepsilon)$  are differentiable at zero. These proposed properties will be essential in the algorithm design as well as in the proof of main results.

Our proposed algorithm calls for solving a sequence of minimization subproblems. The objective of each subproblem is  $\|x\|_1$  plus a linear term, which is convex and tractable. We further validate robustness of this framework, by showing theoretically and numerically that it performs at least as well as basis pursuit in terms of uniform sparse recovery, independent of the conditioning of  $A$  and sparsity metric.

The paper is organized as follows. In Section 2, we overview RIP and coherence of sensing matrices, as well as descent property of DCA. In Section 3, we provide the iterated  $\ell_1$  framework for non-convex minimization, with worked out examples on representative sparse objectives including the total variation. In Section 4, we prove the main exact recovery results based on unique recovery property of  $\ell_1$  minimization instead of RIP, which forms a theoretical basis of the better performance of DCA. In Section 5, we compare iterative  $\ell_1$  algorithms with two state-of-the-art non-convex CS algorithms, IRLS- $\ell_q$  [27] and IRL<sub>1</sub> [7], and ADMM- $\ell_1$ , in CS test problems with varying degree of coherence. We find that iterative  $\ell_1$  outperforms ADMM- $\ell_1$  independent of the sensing matrix coherence, and leads IRLS- $\ell_q$  [27] and IRL<sub>1</sub> [7] in the highly coherent regime. This is consistent with earlier findings of DCA- $\ell_{1-2}$  algorithm [29, 30, 41] to which our theory also applies. We also evaluate these two non-convex metrics on a two-dimensional example of reconstructing MRI from a small number of projections, our iterative  $\ell_1$  algorithm succeed with 7 projections for both metrics. Using the same objective functions, the state-of-the-art algorithms need at least 10 projections. Concluding remarks are in Section 6.

**Notations.** Let us fix some notations. For any  $x, y \in \mathbb{R}^n$ ,  $\langle x, y \rangle = x^T y$  is their inner product.  $\mathbf{0} \in \mathbb{R}^n$  is the vector of zeros, and similar to  $\mathbf{1}$ .  $\circ$  is Hadamard (entry-wise) product, meaning that  $x \circ y = \sum_i x_i y_i$ .  $I_m$  is the identity matrix of dimension  $m$ . For any function  $g$  on  $\mathbb{R}^n$ ,  $\nabla g(x) \in \partial g(x)$  is a subgradient of  $g$  at  $x$ . The  $\text{sgn}(x)$  is the signum function on  $\mathbb{R}^n$  defined as

$$(\text{sgn}(x))_i := \begin{cases} \frac{x_i}{|x_i|} & \text{if } x_i \neq 0, \\ 0 & \text{if } x_i = 0. \end{cases}$$

For any set  $\Omega \subseteq \mathbb{R}^n$ ,  $\iota_\Omega(x)$  is given by

$$\iota_\Omega(x) := \begin{cases} 0 & \text{if } x \in \Omega, \\ \infty & \text{if } x \notin \Omega. \end{cases}$$

## 2. Preliminaries

The well-known CS concept during the past decade is the restricted isometry property (RIP) introduced by Candès *et al.* [8], which is used to characterize matrices that are nearly orthonormal.

**Definition 2.1.** For each number  $s$ ,  $s$ -restricted isometry constant of  $A$  is the smallest  $\delta_s \in (0, 1)$  such that

$$(1 - \delta_s) \|x\|_2^2 \leq \|Ax\|_2^2 \leq (1 + \delta_s) \|x\|_2^2$$

for all  $x \in \mathbb{R}^n$  with sparsity of  $s$ . The matrix  $A$  is said to satisfy the  $s$ -RIP with  $\delta_s$ .

Mutual coherence [15] is another commonly-used concept closely related to the success of CS task.

**Definition 2.2.** The coherence of a matrix  $A$  is the maximum absolute value of the cross-correlations between the columns of  $A$ , namely

$$\mu(A) := \max_{i \neq j} \frac{|A_i^T A_j|}{\|A_i\|_2 \|A_j\|_2}.$$

When matrix  $A$  has small mutual coherence (incoherent) or small RIP constant, its columns tend to be more separated or distinguishable, which is intuitively favorable to identification of the supports of target signal. On the other hand, a highly coherent matrix with large coherence poses challenge to the reconstruction.

Next we give a brief review on the difference of convex functions algorithm (DCA). DCA has been widely applied to sparse optimization problems in several works [19, 23, 28, 30, 41]. For an objective function  $F(x) = G(x) - H(x)$  on the space  $\mathbb{R}^n$ , where  $G(x)$  and  $H(x)$  are lower semicontinuous proper convex functions, we call  $G - H$  a DC decomposition of  $F$ .

DCA takes the following form

$$\begin{cases} y^{(k)} \in \partial H(x^{(k)}) \\ x^{(k+1)} = \arg \min_{x \in \mathbb{R}^n} G(x) - (H(x^{(k)}) + \langle y^k, x - x^k \rangle) \end{cases}$$

Since  $y^{(k)} \in \partial H(x^{(k)})$ , by the definition of subgradient, we have

$$H(x^{k+1}) \geq H(x^k) + \langle y^k, x^{k+1} - x^k \rangle.$$

Consequently,

$$G(x^{(k)}) - H(x^{(k)}) \geq G(x^{(k+1)}) - (H(x^{(k)}) + \langle y^k, x^{(k+1)} - x^{(k)} \rangle) \geq G(x^{(k+1)}) - H(x^{(k+1)}).$$

The fact that  $x^{(k+1)}$  minimizes  $G(x) - (H(x^{(k)}) + \langle y^k, x - x^{(k)} \rangle)$  was used in the first inequality above. Therefore, DCA permits a decreasing sequence  $\{F(x^{(k)})\}$ , leading to its convergence provided  $F(x)$  is bounded from below.

### 3. Iterative $\ell_1$ Framework

Our proposed iterative  $\ell_1$  framework for solving (1.3) is built on  $\ell_1$  minimization and DCA. Note that (1.3) can be equivalently written as

$$\min_{x \in \mathbb{R}^n} P(|x|) + \iota_{\{x: Ax=A\bar{x}\}}(x).$$

We then rewrite the above objective in DC decomposition form:

$$P(|x|) + \iota_{\{x: Ax=A\bar{x}\}}(x) = (p'(0+)\|x\|_1 + \iota_{\{x: Ax=A\bar{x}\}}(x)) - (p'(0+)\|x\|_1 - \sum_{i=1}^n p(|x_i|))$$

Clearly the first term on the right-hand side is convex in terms of  $x$ . We show below that the second term is also a convex function.

**Proposition 3.1.**  $p'(0+)\|x\|_1 - \sum_{i=1}^n p(|x_i|)$  is convex in  $x$ .

*Proof.* For notational convenience, define  $f(t) := p'(0+)t - p(t)$  on  $[0, \infty)$ . Since  $p$  is concave on  $[0, \infty)$ , we have that  $f$  is convex on  $[0, \infty)$ . We only need to show that  $f(|\cdot|)$  is convex on  $\mathbb{R}$ , or equivalently, for all  $t_1, t_2 \in \mathbb{R}$ ,  $a \in (0, 1)$ ,

$$f(|at_1 + (1-a)t_2|) \leq af(|t_1|) + (1-a)f(|t_2|).$$

**Case 1.** If  $t_1$  and  $t_2$  have the same sign or one of them is 0. Since  $f(|at_1 + (1-a)t_2|) = f(a|t_1| + (1-a)|t_2|)$  and  $f$  is convex on  $[0, \infty)$ , then the above inequality holds.

**Case 2.** If  $t_1$  and  $t_2$  are of the opposite sign. By the concavity of  $p$  on  $[0, \infty)$ , we have

$$p(t) \leq p(0) + p'(0+)t, \quad \forall t > 0,$$

that is,  $f(t) \geq f(0)$  for all  $t > 0$ . Without loss of generality, we suppose  $a|t_1| \geq (1-a)|t_2|$ . Then

$$\begin{aligned} f(|at_1 + (1-a)t_2|) &= f(a|t_1| - (1-a)|t_2|) \\ &\leq \frac{(1-a)(|t_1| + |t_2|)}{|t_1|} f(0) + \frac{a|t_1| - (1-a)|t_2|}{|t_1|} f(|t_1|) \\ &\leq (1-a)f(|t_2|) + \frac{(1-a)|t_2|}{|t_1|} f(|t_1|) + \frac{a|t_1| - (1-a)|t_2|}{|t_1|} f(|t_1|) \\ &= af(|t_1|) + (1-a)f(|t_2|) \end{aligned}$$

In the first inequality above, we used the convexity of  $f$  on  $[0, \infty)$ , whereas in the second one, we used the fact that  $f(t) \geq f(0)$  for  $t > 0$ .  $\square$

At the  $(k+1)^{\text{th}}$  iteration, DCA calls for linearization of the second convex term at the current guess  $x^{(k)}$ , and solving the resulting convex subproblem for  $x^{(k+1)}$ . After converting back the linear constraint and removing the constant and the factor of  $p'(0+)$ , we iterate:

$$x^{(k+1)} = \arg \min_x \|x\|_1 - \langle R(x^{(k)}), x \rangle \quad \text{s.t.} \quad Ax = A\bar{x}, \quad (3.1)$$

where

$$R(x) := \text{sgn}(x) \circ \left( \mathbf{1} - \frac{P'(|x|)}{p'(0+)} \right) \in \partial(\|\cdot\|_1 - \frac{P(|\cdot|)}{p'(0+)})(x).$$

Be aware that  $P'(|x|) \in \partial P(\cdot)(|x|)$  denotes subgradient of  $P$  at  $|x|$  rather than subgradient of  $P(|\cdot|)$  at  $x$ . In this way, the subproblem reduces to minimizing  $\|x\|_1$  plus a linear term of  $x$ , which can be efficiently solved by a variety of state-of-the-art algorithms for basis pursuit (with minor modifications). In Table 3.1, we list some non-convex metrics and the corresponding iterative  $\ell_1$  algorithm.

Table 3.1: Examples of sparse metrics and associated iterative  $\ell_1$  scheme.

sparse metric	$p(t)$	$p'(0+)$	$(R(x))_i$
capped- $\ell_1$	$\min\{t, \theta\}, \theta > 0$	1	$\text{sgn}(x_i) \iota_{ x_i  \geq \theta}$
transformed- $\ell_1$	$\frac{(\theta+1)t}{t+\theta}, \theta > 0$	$\frac{\theta+1}{\theta}$	$\text{sgn}(x_i) (1 - (\frac{\theta}{ x_i +\theta})^2)$
smoothed log	$\log(t + \varepsilon), \varepsilon > 0$	$\frac{1}{\varepsilon}$	$\text{sgn}(x_i) (1 - \frac{\varepsilon}{ x_i +\varepsilon})$
smoothed $\ell_q$	$(t + \varepsilon)^q, \varepsilon > 0$	$q\varepsilon^{q-1}$	$\text{sgn}(x_i) (1 - (\frac{\varepsilon}{ x_i +\varepsilon})^{1-q})$

For initialization, we take  $x^{(0)} = R(x^{(0)}) = 0$ , which is basically  $\ell_1$  minimization. The proposed algorithm is thus summarized in Algorithm 3.1 below. Due to the descending property of DCA, Algorithm 3.1 produces a convergent sequence  $\{P(x^{(k)})\}$ . Beyond that, we shall not prove any stronger convergence result on the iterates  $\{x^{(k)}\}$  itself in this paper. The reason is that the convergence analysis may vary individually by choice of sparse metric. We refer the readers to [41] and [44], in which subsequential convergence of  $\{x^{(k)}\}$  is established for DCA- $\ell_{1-2}$  and DCA-transformed- $\ell_1$  respectively.

**Algorithm 3.1** Iterative  $\ell_1$  minimizationInitialize:  $x^{(0)} = \mathbf{0}$ .

```

for  $k = 1, 2, \dots$  do
   $y^{(k)} = \text{sgn}(x^{(k)}) \circ (\mathbf{1} - \frac{P'(|x^{(k)}|)}{p'(0+)})$ 
   $x^{(k+1)} = \arg \min_x \|x\|_1 - \langle y^{(k)}, x \rangle \quad \text{s.t.} \quad Ax = A\bar{x}$ 
end for

```

**Extensions.** Two natural extensions of (1.3) are regularized model:

$$\min_{x \in \mathbb{R}^n} \frac{1}{2} \|Ax - b\|_2^2 + \lambda P(|Dx|), \quad (3.2)$$

and denoising model:

$$\min_{x \in \mathbb{R}^n} P(|Dx|) \quad \text{s.t.} \quad \|Ax - b\|_2 \leq \sigma, \quad (3.3)$$

where  $b$  is the measurement,  $D$  is a general matrix, and  $\lambda, \sigma > 0$  are parameters. They find applications in magnetic resonance imaging [5], total variation denoising [33] and so on. We can show that DC decomposition of  $P(|Dx|)$  is

$$P(|Dx|) = p'(0+) \|Dx\|_1 - (p'(0+) \|Dx\|_1 - P(|Dx|)). \quad (3.4)$$

The iterative  $\ell_1$  frameworks are detailed in Algorithms 3.2 and 3.3 respectively.**Algorithm 3.2** Iterative  $\ell_1$  regularizationInitialize:  $x^{(0)} = \mathbf{0}$ .

```

for  $k = 1, 2, \dots$  do
   $y^{(k)} = D^T \left( \text{sgn}(Dx^{(k)}) \circ (\mathbf{1} - \frac{P'(|Dx^{(k)}|)}{p'(0+)}) \right)$ 
   $x^{(k+1)} = \arg \min_x \frac{1}{2} \|Ax - b\|_2^2 + \lambda p'(0+) (\|Dx\|_1 - \langle y^{(k)}, x \rangle)$ 
end for

```

**Algorithm 3.3** Iterative  $\ell_1$  denoisingInitialize:  $x^{(0)} = \mathbf{0}$ .

```

for  $k = 1, 2, \dots$  do
   $y^{(k)} = D^T \left( \text{sgn}(Dx^{(k)}) \circ (\mathbf{1} - \frac{P'(|Dx^{(k)}|)}{p'(0+)}) \right)$ 
   $x^{(k+1)} = \|Dx\|_1 - \langle y^{(k)}, x \rangle \quad \text{s.t.} \quad \|Ax - b\|_2 \leq \sigma$ 
end for

```

## 4. Recovery Results

Although in general global minimum is not guaranteed in minimization, we can show that its performance is provably robust to the conditioning of measurement matrix  $A$ , by proving that it always tends to sharpen  $\ell_1$  solution.

Let us take another look at the assumptions on  $p$  which were crucial in the proof of Proposition 3.1. Since  $p$  is concave and increasing on  $[0, \infty)$ , we have

$$0 \leq \left( \frac{P'(|x|)}{p'(0+)} \right)_i \leq 1, \quad \forall x \in \mathbb{R}, 1 \leq i \leq n,$$

and thus  $\|R(x)\|_\infty \leq 1$ . Now we are ready to show the main results.

**Theorem 4.1 (Support-wise uniform recovery)** *Let  $T \subseteq \{1, \dots, n\}$  be an arbitrary but fixed index set. If basis pursuit uniquely recovers all  $\bar{x}$  supported on  $T$ , so does (3.1).*

*Proof.* By the assumption that basis pursuit uniquely recovers all  $\bar{x}$  supported on  $T$ , and by the well-known null space property [22] for  $\ell_1$  minimization, we must have

$$\|h_T\|_1 < \|h_{T^c}\|_1, \quad \forall h \in \text{Ker}(A) \setminus \{\mathbf{0}\},$$

and  $x^{(1)} = \bar{x}$  in (3.1). The 2<sup>nd</sup> step of DCA reads

$$x^{(2)} = \arg \min \|x\|_1 - \langle R(\bar{x}), x \rangle \quad \text{s.t.} \quad Ax = A\bar{x}.$$

Let  $x^{(2)} = \bar{x} + h^{(2)}$ , then

$$\begin{aligned} \|\bar{x}\|_1 - \langle R(\bar{x}), \bar{x} \rangle &\geq \|\bar{x} + h^{(2)}\|_1 - \langle R(\bar{x}), \bar{x} + h^{(2)} \rangle \\ \implies \|\bar{x}\|_1 - \langle R(\bar{x}), \bar{x} \rangle &\geq \|\bar{x}\|_1 + \langle \text{sgn}(\bar{x}), h_T^{(2)} \rangle + \|h_{T^c}^{(2)}\|_1 - \langle R(\bar{x}), \bar{x} + h^{(2)} \rangle \\ \iff -\langle \text{sgn}(\bar{x}) - R(\bar{x}), h_T^{(2)} \rangle &\geq \|h_{T^c}^{(2)}\|_1 \\ \iff -\langle \text{sgn}(\bar{x}) \circ \frac{P'(|\bar{x}|)}{p'(0+)}, h_T^{(2)} \rangle &\geq \|h_{T^c}^{(2)}\|_1 \end{aligned}$$

Since  $\|\frac{P'(|\bar{x}|)}{p'(0+)}\|_\infty \leq 1$ , we have

$$\|h_T^{(2)}\|_1 \geq \|h_{T^c}^{(2)}\|_1.$$

As a result,  $h^{(2)}$  must be  $\mathbf{0}$ .

If nonzero entries of  $\bar{x}$  have the same magnitude, a stronger result holds that (3.1) recovers any fixed signal whenever basis pursuit does.  $\square$

**Theorem 4.2 (Recovery of equal-height signals)** *Let  $\bar{x}$  be a signal with equal-height peaks supported on  $T$ , i.e.*

$$|x_i| = |x_j|, \quad \forall i, j \in T.$$

*If the basis pursuit uniquely recovers  $\bar{x}$ , so does (3.1).*

*Proof.* If basis pursuit uniquely recovers  $\bar{x}$ , then for all  $h \in \text{Ker}(A) \setminus \{\mathbf{0}\}$ ,

$$\|\bar{x}\|_1 < \|\bar{x} + h\|_1 = \|\bar{x} + h_T\|_1 + \|h_{T^c}\|_1.$$

This implies that for all  $h \in \text{Ker}(A) \setminus \{\mathbf{0}\}$  and

$$\|h\|_\infty \leq \min_{i \in T} |\bar{x}_i|, \quad \|\bar{x}\|_1 < \|\bar{x} + h_T\|_1 + \|h_{T^c}\|_1 = \|\bar{x}\|_1 + \langle \text{sgn}(\bar{x}), h_T \rangle + \|h_{T^c}\|_1.$$

So for all  $h \in \text{Ker}(A) \setminus \{\mathbf{0}\}$  and  $\|h\|_\infty \leq \min_{i \in T} |\bar{x}_i|$ , we have  $-\langle \text{sgn}(\bar{x}), h_T \rangle < \|h_{T^c}\|_1$ .

Therefore,

$$-\langle \text{sgn}(\bar{x}), h_T \rangle < \|h_{T^c}\|_1, \quad \forall h \in \text{Ker}(A) \setminus \{\mathbf{0}\}, \quad (4.1)$$

and also  $x^{(1)} = \bar{x}$ .

We let  $x^{(2)} = \bar{x} + h^{(2)}$ , and suppose that  $h^{(2)} \neq \mathbf{0}$ . Repeating the argument in Theorem 4.1 and by (4.1), we arrive at

$$-\langle \text{sgn}(\bar{x}) \circ \frac{P'(|\bar{x}|)}{p'(0+)}, h_T^{(2)} \rangle \geq \|h_{T^c}^{(2)}\|_1 > -\langle \text{sgn}(\bar{x}), h_T^{(2)} \rangle.$$

Since peaks of  $\bar{x}$  have equal height,  $(\frac{P'(|\bar{x}|)}{p'(0+)})_i \in [0, 1)$  is a constant for all  $i \in T$ . So  $-\langle \text{sgn}(\bar{x}) \circ \frac{P'(|\bar{x}|)}{p'(0+)}, h_T^{(2)} \rangle$  is non-negative and less than  $-\langle \text{sgn}(\bar{x}), h_T^{(2)} \rangle$ , which leads to a contradiction.  $\square$

**Remark 4.1.** Although the conditions proposed in Section 1 are not applicable to the metric  $\ell_{1-2}$  since it is not separable, it is not hard to generalize these conditions and iterative  $\ell_1$  algorithm to accommodate this case. The resulting algorithm is exactly DCA- $\ell_{1-2}$  in [41]. We can also readily extend the recovery theory to DCA- $\ell_{1-2}$ , with  $P(x) = \|x\|_1 - \|x\|_2$  and

$$R(x) = \begin{cases} \frac{x}{\|x\|_2} & \text{if } x \neq \mathbf{0}, \\ \mathbf{0} & \text{if } x = \mathbf{0}. \end{cases}$$

Theorem 4.1 provides a theoretical explanation for the experimental observations made in [29, 30, 41] that DCA- $\ell_{1-2}$  performs consistently better than  $\ell_1$  minimization.

## 5. Numerical Experiments

### 5.1. Exact recovery of sparse vectors

We reconstruct sparse vector  $\bar{x}$  using iterative  $\ell_1$  algorithm (Algorithm 3.2 with  $D = I_n$ ) for minimizing the regularized model (3.2) with smoothed  $\ell_q$  norm (IL $_1$ - $\ell_q$ ) and smoothed logarithm functional (IL $_1$ -log), and compare them with two state-of-the-art non-convex CS algorithms, namely IRLS- $\ell_q$  [27] and IRL $_1$  [7]. Note that IRLS- $\ell_q$  and IRL $_1$  attempt to minimize  $\ell_q$  and logarithm, respectively, and both involve a smoothing strategy in minimization. So it would be particularly interesting to compare IL $_1$ - $\ell_q$  with IRLS- $\ell_q$ , and IL $_1$ -log with IRL $_1$ .  $q = 0.5$  is chosen for IRLS- $\ell_q$  and IL $_1$ - $\ell_q$  in all experiments. We shall also include ADMM- $\ell_1$  [3] for solving  $\ell_1$  regularization (LASSO) in comparison.

Experiments are carried out as follows. We first sample a sensing matrix  $A \in \mathbb{R}^{m \times n}$ , and generate a test signal  $\bar{x} \in \mathbb{R}^n$  of sparsity  $s$  supported on a random index set with i.i.d. Gaussian entries. We then compute the measurement  $A\bar{x}$  and apply each solver to produce a reconstruction  $x^*$  of  $\bar{x}$ . The reconstruction is called a success if

$$\frac{\|x^* - \bar{x}\|_2}{\|\bar{x}\|_2} < 10^{-3}.$$

We run 100 independent realizations and record the corresponding success rates at different sparsity levels.

**Matrix for test.** We test on random Gaussian matrix whose columns satisfy

$$A_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, I_m/m), \quad i = 1, \dots, n$$

Gaussian matrices are RIP-like and have uncorrelated (incoherent) columns. For Gaussian matrix, we choose  $m = 64$  and  $n = 256$ .



We also use more ill-conditioned sensing matrix of significantly higher coherence. Specifically, a randomly oversampled partial DCT matrix  $A$  is defined as

$$A_i = \frac{1}{\sqrt{m}} \cos(2i\pi\xi/F), \quad i = 1, \dots, n$$

where  $\xi \in \mathbb{R}^m \sim \mathcal{U}([0, 1]^m)$  whose components are uniformly and independently sampled from  $[0, 1]$ .  $F \in \mathbb{N}$  is the refinement factor. Coherence  $\mu(A)$  goes up as  $F$  increases. In this setting, it is still possible to recover the sparse vector  $\bar{x}$  if its spikes are sufficiently separated. Specifically, we randomly select a  $T$  (support of  $\bar{x}$ ) so that

$$\min_{j, k \in T} |j - k| \geq L,$$

where  $L$  is called the minimum separation. It is necessary for  $L$  to be at least 1 Rayleigh length (RL) which is unity in the frequency domain [16, 21]. In our case, the value of 1 RL equals  $F$ . The testing matrix  $A \in \mathbb{R}^{100 \times 1500}$ , i.e.  $m = 100$ ,  $n = 1500$ . We test at three coherence levels with  $F = 5, 10, 15$ . Note that  $\mu(A) \approx 0.95$  for  $F = 5$ ,  $\mu(A) \approx 0.998$  for  $F = 10$ , and  $\mu(A) \approx 0.9996$  for  $F = 15$ . We also set  $L = 2F$  in experiments.

**Algorithm implementation.** For ADMM- $\ell_1$ , we let  $\lambda = 10^{-6}$ ,  $\beta = 1$ ,  $\rho = 10^{-5}$ ,  $\epsilon^{\text{abs}} = 10^{-7}$ ,  $\epsilon^{\text{rel}} = 10^{-5}$ , and the maximum number of iterations `maxiter` = 5000 [3, 41]. For IRLS- $\ell_q$ , `maxiter` = 1000, `tol` =  $10^{-8}$ . For reweighted  $\ell_1$ , the smoothing parameter  $\varepsilon$  is adaptively updated as introduced in [7], and the outer iteration criterion is stopped if the relative error between two consecutive iterates is less than  $10^{-2}$ . The weighted  $\ell_1$  minimization subproblems is solved by the YALL1 solver (available at <http://yall1.blogs.rice.edu/>). The tolerance for YALL1 was set to  $10^{-6}$ . All other settings of the algorithms are set to default ones.

For  $\text{IL}_1$ - $\ell_q$ , we let  $\lambda = 10^{-6}$ , and the smoothing parameter  $\varepsilon = \max\{\frac{|x^{(1)}|_{(d)}}{3}, 0.01\}$ , where  $x^{(1)}$  is the output from the first iteration, which is also the solution to LASSO.  $|x|_{(d)}$  denotes the  $d^{\text{th}}$  largest entry of  $|x|$ . We set  $d$  to  $\lfloor \frac{m}{4} \rfloor$ . For  $\text{IL}_1$ -log,  $\varepsilon = \max\{|x^{(1)}|_{(d)}, 0.01\}$ . The subproblems are solved by alternating direction method of multipliers (ADMM), which is detailed in [41]. The parameters for solving subproblems are the same as that for ADMM- $\ell_1$ .

**Interpretation of results.** The plot of success rates is shown in Figure 5.1. When  $A$  is Gaussian, we see that all non-convex CS solvers are comparable and much better than ADMM- $\ell_1$ , with IRLS- $\ell_q$  being the best. For oversampled DCT matrices, we see that the success rates of IRLS- $\ell_q$  and  $\text{IRL}_1$  drop as  $F$  increases, whereas the proposed  $\text{IL}_1$ - $\ell_q$  and  $\text{IL}_1$ -log are robust and consistently outperform ADMM- $\ell_1$ .

## 5.2. MRI reconstruction

We present an example of reconstructing the shepp-Logan phantom image of size  $256 \times 256$ , to further demonstrate effectiveness of  $\text{IL}_1$  algorithm. In this application, the sparsity of the gradient of the image/signal denoted by  $u$  is exploited, which leads to the following minimization problem:

$$\min_u P(|\nabla u|) \quad \text{s.t.} \quad \mathcal{F}Su = b,$$

where  $S$  denotes the sampling mask in the frequency domain, and  $\mathcal{F}$  is the Fourier transform and  $b$  the acquired data. With  $P(|\cdot|)$  being the  $\ell_1$  norm, the above formulation reduces to the

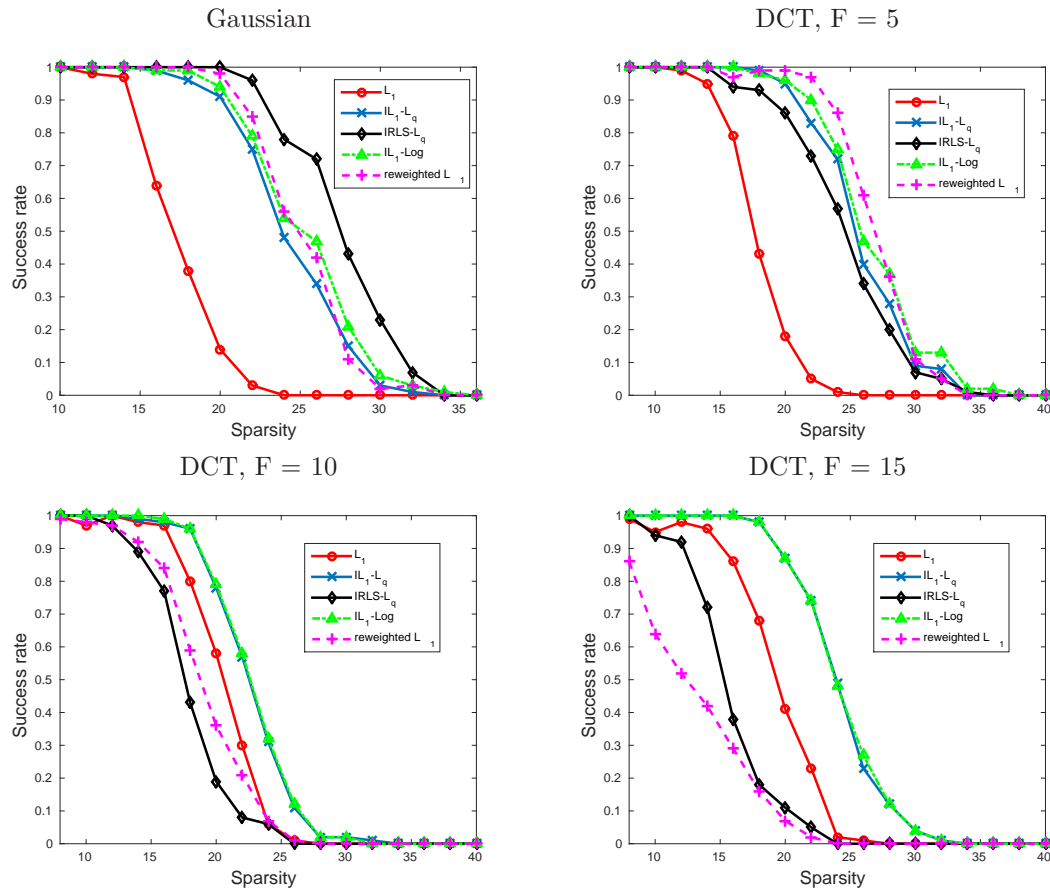


Fig. 5.1. Plots of success rates for comparing the iterative  $\ell_1$  with other CS algorithms under the increasing coherence of the sensing matrices.

celebrated total variation (TV) minimization:

$$\min_u \|\nabla u\|_1 \quad \text{s.t.} \quad \mathcal{S}\mathcal{F}u = b.$$

The above unconstrained problem together with its regularized problem

$$\min_u \frac{1}{2} \|\mathcal{S}\mathcal{F}u - b\|_2^2 + \lambda \|\nabla u\|_1, \quad (5.1)$$

can be solved efficiently by split Bregman method [24], known to be equivalent to ADMM [18]. For general sparse metric  $P$  (or  $p$ ), (3.4) gives the DC decomposition

$$P(|\nabla u|) = p'(0+) \|\nabla u\|_1 - (p'(0+) \|\nabla u\|_1 - P(|\nabla u|)),$$

and thus its  $\text{IL}_1$  algorithm for solving the regularized model takes the following form

$$\begin{cases} w^{(k)} = \nabla^T \left( \text{sgn}(\nabla u^{(k)}) \circ \left( \mathbf{1} - \frac{P'(|\nabla u^{(k)}|)}{p'(0+)} \right) \right) \\ u^{(k+1)} = \arg \min_u \frac{1}{2} \|\mathcal{S}\mathcal{F}u - b\|_2^2 + \lambda p'(0+) (\|\nabla u\|_1 - \langle w^{(k)}, u \rangle). \end{cases}$$

Likewise the subproblem for updating  $u^{(k+1)}$  above can also be solved by split Bregman, as the objective only differs by a linear term compared with (5.1).

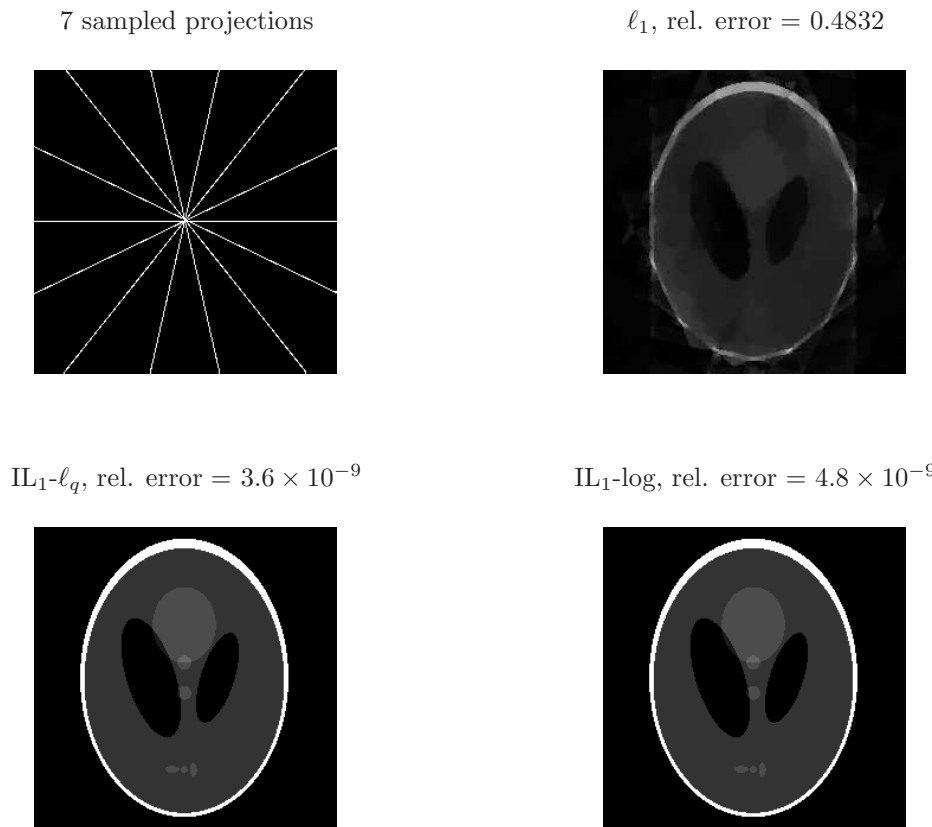


Fig. 5.2. Sampled lines and reconstructions for Shepp-Logan phantom image.

**Numerical results.** In the experiment, we again choose  $p$  to be the smoothed  $\ell_q$  ( $q = \frac{1}{2}$ ) and smoothed log respectively for the  $\text{IL}_1$  algorithm, and set smoothing parameter  $\varepsilon = 0.1$  and regularization parameter  $\lambda = 10^{-6}$  for both implementations. The reconstruction results are shown in Figure 5.2. We find that 7 sampled projections are sufficient for both of the two penalties to recover the phantom image perfectly, in comparison to  $\ell_1$  (TV) minimization which needs 10 projections for perfect image reconstruction by split Bregman. To the best of our knowledge, the other existing non-convex solvers for minimizing either  $\ell_q$  or log penalties did no better than 10 projections [7, 9, 10].

## 6. Conclusions

We developed an iterative  $\ell_1$  framework for a broad class of Lipschitz continuous non-convex sparsity promoting objectives, including those arising in statistics. The iterative  $\ell_1$  algorithm is shown via theory and computation to improve on the  $\ell_1$  minimization for CS problems independent of the coherence of the sensing matrices.

**Acknowledgments.** The authors would like to thank Yifei Lou (University of Texas at Dallas) and Jong-Shi Pang (University of Southern California) for helpful discussions. The authors

would also like to thank anonymous reviewers for their helpful comments. The work was partially supported by NSF grants DMS-1222507 and DMS-1522383.

## References

- [1] A. Beck and M. Teboulle, A Fast Iterative Shrinkage-Thresholding Algorithm for Linear Inverse Problems, *SIAM J. Imaging Sci.*, **2**:1, 2009, 183-202.
- [2] T. Blumensath and M. Davies, Iterative hard thresholding for compressed sensing, *Applied and Computational Harmonic Analysis*, **27**:3 (2009), 265-274.
- [3] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, Distributed Optimization and Statistical Learning via the Alternating Direction Method of Multipliers, *Foundations and Trends in Machine Learning*, 3:1 (2011), 1-122.
- [4] E. Candès and C. Fernandez-Granda, Towards a mathematical theory of super-resolution, *Communication on Pure and Applied Mathematics*, **67**:6 (2014), 906-956.
- [5] E. Candès, J. Romberg, and T. Tao, Robust uncertainty principles: Exact signal reconstruction from highly incomplete Fourier information, *IEEE Transactions on Information Theory*, **52**:2 (2006), 489-509.
- [6] E. Candès, J. Romberg, and T. Tao. Stable signal recovery from incomplete and inaccurate measurements, *Communications on Pure and Applied Mathematics*, **59**:8 (2006), 1207-1223.
- [7] E. Candès, M. Wakin, and S. Boyd, Enhancing Sparsity by Reweighted  $\ell_1$  Minimization, *J. Fourier Anal. Appl.*, **14**:5 (2008), 877-905.
- [8] E. Candès and T. Tao, Decoding by Linear Programming, *IEEE Trans. Info. Theory*, **51** (2005), 4203-4215.
- [9] R. Chartrand, Exact reconstruction of sparse signals via nonconvex minimization, *Signal Process. Lett.*, **14**:10 (2007), 707-710.
- [10] R. Chartrand, Fast algorithms for nonconvex compressive sensing: MRI reconstruction from very few data, *IEEE International Symposium on Biomedical Imaging: From Nano to Macro*, 2009.
- [11] R. Chartrand and V. Staneva, Restricted isometry properties and nonconvex compressive sensing, *Inverse Problems*, **24** (2008), 1-14.
- [12] R. Chartrand and W. Yin, Iteratively Reweighted Algorithms for Compressive Sensing, *IEEE international conference on acoustics, speech, and signal processing*, (2008), 3869-3872,
- [13] A. Cohen, W. Dahmen, and R. DeVore, Compressed Sensing and Best  $k$ -Term Approximation, *Journal of the American Mathematical Society*, **22**:1 (2009), 221-231.
- [14] I. Daubechies, M. Defrise, C. De Mol, An iterative thresholding algorithm for linear inverse problems with a sparsity constraint, *Comm. Pure Applied Math.*, **57**(2004), 1413-1457.
- [15] D. Donoho, X. Huo, Uncertainty principles and ideal atomic decomposition, *IEEE Trans. Inform. Theory*, **47** (2001), 2845-2862.
- [16] D. Donoho. Super-resolution via sparsity constraints, *SIAM Journal on Mathematical Analysis*, **23**:5 (1992), 1309-1331.
- [17] D. Donoho. Compressed sensing, *IEEE Transactions on Information Theory*, **52**:4, (2006), 1289-1306.
- [18] E. Esser, Applications of Lagrangian-Based Alternating Direction Methods and Connections to Split Bregman, CAM-report 09-31, UCLA, 2009.
- [19] E. Esser, Y. Lou, and J. Xin, A Method for Finding Structured Sparse Solutions to Non-negative Least Squares Problems with Applications, *SIAM J. Imaging Sciences*, **6**:4 (2013), 2010-2046.
- [20] J. Fan and R. Li, Variable selection via nonconcave penalized likelihood and its oracle properties, *Journal American Statistical Association*, **96**:456 (2001), 1348-1360.
- [21] A. Fannjiang and W. Liao, Coherence Pattern-Guided Compressive Sensing with Unresolved Grids, *SIAM J. Imaging Sciences*, **5**:1 (2012), 179-202.
- [22] S. Foucart and H. Rauhut, A Mathematical Introduction to Compressive Sensing, Springer, 2013.

- [23] G. Gasso, A. Rakotomamonjy, and S. Canu, Recovering Sparse Signals With a Certain Family of Nonconvex Penalties and DC Programming, *IEEE Transactions on Signal Processing*, **57**:12 (2009), 4686-4698.
- [24] T. Goldstein and S. Osher, The split Bregman method for L1 regularized problems, *SIAM Journal on Imaging Sciences*, **2**:1 (2009), 323-343.
- [25] E. Hale, W. Yin, and Y. Zhang, Fixed-point continuation for  $\ell_1$ -minimization: Methodology and convergence, *SIAM J. Optim.*, **19** (2008), 1107-1130.
- [26] X. Huang, L. Shi, and M. Yan, Nonconvex Sorted  $\ell_1$  Minimization for Sparse Approximation, *Journal of the Operations Research Society of China*, **3**:2 (2015), 207-229.
- [27] M. Lai, Y. Xu, and W. Yin, Improved Iteratively reweighted least squares for unconstrained smoothed  $\ell_q$  minimization, *SIAM J. Numer. Anal.*, **51**:2 (2013), 927-957.
- [28] J. Lv, and Y. Fan, A unified approach to model selection and sparse recovery using regularized least squares, *Annals of Statistics*, **37**:6A (2009), 3498-3528.
- [29] Y. Lou, P. Yin, Q. He, and J. Xin, Computing Sparse Representation in a Highly Coherent Dictionary Based on Difference of  $L_1$  and  $L_2$ , *J. Sci. Comput.*, **64** (2015), 178-196.
- [30] Y. Lou, P. Yin, and J. Xin, Point Source Super-resolution via Non-convex  $L_1$  Based Methods, *J. Sci. Computing*, **68**:3 (2016), 1082-1100.
- [31] D. Needell and J. Tropp, CoSaMP: Iterative signal recovery from incomplete and inaccurate samples, *Appl. Comput. Harmon. Anal.*, **26** (2009), 301-221.
- [32] H. Rauhut, Compressive Sensing and Structured Random Matrices, *Radon Series Comp. Appl. Math.*, **9** (2010), 1-92.
- [33] L. Rudin, S. Osher, and E. Fatemi, Nonlinear total variation based noise removal algorithms, *Physica D*, **60** (1992), 259-268.
- [34] P.D. Tao and L.T.H. An, Convex analysis approach to d.c. programming: Theory, algorithms and applications, *Acta Mathematica Vietnamica*, **22**:1 (1997), 289-355.
- [35] P.D. Tao and L.T.H. An, A DC optimization algorithm for solving the trust-region subproblem, *SIAM Journal on Optimization*, **8**:2 (1998), 476-505.
- [36] R. Tibshirani, Regression shrinkage and selection via the lasso, *J. Royal. Statist. Soc.*, **58**:1 (1996), 267-288.
- [37] J. Tropp and A. Gilbert, Signal recovery from partial information via orthogonal matching pursuit, *IEEE Trans. Inform. Theory*, **53**:12 (2007), 4655-4666.
- [38] Y. Wang and W. Yin, Sparse signal reconstruction via iterative support detection, *SIAM Journal on Imaging Sciences*, **3**:3 (2010), 462-491.
- [39] Z. Xu, X. Chang, F. Xu, H. Zhang,  $L_{1/2}$  regularization: an iterative thresholding method, *IEEE Transactions on Neural Networks and Learning Systems*, **23** (2012), 1013-1027.
- [40] J. Yang and Y. Zhang, Alternating direction algorithms for  $l_1$  problems in compressive sensing, *SIAM J. Sci. Comput.*, **33**:1 (2011), 250-278.
- [41] P. Yin, Y. Lou, Q. He, and J. Xin, Minimization of  $\ell_{1-2}$  for Compressed Sensing, *SIAM J. Sci. Comput.*, **37** (2015), A536-A563.
- [42] P. Yin, E. Esser, J. Xin, Ratio and Difference of  $L_1$  and  $L_2$  Norms and Sparse Representation with Coherent Dictionaries, *Commun. Information and Systems*, **14**:2 (2014), 87-109.
- [43] W. Yin, S. Osher, D. Goldfarb, and J. Darbon, Bregman iterative algorithms for  $\ell_1$  minimization with applications to compressed sensing, *SIAM J. Imaging Sci.*, **1** (2008), 143-168.
- [44] S. Zhang and J. Xin, Minimization of Transformed  $L_1$  Penalty: Theory, Difference of Convex Function Algorithm, and Robust Application in Compressed Sensing, preprint, arXiv:1411.5735.
- [45] T. Zhang, Multi-stage convex relaxation for learning with sparse regularization, *NIPS proceedings*, 2008.
- [46] S. Zhou, N. Xiu, Y. Wang, L. Kong, and H. Qi, A null-space-based weighted  $\ell_1$  minimization approach to compressed sensing, *Information and Inference*, 2016: iaw002v1-iaw002.