

Computing Sparse Representation in a Highly Coherent Dictionary Based on Difference of L_1 and L_2

Yifei Lou · Penghang Yin · Qi He · Jack Xin

Received: 7 January 2014 / Revised: 19 July 2014 / Accepted: 24 September 2014 /
Published online: 16 October 2014
© Springer Science+Business Media New York 2014

Abstract We study analytical and numerical properties of the $L_1 - L_2$ minimization problem for sparse representation of a signal over a highly coherent dictionary. Though the $L_1 - L_2$ metric is non-convex, it is Lipschitz continuous. The difference of convex algorithm (DCA) is readily applicable for computing the sparse representation coefficients. The L_1 minimization appears as an initialization step of DCA. We further integrate DCA with a non-standard simulated annealing methodology to approximate globally sparse solutions. Non-Gaussian random perturbations are more effective than standard Gaussian perturbations for improving sparsity of solutions. In numerical experiments, we conduct an extensive comparison among sparse penalties such as L_0 , L_1 , L_p for $p \in (0, 1)$ based on data from three specific applications (over-sampled discrete cosine basis, differential absorption optical spectroscopy, and image denoising) where highly coherent dictionaries arise. We find numerically that the $L_1 - L_2$ minimization persistently produces better results than L_1 minimization, especially when the sensing matrix is ill-conditioned. In addition, the DCA method outperforms many existing algorithms for other nonconvex metrics.

Keywords Highly coherent dictionary · Sparse representation · $L_1 - L_2$ minimization · Difference of convex programming · Simulated annealing · Comparison with L_p $p \in (0, 1)$

1 Introduction

Sparse representation in an overcomplete basis appears frequently in signal processing and imaging applications such as oversampled discrete Fourier transform, Gabor frames, curvelet

The work was partially supported by NSF grants DMS- 0928427 and DMS-1222507.

Y. Lou · P. Yin · Q. He · J. Xin
Department of Mathematics, UC Irvine, Irvine, CA 92697, USA

Present Address:

Y. Lou (✉)
Department of Mathematical Sciences, University of Texas at Dallas, Richardson, TX 75080, USA
e-mail: louyifei@gmail.com

frames, concatenation of different orthonormal bases [3]. It is known to be related to sparse coding in visual systems [19]. The advantage is robustness and reliability.

Mathematically it amounts to finding the sparsest solution to an under-determined linear system

$$b = Ax + n, \tag{1}$$

where b is the observed data, A is a $M \times N$ ($M < N$) matrix and n is noise. A fundamental issue in compressed sensing (CS) is how to enforce sparsity when solving the linear system (1). A natural strategy is to minimize L_0 norm, $\|x\|_0$, which is the number of nonzero elements. Unfortunately, the L_0 minimization is NP hard [18]. There are two methods that approach L_0 directly. One is the greedy approach, the so called orthogonal matching pursuit (OMP) [23]. Its basic idea is to select one best column from the matrix A at a time, followed by an orthogonal projection to avoid selecting the same vector multiple times. The other is the penalty decomposition method [17], which solves

$$\min_{x \in \mathbb{R}^N} \frac{1}{2} \|Ax - b\|_2^2 + \lambda \|x\|_0, \tag{2}$$

by a series of minimization problems with an increasing sequence $\{\rho^k\}$:

$$\begin{aligned} (x^{k+1}, y^{k+1}) &= \arg \min \frac{1}{2} \|Ax - b\|^2 + \frac{\rho^k}{2} \|x - y\|^2 + \lambda \|y\|_0 \\ \rho^{k+1} &= \tau \rho^k \quad (\tau > 1). \end{aligned} \tag{3}$$

In general, the two approaches only provide sub-optimal sparse solutions to the original problem.

The convex relaxation of L_1 in lieu of L_0 attracts considerable attention in CS. There are numerous algorithms, such as LASSO [22], Bregman iterations [27], and alternating direction method of multipliers (ADMM) [2], devoted to solving L_1 regularized problems efficiently and accurately. The theoretical aspect of L_1 relaxation is studied in [5,8] and elsewhere. A deterministic result in [8] says that exact L_1 recovery is possible if

$$\|x\|_0 < \frac{1 + 1/\mu}{2}, \tag{4}$$

where μ is mutual coherence of the matrix $A = [\mathbf{a}_1, \dots, \mathbf{a}_N]$, defined as

$$\mu = \max_{i \neq j} \frac{|\mathbf{a}_i^T \mathbf{a}_j|}{\|\mathbf{a}_i\| \|\mathbf{a}_j\|}.$$

The inequality (4) suggests that L_1 may not perform well for highly coherent sensing matrices in that if $\mu \sim 1$, then $\|x\|_0$ can be at most 1. Though the theoretical estimate is far from sharp, we shall show numerical examples of such phenomenon later.

In this paper we study a non-convex but Lipschitz continuous metric $L_1 - L_2$ for sparse signal recovery in the highly coherent regime of CS, and compare with the concave and Hölder continuous sparsity metric L_p ($p \in (0, 1)$).

Recently, nonconvex measures, such as L_p for $p \in (0, 1)$ in [7], L_1/L_2 and $L_1 - L_2$ in [10,25], have been proposed as alternatives to L_1 . As illustrated in Fig. 1 in \mathbb{R}^2 , the level curves of L_p and $L_1 - L_2$ are closer to L_0 than those of L_1 . Geometrically, minimizing a sparse measure subject to linear constraints is equivalent to finding an interception of an affine subspace (corresponding to solutions of the linear constraints) with a level set of that measure so that the intersection is closest to a coordinate axis/plane. For L_p and $L_1 - L_2$, due to their curved level set, the interception is more likely to occur at the coordinate axis/plane, giving

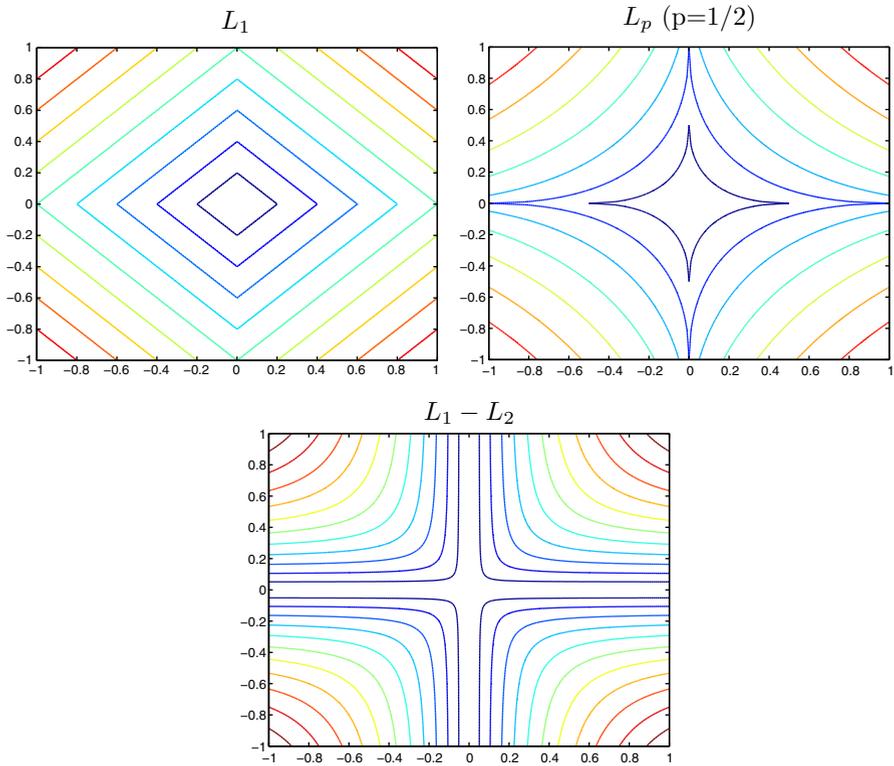


Fig. 1 Level lines of three sparsity metrics. Compared with L_1 , the level lines of L_p and $L_1 - L_2$ are closer to the axes when minimized (closer to those of L_0)

a sparse solution. For L_1 , it is possible that the affine subspace coincides with a segment of a level set (a resonance phenomenon), i.e., any point on that segment is a solution of L_1 minimization. If such resonance occurs, L_1 minimization fails to find a sparse solution. There may be other degenerate scenarios in three and above dimensions when the sensing matrix is highly coherent.

Though L_p and $L_1 - L_2$ measures are theoretically better than L_1 to promote sparsity, the non-convexity poses a challenge to computation. For L_p minimization, iterative reweighted least-square [7] is considered for the constrained problem, while an unconstrained formulation is discussed in [16]. The minimizing sequence may get stuck at a stationary point. The L_p metric has an a-priori unknown parameter p and is non-Lipschitz. The $L_1 - L_2$ is however Lipschitz continuous and free of parameter. It can be minimized by the difference of convex algorithm (DCA) [21] where linearization convexifies the objective without additional smoothing or regularization. The DCA minimizing sequence converges to a stationary point theoretically, which empirically often (yet not always) turns out to be a global minimizer. To avoid trapping in a local basin, we further incorporate a variant of the so called simulated annealing (SA) technique in global optimization. Here we found that non-Gaussian random perturbations are better at improving sparsity of solutions than standard Gaussian perturbations. A hybrid method integrating SA for L_p minimization is discussed in [24].

The rest of the paper is organized as follows. In Sect. 2.1, we show a toy example when L_1 minimization fails to find the sparsest solution, while L_p and $L_1 - L_2$ minimization succeed. In Sect. 2.2, we then present some nice properties of $L_1 - L_2$ metric as a sparsity measure. In Sect. 3, we discuss the algorithms for computing sparse representation based on $L_1 - L_2$ minimization, where both constrained and unconstrained formulations are given. In order to find a global solution, we further integrate DCA and the simulated annealing method. As numerical experiments, we investigate three specific applications (over-sampled discrete cosine transform, optical spectroscopy, and image denoising) where highly coherent matrices are encountered. We demonstrate that $L_1 - L_2$ minimization with DCA solver is robust in finding sparse solutions, and outperforms some state-of-the-art algorithms. Finally, discussion and conclusion are given in Sect. 5.

2 Sparsity Measures

2.1 A Toy Example

We study a toy example where L_1 fails to find the sparsest solution, while both $L_1 - L_2$ and L_p can. Consider a linear system

$$\begin{bmatrix} 2 & 1 & 0 \\ 1 & 1 & 1 \\ 0 & 1 & 2 \end{bmatrix} x = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \tag{5}$$

The sparsest solution is $x_0 = [0, 1, 0]^T$. We find any vector of form $[a, 1 - 2a, a]^T$ for $a \in [0, 0.5]$ is a solution of L_1 minimization subject to (5). In other words, L_1 fails to pick up the sparsest solution.

L_p minimization is equivalent to minimizing $2a^p + (1 - 2a)^p$ for $a \in [0, 0.5]$. As L_p is concave, the minimum is attained at the boundary. Consequently, we only need to evaluate $a = 0$ and $a = 0.5$ to see which one is smaller. By simple calculations, L_p attains its minimum at $a = 0$ for $p < 1$. In other words, L_p minimization yields the sparsest solution.

Let us look at $L_1 - L_2$. A non-zero vector is 1-sparse (only one non-zero element) if and only if its $L_1 - L_2$ is 0. Since x_0 is 1-sparse, then only 1-sparse vectors could be the solution of minimizing $L_1 - L_2$ subject to (5). But the other 1-sparse vectors do not satisfy the linear equation, and as a result we can get x_0 exactly if minimizing $L_1 - L_2$.

2.2 Theoretical Properties of $L_1 - L_2$

To make this paper self-contained, we list some nice properties of $L_1 - L_2$. Please refer to [26] for the proof. Recall that the well-known restricted isometry property (RIP) [5] in compressive sensing is that for all subsets $T \subset \{1, \dots, N\}$ with its cardinality $|T| \leq S$,

$$(1 - \delta_S)\|x\|_2^2 \leq \|A_T x\|_2^2 \leq (1 + \delta_S)\|x\|_2^2 \quad \forall x \in \mathbb{R}^N,$$

where A_T is a submatrix of A with column indices T , and δ_S is a parameter depending on S . The RIP condition for $L_1 - L_2$ exact recovery is given in the following.

Theorem 2.1 *Let x_0 be any vector with sparsity of s and $b = A x_0$. Suppose that the following condition holds*

$$\delta_{3s} + a(s)\delta_{4s} < a(s) - 1, \tag{6}$$

where

$$a(s) = \left(\frac{\sqrt{3s} - 1}{\sqrt{s} + 1} \right)^2,$$

then x_0 is the unique solution to a constrained minimization problem:

$$\min_{x \in \mathbb{R}^N} \|x\|_1 - \|x\|_2 \quad \text{subject to} \quad Ax = b. \tag{7}$$

In fact, Theorem 2.1 does not characterize $L_1 - L_2$ completely, as in practice its assumption can be further relaxed. Due to concavity of the metric, we can prove that even local minimizers of (7) satisfy certain sparsity, no matter whether A satisfies RIP or not.

Theorem 2.2 *Let x^* be a local minimizer of the constrained problem (7) then $A|_{\Lambda^*}$ is of full column rank, i.e. the columns of $A|_{\Lambda^*}$ are linearly independent.*

The same result can be obtained for the unconstrained problem:

$$\min_{x \in \mathbb{R}^N} \frac{1}{2} \|Ax - b\|_2^2 + \lambda(\|x\|_1 - \|x\|_2) \tag{8}$$

Corollary 2.1 *Let x^* be a local minimizer of (8) then the columns of $A|_{\Lambda^*}$ are linearly independent.*

By Theorem 2.2 and Corollary 2.1, we readily conclude the following facts:

- a. Suppose x^* is a local minimizer of (7) or (8) and $A \in \mathbb{R}^{M \times N}$ is of full row rank, i.e. $\text{rank}(A) = M$, then the sparsity of x^* is at most M .
- b. If x^* is a local minimizer of (7), then there is no such $x \in \mathbb{R}^N$ satisfying $Ax = b$ and $\text{support}(x) \subseteq \Lambda^*$, i.e. it is impossible to find a feasible solution whose support is contained in Λ^* .
- c. The number of the local minimizers of (7) or (8) is finite.

3 Algorithm

To compute sparse representation based on $L_1 - L_2$, both constrained and unconstrained minimization problems are discussed in Sects. 3.1 and 3.2 respectively. In Sect. 3.3, we further employ a simulated annealing technique to search for global solutions of these two nonconvex problems.

3.1 Unconstrained Minimization

We start with the unconstrained minimization problem (8). We adopt a DCA, which is to decompose $F(x) = G(x) - H(x)$, where

$$\begin{cases} G(x) = \frac{1}{2} \|Ax - b\|_2^2 + \lambda \|x\|_1 \\ H(x) = \lambda \|x\|_2 \end{cases} \tag{9}$$

By linearizing H , we can design an iterative scheme that starts with $x^1 \neq \mathbf{0}$,

$$x^{n+1} = \arg \min_{x \in \mathbb{R}^N} \frac{1}{2} \|Ax - b\|_2^2 + \lambda \|x\|_1 - \left\langle x - x^n, \lambda \frac{x^n}{\|x^n\|_2} \right\rangle \tag{10}$$

To advance to a new solution, it requires solving a L_1 regularized subproblem of the form

$$\min_{x \in \mathbb{R}^N} \frac{1}{2} x^T (A^T A)x + z^T x + \lambda \|x\|_1, \tag{11}$$

where $z = A^T b + \lambda \frac{x^n}{\|x^n\|_2}$. We consider the augmented Lagrangian

$$L_\delta(x, y, u) = \frac{1}{2} x^T (A^T A)x + z^T x + \lambda \|y\|_1 + u^T (x - y) + \frac{\delta}{2} \|x - y\|_2^2.$$

Alternating direction method of multipliers iterates between minimizing L_δ with respect to x , minimizing with respect to y and updating u . The pseudo-code of solving the unconstrained $L_1 - L_2$ minimization is described in Algorithm 1.

Algorithm 1 A DCA method for unconstrained $L_1 - L_2$ minimization

```

Define  $\epsilon_{outer} > 0, \epsilon_{inner} > 0$  and initialize  $x^0 = \mathbf{0}, x^1 \neq \mathbf{0}, n = 1$ 
while  $\|x^n - x^{n-1}\| > \epsilon_{outer}$  do
  Let  $z = \frac{x^n}{\|x^n\|_2}$  and  $x_0 = \mathbf{0}, x_1 = x^n, i = 1, y_i = x_i, u_i = \mathbf{0}$ 
  while  $\|x_i - x_{i-1}\| > \epsilon_{inner}$  do
     $x_{i+1} = (A^T A + \delta I)^{-1}(\delta y_i - z - u_i)$ 
     $y_{i+1} = shrink(x_{i+1} + u_i/\delta, \lambda/\delta)$ 
     $u_{i+1} = u_i + \delta(x_{i+1} - y_{i+1})$ 
     $i = i + 1$ 
  end while
   $n = n + 1$ 
   $x^n = x_i$ 
end while
    
```

In our experiments we always set the initial value x^1 as the solution of unconstrained L_1 problem, that is to solve (11) with $z = \mathbf{0}$. So basically we are minimizing $L_1 - L_2$ on top of L_1 . In practice the algorithm takes only a few steps to convergence. Theoretically, we can prove that the sequence $\{x^n\}$ is bounded and $\|x^{n+1} - x^n\|_2 \rightarrow 0$, thus limit points of $\{x^n\}$ are stationary points of (8) satisfying the first-order optimality condition.

Theorem 3.1 Let $F(x) = \frac{1}{2} \|Ax - b\|_2^2 + \lambda(\|x\|_1 - \|x\|_2)$ with the D.C. decomposition in (9) and $\{x^n\}$ be the sequence of iterates generated by Algorithm 1, then

- a. $F(x) \rightarrow \infty$ as $\|x\|_2 \rightarrow \infty$. so the level set $\{x \in \mathbb{R}^N : F(x) \leq F(x^0)\}$ is bounded.
- b. $\|x^{n+1} - x^n\|_2 \rightarrow 0$.
- c. Any limit point $x^* \neq \mathbf{0}$ of $\{x^n\}$ satisfies the first-order optimality condition

$$A^T (A x^* - b) + \lambda \left(w^* - \frac{x^*}{\|x^*\|_2} \right) = \mathbf{0}, \quad \text{for some } w^* \in \partial \|x^*\|_1, \tag{12}$$

which means x^* is a stationary point of (8).

Starting from (12), we obtain the following results.

Theorem 3.2 $\forall k \geq 1$, we can choose a regularization parameter $\lambda_k > 0$ for (8) so that $\|x^*\|_0 \leq k$.

Theorem 3.2 suggests that we can control the sparsity of the limit point of the DCA algorithm by choosing a proper λ . Please refer to [26] for the proof of both theorems.

3.2 Constrained Minimization

For the constrained problem (7), we apply a similar trick as the unconstrained problem by considering the following iterative scheme

$$\begin{aligned}
 x^{n+1} &= \arg \min \|x\|_1 - \left\langle \frac{x^n}{\|x^n\|_2}, x \right\rangle \\
 \text{s.t. } Ax &= b
 \end{aligned}
 \tag{13}$$

Each subproblem (13) amounts to solving a constrained L_1 minimization,

$$\min |x| - z^T x \quad \text{s.t. } Ax = b,
 \tag{14}$$

for $z = \frac{x^n}{\|x^n\|_2}$. To solve (14), we introduce two Lagrange multipliers u, v and define an augmented Lagrangian

$$L_\delta(x, y, u, v) = \|y\|_1 - z^T x + u^T(x - y) + v^T(Ax - b) + \frac{\delta}{2}\|x - y\|^2 + \frac{\delta}{2}\|Ax - b\|^2,$$

for $\delta > 0$. ADMM finds a saddle point

$$L_\delta(x^*, y^*, u, v) \leq L_\delta(x^*, y^*, u^*, v^*) \leq L_\delta(x, y, u^*, v^*) \quad \forall x, y, u, v$$

by alternately minimizing L_δ with respect to x , minimizing with respect to y and updating the dual variables u and v . The saddle point (x^*, y^*) will be a solution to (14) and we can take x^* be the solution to (13), i.e., $x^{n+1} = x^*$. The overall algorithm for solving the constrained $L_1 - L_2$ is described in Algorithm 2.

Algorithm 2 A DCA method for constrained $L_1 - L_2$ minimization

```

Define  $\epsilon_{outer} > 0, \epsilon_{inner} > 0$  and initialize  $x^0 = \mathbf{0}, x^1 \neq \mathbf{0}, n = 1$ 
while  $\|x^n - x^{n-1}\| > \epsilon_{outer}$  do
  Let  $z = \frac{x^n}{\|x^n\|_2}$  and  $x_0 = \mathbf{0}, x_1 = x^n, i = 1, y_i = x_i, u_i = v_i = \mathbf{0}$ 
  while  $\|x_i - x_{i-1}\| > \epsilon_{inner}$  do
     $x_{i+1} = (A^T A + I)^{-1}(A^T b + y_i + (z - u_i - A^T v_i)/\delta)$ 
     $y_{i+1} = shrink(x_{i+1} + u_i/\delta, 1/\delta)$ 
     $u_{i+1} = u_i + \delta(x_{i+1} - y_{i+1})$ 
     $v_{i+1} = v_i + \delta(Ax_{i+1} - b)$ 
     $i = i + 1$ 
  end while
   $n = n + 1$ 
   $x^n = x_i$ 
end while

```

3.3 Simulated Annealing

The DCA method does not guarantee a global minimum in general. We further employ a technique, called SA, to traverse a stationary point to a global solution. SA has drawn much attentions dealing with global optimization. There are many generic SA algorithms, see Kirkpatrick [15], Geman and Geman [13], Gidas [14], and the reference therein. In addition, Carnevali et al. [6] apply this technique to many applications in image processing.

Here is a brief description of simulated annealing. The term “annealing” is analogous to the cooling of a liquid or solid in a physical system. Consider the problem of minimizing the cost

function $f(x)$. Simulated annealing algorithm begins with an initial solution and iteratively generates new ones, each of which is randomly selected among the “neighborhood” of the previous state. If the new solution is better than the previous one, it is accepted; otherwise, it is accepted with certain probability. The probability of accepting a new state is given by $\exp(-\frac{f_{new}-f_{curr}}{T}) > R$, where R is a random number between 0 and 1, and T is a temperature parameter. The algorithm usually starts with a high temperature, and then gradually goes down to 0. The cooling must be slow enough so that the system does not get stuck into saddle points or local minima of $f(x)$.

There are two important aspects in implementing simulated annealing. One is how to lower the temperature T . Kirkpatrick et al. [15] suggest T decays geometrically in the number of cooling phases. Geman and Geman [13] prove that if T decreases at the rate of $\frac{1}{\log k}$, where k is the number of iterations, then the probability distribution for the algorithm converges to the uniform distribution over all the global minimum points. In our algorithm, we follow Geman and Geman’s suggestion by decreasing T at the rate of $\frac{1}{\log k}$. Another aspect is how to advance to a new state based on the current one. One of the most common methods is to add Gaussian random noise. However, due to the presence of a large number of saddle points and local minima, this perturbation method yields slow convergence of the SA algorithm. To overcome this difficulty, we propose two perturbative strategies. Together with Gaussian perturbation, we list the three SA methods as follows,

- SA1. Define the support set of a vector x by $support(x) := \{i | x_i \neq 0\}$. We randomly choose the new state x_{new} such that $support(x_{new}) \subset support(x_{curr})$ and $|x_{new}|_0 < |x_{curr}|_0$.
- SA2. We randomly choose the new state x_{new} such that $|x_{new}|_0 < \gamma|x_{curr}|_0$ with some constant $\gamma < 1$.
- SA3. We choose the new state x_{new} by Gaussian perturbations, i.e. $x_{new} = x_{curr} +$ Gaussian noise.

The idea of SA1 and SA2 is to help maintain the monotone (non-increasing) sparsity property of the iterates. The pseudo-code of the SA method in combination of DCA is given in Algorithm 3.

Algorithm 3 L1-L2 DCA with simulated annealing

```

Define  $x_{curr}, x_{new}, \epsilon, T, \gamma < 1, \maxIter$  and  $AcceptMax, Accept = 0$ .
while  $T > \epsilon$  or  $k \leq \maxIter$  do
  SA1:  $x_{new} = \text{randsample}(x_{curr})$  such that  $support(x_{new}) \subset support(x_{curr})$  and  $|x_{new}|_0 < |x_{curr}|_0$ .
  SA2:  $x_{new} = \text{randsample}(x_{curr})$  such that  $|x_{new}|_0 < \gamma|x_{curr}|_0$ , with  $\gamma < 1$ .
  SA3:  $x_{new} = x_{curr} +$  Gaussian noise.
  Update  $x_{new}$  by the DCA solution of  $L_1 - L_2$  using  $x_{new}$  as initial guess
  if  $f(x_{new}) \leq f(x_{curr})$  then
     $x_{curr} = x_{new}$ 
  else
    if  $\exp(-\frac{f(x_{new})-f(x_{curr})}{T}) > \text{rand}(1)$  then
       $x_{curr} = x_{new}$ 
    end if
  end if
   $k = k + 1$ 
   $Accept = Accept + 1$ 
  if  $Accept \geq AcceptMax$  then
     $T = \frac{1}{\log k}$ 
     $Accept = 0$ 
  end if
end while

```

4 Applications

In this section we examine three specific applications where the sensing matrix is highly coherent. They are compressive sensing based on oversampled DCT matrices, wavelength misalignment in differential optical absorption spectroscopy (DOAS) analysis, and image denoising via sparse representation in an overcomplete dictionary. For each problem, we compare the proposed method for minimizing $L_1 - L_2$ with some state-of-the-art algorithms for L_0 , L_1 , and L_p . Experiments show promising results of using $L_1 - L_2$ as a sparse measure and solving sparse coefficients by DCA and SA.

4.1 Over-Sampled DCT

We consider an over-sampled DCT matrix $A = [\mathbf{a}_1, \dots, \mathbf{a}_N] \in \mathbb{R}^{M \times N}$ with

$$\mathbf{a}_j = \frac{1}{\sqrt{N}} \cos\left(\frac{2\pi \mathbf{w} j}{F}\right), j = 1, \dots, N \quad (15)$$

where \mathbf{w} is a random vector of length M . This matrix is derived from the problem of spectral estimation [11] in signal processing, if we replace the cosine function in (15) by exponential. The matrix is highly coherent. For a $100 \times 1,000$ matrix with $F = 10$, the coherence is .9981, while the coherence of a same size matrix with $F=20$ is .9999.

The sparse recovery under such matrices is possible only if the non-zero elements of solution x are sufficiently separated. This phenomenon is characterized as minimum separation in [4], and this minimum length is referred as the Rayleigh length (RL). The value of RL in (15) is equal to F . It is closely related to the coherence in the sense that larger F corresponds to larger coherence of a matrix. We find empirically that at least $2RL$ is necessary to ensure optimal sparse recovery. Intuitively, we need sparse spikes to be further apart for more coherent matrices.

Under the assumption of sparse signal with $2RL$ separated spikes, we compare algorithms in terms of success rate. Denote x_r as a reconstructed solution by a certain algorithm. We consider the algorithm successful, if the relative error of x_r to the ground truth x_g is less than .001, i.e., $\frac{\|x_r - x_g\|}{\|x_g\|} < .001$. The success rate is based on 100 random realizations.

We first investigate the unconstrained algorithms for sparse penalties $L_1 - L_2$, L_1 , L_p ($p=1/2$) and a direct L_0 solver, penalty decomposition (3). The unconstrained L_p minimization is solved by Lai et al. in [16]. The sensing matrix is of size $100 \times 1,000$. The success rate of each measure is plotted in Fig. 2. For smaller $F = 5$, each measure performs relatively well. When $F = 20$, which corresponds to highly coherent matrices, the proposed $L_1 - L_2$ outperforms L_1 and L_p for $p = 1/2$. For both cases, $L_1 - L_2$ is consistently better than L_1 .

Figure 3 illustrates that the success rate of $L_1 - L_2$ increases with the help of SA. The matrix size is $100 \times 1,500$ and $F = 20$. We also compare three different random generations of the new state, referred to as SA1–SA3 in Algorithm 3. All of these SA methods can improve the accuracy of the original DCA algorithm for $L_1 - L_2$ minimization. Both of SA1 and SA2 have better performance than the regular Gaussian perturbation method SA3. Apparently, SA2 has the best performance out of the three, especially when the number of non-zero elements is large.

For the constrained versions of each measure, we observe the similar behavior compared to the unconstrained ones. The plots are presented in Fig. 4 for matrices of size $100 \times 1,000$. The iterative reweighted least square [7] is applied to solve the constrained L_p . $L_1 - L_2$ is the best for both incoherent and coherent matrices.

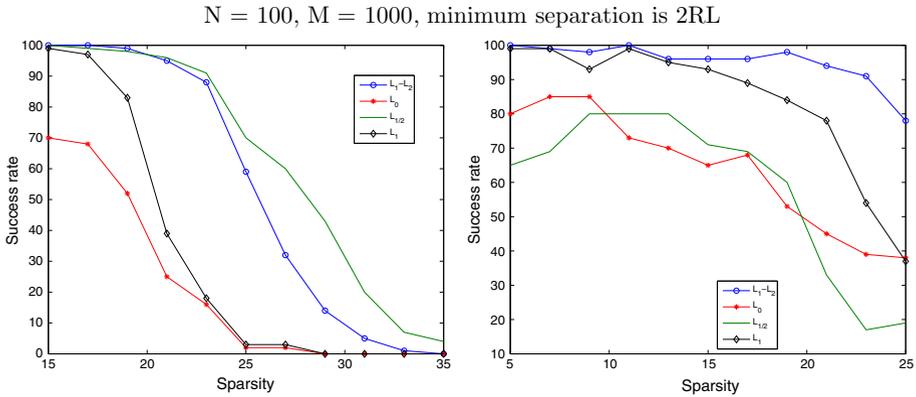


Fig. 2 Plots of success rates as a function of sparsity for $F = 5$ (left) and $F = 20$ (right). Each metric is solved in an unconstrained optimization framework

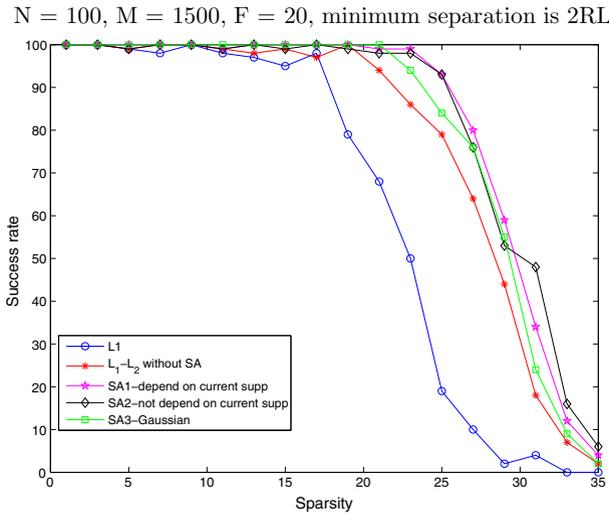


Fig. 3 Plots of success rates of three SA algorithms: SA1-the support of the new state is within the current one, SA2-the new state does not depend on the current one, and SA3-the new state is obtained by adding Gaussian noise to the current one. The results of L_1 and $L_1 - L_2$ without SA are plotted as reference

We further look at the success rates of $L_1 - L_2$ with different combinations of sparsity and RL. The rates are recorded in Table 1, which shows it is possible to deal with large RL, but with the sacrifice of high sparsity, or small number of non-zero elements.

The computation time for each method is listed in Table 2. The reported time is the average of 100 realizations. Since all the nonconvex optimization methods $L_0, L_p, L_1 - L_2$ use the L_1 solution as initial guess, we add L_1 run time on top of each of them. Our $L_1 - L_2$ method is slower than L_p , but with better accuracy.

4.2 DOAS

Differential optical absorption spectroscopy analysis [20] is a technique that uses Beer’s law to estimate chemical contents and concentrations of a mixture of gases by decompos-

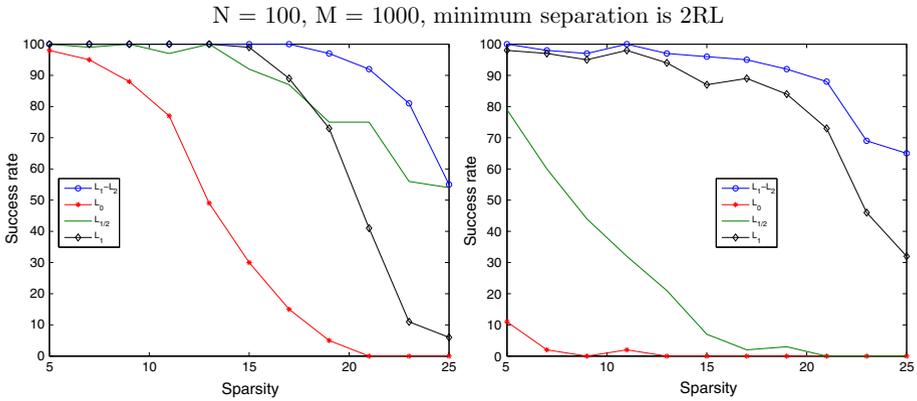


Fig. 4 Plots of success rates as a function of sparsity for $F = 5$ (left) and $F = 20$ (right). Each metric is solved in a constrained optimization framework

Table 1 The success rates (%) of $L_1 - L_2$ for different combinations of sparsity and minimum separation

Sparsity	6	9	12	15	18	21	24
1RL	100	100	100	97	85	37	3
2RL	100	100	100	100	88	44	11
3RL	100	100	100	99	90	41	8
4RL	100	100	100	98	85	46	8

Table 2 Average computation time for each method at different settings

	L_0	$L_{1/2}$	L_1	$L_1 - L_2$
Unconstrained $F = 5$	12.26	1.94	1.53	3.56
Unconstrained $F = 20$	12.23	2.01	1.54	3.55
Constrained $F = 5$	4.37	1.97	3.68	9.94
Constrained $F = 20$	6.35	7.44	3.67	10.91

ing a measured characteristic absorption spectra of all the gases into a set of individual ones. A mathematical model is to estimate fitting coefficients $\{a_j\}$ from a linear model $J(\lambda) = \sum_j^M a_j \cdot y_j(\lambda) + \eta(\lambda)$, where the data $J(\lambda)$ and reference spectra $\{y_j(\lambda)\}$ are given at each wavelength λ and $\eta(\lambda)$ is noise. A challenging complication in practice is wavelength misalignment, *i.e.*, the nominal wavelengths in the measurement $J(\lambda)$ may not correspond exactly to those in the basis $y_j(\lambda)$. We must allow for small deformations $v_j(\lambda)$ so that $y_j(\lambda + v_j(\lambda))$ are all aligned with the data $J(\lambda)$. Taking into account wavelength misalignment, the data model becomes

$$J(\lambda) = \sum_j^M a_j \cdot y_j(\lambda + v_j(\lambda)) + \eta(\lambda). \tag{16}$$

Esser et al. [10] construct an incremented dictionary by deforming each y_j with a set of possible deformations for the DOAS problem. Specifically, since it has been discovered that the deformations can be well approximated by linear functions, *i.e.*, $v_j(\lambda) = p_j \lambda + q_j$, all the possible deformations are enumerated by choosing p_j, q_j from two pre-determined sets $\{P_1, \dots, P_K\}, \{Q_1, \dots, Q_L\}$. Let Y_j be a matrix with each column be a deformed basis

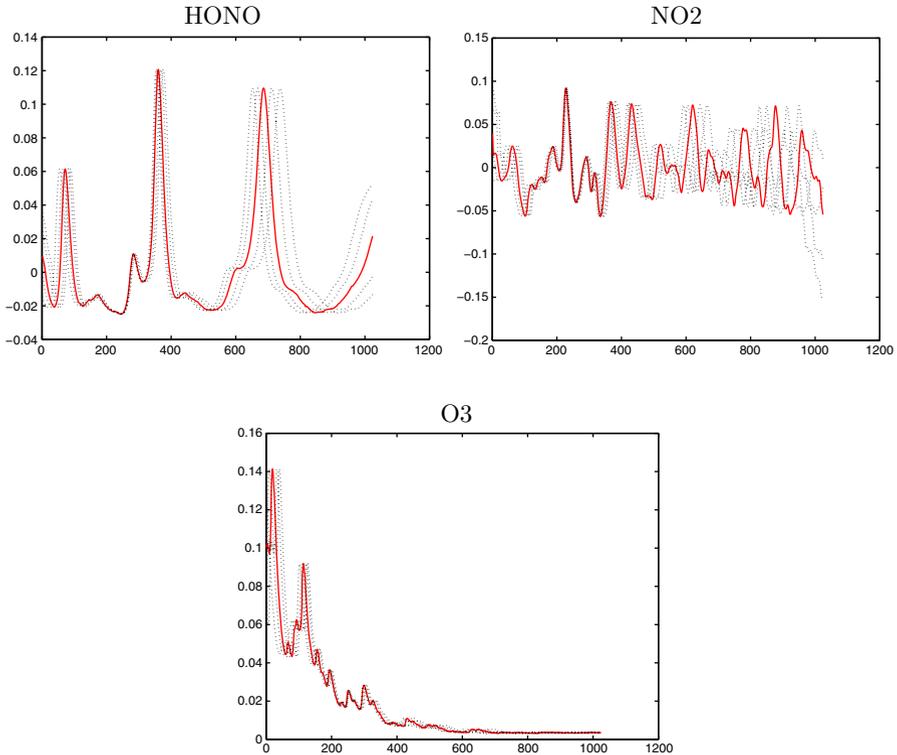


Fig. 5 Each gas spectrum (a dictionary element) is plotted in *red*, while four deformed spectra (nearby dictionary elements) are in *dotted black*

corresponding to $y_j(\lambda)$, i.e., $y_j(\lambda + P_k\lambda + Q_l)$ for $k = 1, \dots, K$ and $l = 1, \dots, L$. Then the model (16) can be rewritten in terms of a matrix-vector form,

$$J = [Y_1, \dots, Y_M] \begin{bmatrix} \mathbf{a}_1 \\ \vdots \\ \mathbf{a}_M \end{bmatrix} + \eta, \tag{17}$$

where \mathbf{a}_j is a (KL) -dimensional column vector.

In our experiments, we generate such a dictionary by taking three given reference spectra $y_j(\lambda)$ for the gases HONO, NO2 and O3 and deforming each by a set of linear functions. The represented wavelengths in nanometers are $\lambda = 340 + 0.04038w$, $w = 0, \dots, 1023$, thus yielding each $y_j \in \mathbb{R}^{1024}$. We choose two pre-determined sets $P_k = -1.01 + 0.01k$ ($k = 1, \dots, 21$) and $Q_l = -1.1 + 0.1l$ ($l = 1, \dots, 21$), and hence there are a total of 441 linearly deformed references for each of the three groups. In Fig. 5, we plot the reference spectra of these three gases together with four deformed examples. The coherence of the resulting dictionary is .9996.

To generate synthetic data $J(\lambda) \in \mathbb{R}^W$, we randomly select one element for each group with random magnitude plus additive zero mean Gaussian noise. Mimicking the relative magnitudes of a real DOAS dataset [12] after normalization of the dictionary, the random magnitudes are chosen to be at different orders with mean values of 1, 0.1, 1.5 for HONO, NO2 and O3 respectively. We consider four different noise levels, whose standard deviations σ are 0, .001, .005 and .01 respectively.

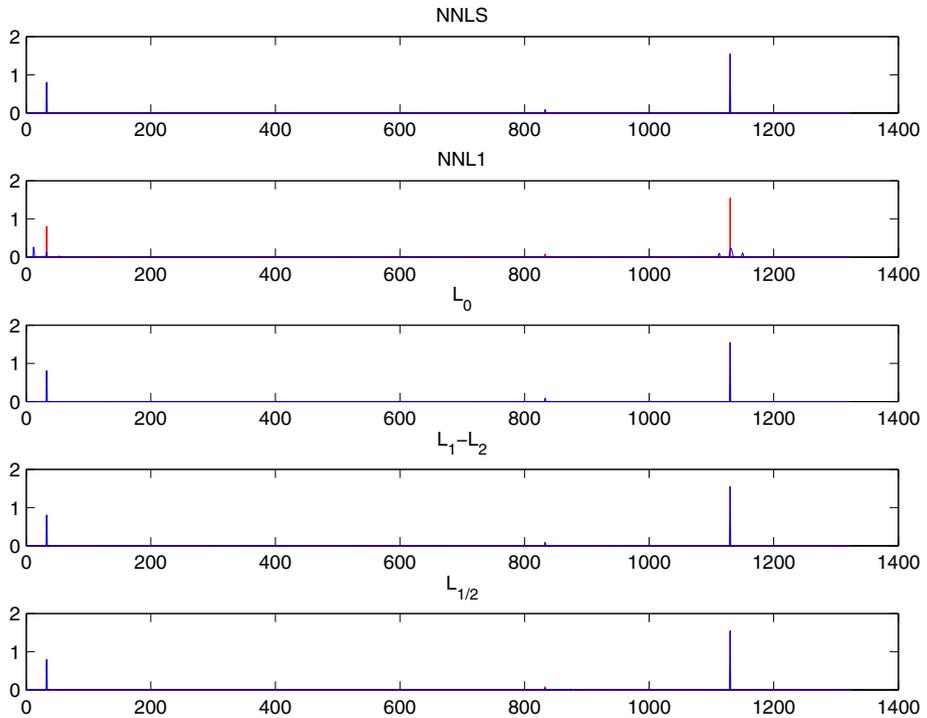


Fig. 6 Method comparisons on synthetic DOAS data without noise. Computed coefficients (blue) are plotted on top of the ground truth (red)

To solve the wavelength misalignment, the following minimization model is considered,

$$\arg \min_{\mathbf{a}_j} \|J - [Y_1, \dots, Y_M] \begin{bmatrix} \mathbf{a}_1 \\ \vdots \\ \mathbf{a}_M \end{bmatrix}\|^2, \tag{18}$$

$$\text{s.t. } \mathbf{a}_j \geq 0, \|\mathbf{a}_j\|_0 \leq 1 \quad j = 1, \dots, M. \tag{19}$$

The second constraint in (19) is to enforce each \mathbf{a}_j to have at most one non-zero element. $\|\mathbf{a}_j\|_0 = 1$ indicates the existence of the gas with a spectrum y_j and its non-zero component corresponds to the selected deformation.

We consider a generic method to find sparse coefficients from the least square model (18). To enforce sparsity, we use $L_1 - L_2$ and compare it with L_p for $p = 1/2$. When taking the non-negativity into account, we look at non-negative least square (NNLS) and non-negative constrained L_1 (NNL1) as comparison. We use MATLAB’s `lsqnonneg` function, which is parameter free, to solve the NNLS. The constrained NNL1 is modelled as,

$$\min_{x \geq 0} \|x\|_1 \quad \text{such that} \quad \|Ax - b\| \leq \tau, \tag{20}$$

which can be solved by Bregman iteration [27]. We also compare with a direct L_0 approach that takes advantages of the structured sparsity. This method is based on the idea of penalty decomposition [17], and therefore it requires a good initialization and slowly increases penalty parameter ρ^k in (3).

In Figs. 6 and 7, we plot the results of different methods in blue along with the ground truth solution in red for $\sigma = 0$ and 0.005 respectively. Table 3 shows the relative errors between the

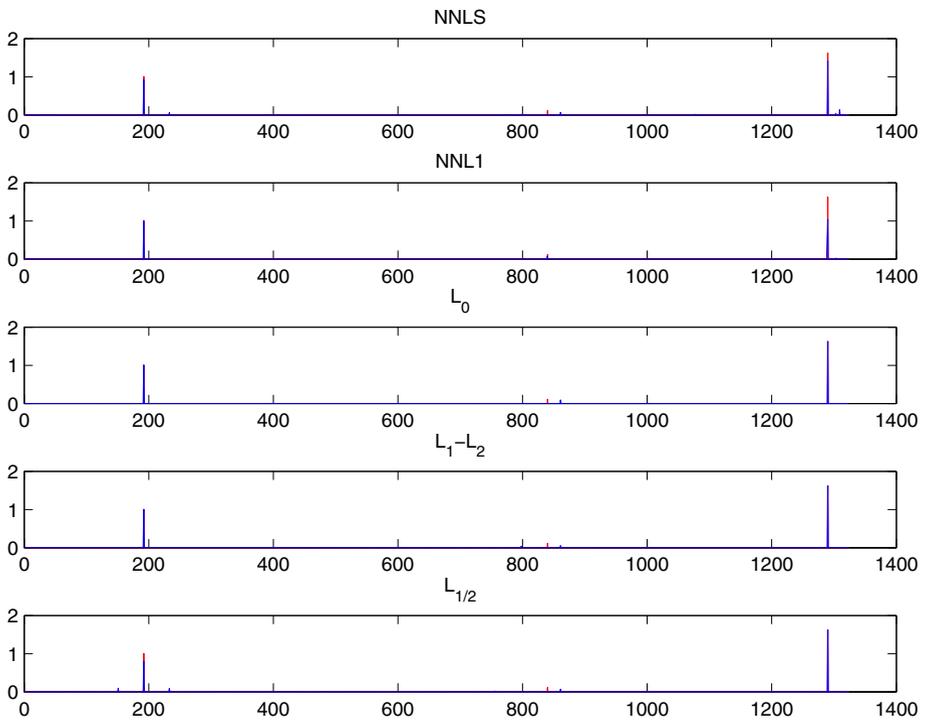


Fig. 7 Method comparisons on synthetic DOAS data with additive noise $\sigma = .005$. Computed coefficients (blue) are plotted on top of the ground truth (red)

Table 3 Relative errors for each method under different amounts of noise. Each recorded value is the mean of 100 random realizations

Noise std	NNLS	NN L1	L0	$L_1 - L_2$	$L_{1/2}$
0	0.00	0.90	0.00	0.0000	0.04
0.001	0.04	0.87	0.005	0.003	0.06
0.005	0.19	0.16	0.06	0.058	0.18
0.01	0.39	0.33	0.20	0.34	0.40

The bold number indicates the smallest error among all the methods, which corresponds to the best reconstruction result

reconstructed vector and the ground-truth under different amounts of noise. Each recorded value is the average of 100 random realizations. All the results demonstrate that $L_1 - L_2$ is comparable to other methods without additional assumption on non-negativity or structured sparsity.

4.3 Image Denoising

An image denoising model [9] interprets an image as linear combinations of local overcomplete bases, where the vector of coefficients is *sparse*, so at any location only few bases

contribute to the approximation. Suppose the discrete image patches of size $\sqrt{n} \times \sqrt{n}$ pixels, ordered lexicographically as column vectors $x \in \mathbb{R}^n$, then the sparsity assumption corresponds to assuming the existence of a matrix $D \in \mathbb{R}^{n \times K}$, the “dictionary,” such that every image patch x can be represented as a linear combination of its columns with a vector of coefficients that has small L_0 norm. If we measure y , a version of x corrupted by additive Gaussian noise with standard deviation σ , then the maximum a-posteriori estimator of the “denoised” patch x is given by $D\hat{\alpha}$, where

$$\hat{\alpha} = \arg \min_{\alpha} \|\alpha\|_0 \quad \text{s.t.} \quad \|D\alpha - y\|_2^2 \leq T, \tag{21}$$

where T is dictated by σ . If one wishes to encode a larger image X of size $\sqrt{N} \times \sqrt{N}$ ($N \gg n$), with a given dictionary $D \in \mathbb{R}^{n \times K}$, a natural approach is to use a block-coordinate relaxation.

$$\hat{X} = \arg \min_{X, \alpha_{ij}} \|X - Y\|_2^2 + \lambda \sum_{i,j} \|\alpha_{ij}\|_0 + \mu \sum_{i,j} \|D\alpha_{ij} - R_{ij}X\|_2^2. \tag{22}$$

The first term measures the fidelity between the measured image Y and its denoised (and unknown) version X . The second term enforces sparsity of each patch; the $n \times N$ matrix R_{ij} extracts the (i, j) th block from the image. A simple denoising algorithm [9] goes as follows,

1. Given an overcomplete dictionary D and let X be noisy data Y
2. Compute the coefficients α_{ij} for each patch $R_{ij}X$

$$\hat{\alpha}_{ij} = \arg \min_{\alpha} \|\alpha\|_0 \quad \text{s.t.} \quad \|D\alpha - R_{ij}X\|_2^2 \leq T. \tag{23}$$

3. Update X by

$$X = \frac{Y + \mu \sum_{i,j} R_{ij}^T D \alpha_{ij}}{I + \mu \sum_{i,j} R_{ij}^T R_{ij}}, \tag{24}$$

which is a simple averaging of shifted patches.

The aforementioned DCT basis (15) can be a candidate for such dictionary, but it is not suitable to represent natural image patches. Aharon et al. propose a dictionary learning technique called K-SVD [1]. They construct a global dictionary that is trained from a large number of natural images. They also consider an adaptive dictionary by training on random samples of the noisy data so that the dictionary is more tailored to the data. The global dictionary is presented in Fig. 8, which is used in our denoising experiments.

In the denoising model [9], the sparse coding step (23) is solved by OMP [23], which is described in Algorithm 4. We can replace OMP by enforcing sparse penalties L_1 , L_p or $L_1 - L_2$. The results are presented in Figs. 9 and 10 for $\sigma = 20, 30$ respectively. The peak-signal-to-noise (PSNR) is provided for quantitative comparison. We find L_0 by OMP and its approximation L_p for $p=1/4$ outperforms L_1 and $L_1 - L_2$ objectively in terms of PSNR. Perceptually, $L_1 - L_2$ appears to leave fewer defects on Lena’s face and elsewhere than L_0 and L_p with a rather close PSNR value. The reasons can be twofold. First, the dictionary is not as coherent as the previous examples, (its coherence is .9559). Second, there is no



Fig. 8 A global dictionary shown on the *left* is obtained from [1]. Each small blob is an 8×8 patch, which corresponds to one column in the dictionary. A test image is shown on the *right*. It contains 4 sub images with different features

Algorithm 4 Orthogonal matching pursuit [23]

1. Start from vector b and initialize the residual $R^i = b$ for $i = 1$.
2. Select the atom that maximizes the absolute value of the inner product of columns of A with R^i .
3. Form a matrix, Φ , with previously selected atoms as the columns. Define the orthogonal projection operator onto the span of the columns of Φ

$$P = \Phi(\Phi^* \Phi)^{-1} \Phi^*.$$

4. Apply the orthogonal projection operator to the residual and update

$$R^{i+1} = (I - P)R^i.$$

5. Let $i = i + 1$ and go to Step 2; stop if s atoms are chosen.
-

ground-truth sparsest solution in this case. Many solutions may look reasonable to human perception. Visually and objectively via PSNR, $L_1 - L_2$ is always better than L_1 .

5 Conclusions and Future Work

In this paper, we studied $L_1 - L_2$ as an alternative to L_1 for sparse representation. We addressed several analytical properties of $L_1 - L_2$ to promote sparsity. We proposed to compute sparse coefficients based on the difference of convex algorithm. Due to its nonconvex nature, we further considered a simulated annealing framework to approach a global solution. We have conducted an extensive study comparing sparse penalties, $L_0, L_1, L_p, L_1 - L_2$, and their numerical algorithms. Experiments have demonstrated that $L_1 - L_2$ is better than L_1 as a sparse regularization, especially when the sensing matrix or dictionary is highly coherent, and the DCA of $L_1 - L_2$ is better than iterative reweighed strategies for L_p minimization.



Fig. 9 Denoising comparison with additive noise whose standard deviation is 20

For future work, we plan to pursue the following directions. As the DCA of $L_1 - L_2$ is built upon L_1 , it is interesting to study whether the support of $L_1 - L_2$ solution is within the one of L_1 . It has been noticed that not only does exact recovery of a sparse signal depend on the matrix, it also depends on the signal itself. We want to characterize exact sparse recovery in terms of minimum separation and/or Rayleigh length. Finally we want to investigate simulated annealing for this specific algorithm, such as how to advance to next step, how to determine the cooling strategy, and how to converge faster.

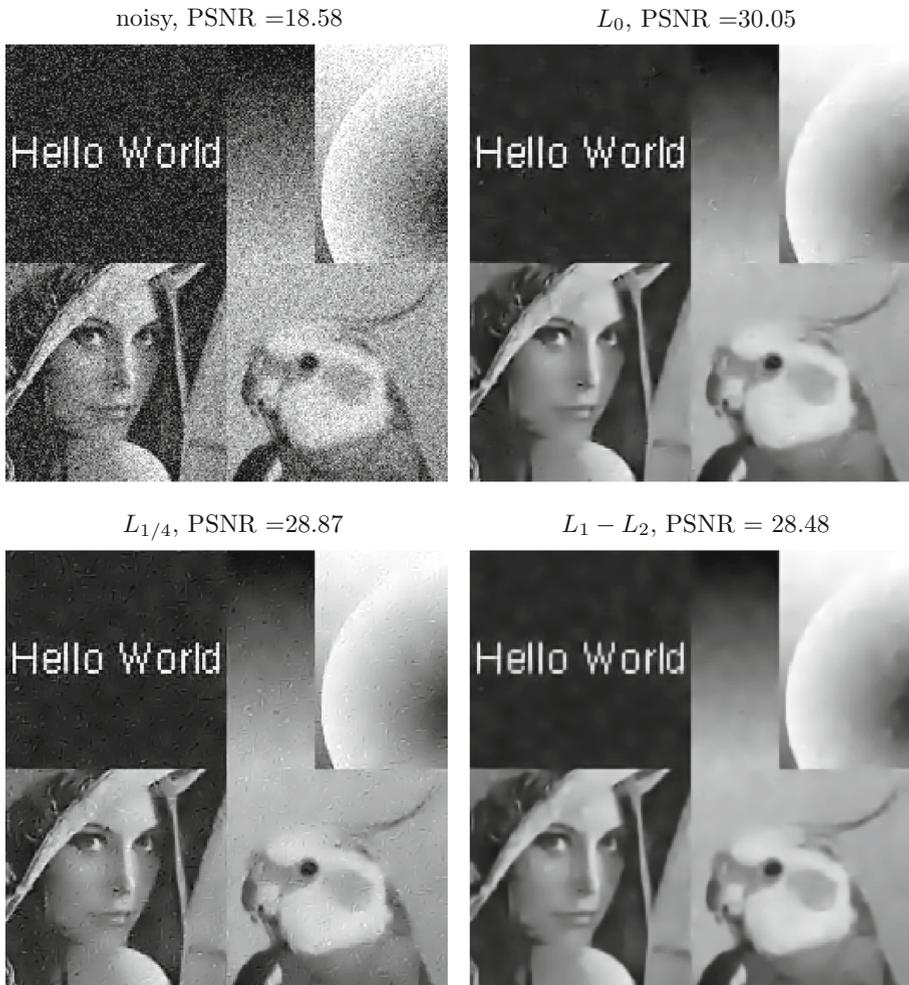


Fig. 10 Denoising comparison with additive noise whose standard deviation is 30. There are more noticeable defects in L_0 image than in $L_1 - L_2$ image

Acknowledgments We thank Dr. Wotao Yin of the Department of Mathematics, UCLA, for providing us with Matlab codes of L_p minimization algorithms published in [7, 16].

References

1. Aharon, M., Elad, M., Bruckstein, A.: K-SVD: an algorithm for designing overcomplete dictionaries for sparse representation. *IEEE Trans. Signal Process.* **54**(11), 4311–4322 (2006)
2. Boyd, S., Parikh, N., Chu, E., Peleato, B., Eckstein, J.: Distributed optimization and statistical learning via the alternating direction method of multipliers. *Found. Trends Mach. Learn.* **3**(1), 1–122 (2011)
3. Candès, E., Elder, Y., Needle, D., Randall, P.: Compressed sensing with coherent and redundant dictionaries. *Appl. Comput. Harmon. Anal.* **31**, 59–73 (2011)
4. Candès, E.J., Fernandez-Granda, C.: Super-resolution from noisy data. *J. Fourier Anal. Appl.* **19**(6), 1229–1254 (2013)

5. Candés, E.J., Tao, T.: Decoding by linear programming. *IEEE Trans. Inf. Theory* **51**(2), 4203–4215 (2005)
6. Carnevali, P., Coletti, L., Patarnello, S.: Image processing by simulated annealing. *IBM J. Res. Dev.* **29**(6), 569–579 (1985)
7. Chartrand, R., Yin, W.: Iteratively reweighted algorithms for compressive sensing. In: *International Conference on Acoustics, Speech, and Signal Processing*, pp. 3869–3872, (2008)
8. Donoho, D., Elad, M.: Optimally sparse representation in general (nonorthogonal) dictionaries via l_1 minimization. *Proc. Natl. Acad. Sci. USA* **100**, 2197–2202 (2003)
9. Elad, M., Aharon, M.: Image denoising via sparse and redundant representations over learned dictionaries. *IEEE Trans. Image Process.* **15**(12), 3736–3745 (2006)
10. Esser, E., Lou, Y., Xin, J.: A method for finding structured sparse solutions to non-negative least squares problems with applications. *SIAM J. Imaging Sci.* **6**(4), 2010–2046 (2013)
11. Fannjiang, A., Liao, W.: Coherence pattern-guided compressive sensing with unresolved grids. *SIAM J. Imaging Sci.* **5**(1), 179–202 (2012)
12. Finlayson-Pitts, B.: Unpublished data. Provided by L. Wingen (2000)
13. Geman, S., Geman, D.: Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images. *IEEE Trans. Pattern Anal. Mach. Intell.* **6**, 721–741 (1984)
14. Gidas, B.: Nonstationary Markov chains and convergence of the annealing algorithm. *J. Stat. Phys.* **39**(1–2), 73–131 (1985)
15. Kirkpatrick, S., Gelatt, C.D., Vecchi, M.P.: Optimization by simulated annealing. *Science* **220**(4598), 671–680 (1983)
16. Lai, M.J., Xu, Y., Yin, W.: Improved iteratively reweighted least squares for unconstrained smoothed l_q minimization. *SIAM J. Numer. Anal.* **5**(2), 927–957 (2013)
17. Lu, Z., Zhang, Y.: Penalty decomposition methods for L_0 -norm minimization. *preprint. arXiv:1008.5372v2* [math. OC], 2012
18. Natarajan, B.K.: Sparse approximate solutions to linear systems. *SIAM J. Comput.* **24**(2), 227–234 (1995)
19. Olshausen, B., Field, D.: Sparse coding with an overcomplete basis set: a strategy employed by v1? *Vision Res.* **37**, 3311–3325 (1997)
20. Platt, U., Stutz, J.: *Differential Optical Absorption Spectroscopy: Principles and Applications*. Springer, Berlin (2008)
21. Tao, P.D., An, L.T.H.: Convex analysis approach to d.c. programming: theory, algorithms and applications. *Acta Math. Vietnam.* **22**(1), 289–355 (1997)
22. Tibshirani, R.: Regression shrinkage and selection via the lasso. *J. R. Stat. Soc. Series B Stat. Methodol.* **58**(1), 267–288 (1996)
23. Tropp, J.: Greed is good: algorithmic results for sparse approximation. *IEEE Trans. Inf. Theory* **50**, 2231–2242 (2004)
24. Xu, F., Wang, S.: A hybrid simulated annealing thresholding algorithm for compressed sensing. *Signal Process.* **93**, 1577–1585 (2013)
25. Yin, P., Esser, E., and Xin, J.: Ratio and difference of l_1 and l_2 norms and sparse representation with coherent dictionaries. Technical report, UCLA CAM Report [13-21] (2013)
26. Yin, P., Lou, Y., He, Q., and Xin, J.: Minimization of $l_1 - l_2$ for compressed sensing. Technical report, UCLA CAM Report [14-01] (2014)
27. Yin, W., Osher, S., Goldfarb, D., Darbon, J.: Bregman iterative algorithms for l_1 minimization with applications to compressed sensing. *SIAM J. Imaging Sci.* **1**, 143–168 (2008)