

# Fast L1-Minimization Algorithms and An Application in Robust Face Recognition: A Review

*Allen Yang  
Arvind Ganesh  
Shankar Sastry  
Yi Ma*



Electrical Engineering and Computer Sciences  
University of California at Berkeley

Technical Report No. UCB/EECS-2010-13

<http://www.eecs.berkeley.edu/Pubs/TechRpts/2010/EECS-2010-13.html>

February 5, 2010

Copyright © 2010, by the author(s).  
All rights reserved.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission.

#### Acknowledgement

This work was partially supported by ARO MURI W911NF-06-1-0076 and ARL MAST-CTA W911NF-08-2-0004. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the Army Research Laboratory or the U.S. Government. The U.S. Government is authorized to reproduce and distribute for Government purposes notwithstanding any copyright notation hereon.

# Fast $\ell_1$ -Minimization Algorithms and An Application in Robust Face Recognition: A Review

Allen Y. Yang, Arvind Ganesh, S. Shankar Sastry, and Yi Ma

**Abstract**— $\ell_1$ -minimization solves the minimum  $\ell_1$ -norm solution to an underdetermined linear system  $y = Ax$ . It has recently received much attention, mainly motivated by the new compressive sensing theory that shows that under certain conditions an  $\ell_1$ -minimization solution is also the sparsest solution to that system. Although classical solutions to  $\ell_1$ -minimization have been well studied in the past, including *primal-dual interior-point methods* and *orthogonal matching pursuit*, they suffer from either expensive computational cost or insufficient estimation accuracy in many real-world, large-scale applications. In the past five years, many new algorithms have been proposed. We provide a comprehensive review of five representative approaches, namely, *gradient projection*, *homotopy*, *iterative shrinkage-thresholding*, *proximal gradient*, and *alternating direction*. The repository is intended to fill in a gap in the existing literature to systematically benchmark the performance of these algorithms using a consistent experimental setting. In addition, the experiment will be focused on the application of robust face recognition, where a sparse representation framework has recently been developed to recover human identities from facial images that may be affected by illumination, occlusion, and facial disguise. The paper also provides useful guidelines to practitioners working in similar fields.

## I. INTRODUCTION

$\ell_1$ -minimization ( $\ell_1$ -min) has been one of the hot topics in the signal processing and optimization communities in the last five years or so. In compressive sensing (CS) theory [9], [16], [8], it has been shown to be an efficient approach to recover the sparsest solutions to certain underdetermined systems of linear equations. More specifically, assuming there exists an unknown signal  $x_0 \in \mathbb{R}^n$ , a measurement vector  $b \in \mathbb{R}^d$  can be generated by a linear projection  $b = Ax$ . If we assume the sensing matrix  $A$  to be full-rank and overcomplete, i.e.,  $d < n$ , an  $\ell_1$ -min program solves the following convex optimization problem

$$(P_1): \quad \min \|x\|_1 \text{ subject to } b = Ax. \quad (1)$$

The formulation of  $(P_1)$  constitutes a linear inverse problem, as the number of measurements in  $b$  is smaller than the

This work was partially supported by ARO MURI W911NF-06-1-0076 and ARL MAST-CTA W911NF-08-2-0004. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the Army Research Laboratory or the U.S. Government. The U.S. Government is authorized to reproduce and distribute for Government purposes notwithstanding any copyright notation hereon.

A. Yang and S. Sastry are with the Department of Electrical Engineering and Computer Sciences, University of California, Berkeley, USA. A. Ganesh, and Y. Ma are with the Coordinated Science Laboratory, University of Illinois, Urbana, USA. Y. Ma is also with the Visual Computing Group, Microsoft Research Asia, Beijing, China. Corresponding author: Allen Yang, Cory Hall, University of California, Berkeley, CA 94720. Email: yang@eecs.berkeley.edu. Tel: 1-510-643-5798. Fax: 1-510-643-2356.

number of unknowns in  $x$ . However, CS theory shows that if  $x_0$  is sufficiently sparse and the sensing matrix  $A$  is *incoherent* with the basis under which  $x_0$  is sparse (i.e., the identity matrix in its standard form (1)),  $x_0$  can be exactly recovered. This sparsity-seeking property of  $(P_1)$  has been shown to have tremendous applications in geophysics, data compression, image processing, sensor networks, and most recently computer vision. The reader is referred to [9], [2], [8] for a comprehensive review of these applications.

Traditionally,  $(P_1)$  has been formulated as a linear programming (LP) problem, such as in basis pursuit (BP) [11]. However, one can show that the computational complexity of these general-purpose algorithms is often too high for many real-world, large-scale applications. Alternatively, heuristic greedy algorithms have been developed to approximate  $(P_1)$ , which are also significantly faster than using LP. *Orthogonal matching pursuit* (OMP) [14] and *least angle regression* (LARS) [18] are two well-known algorithms in this category. Empirically, these greedy algorithms often can find sufficiently good solutions to approximate  $(P_1)$ . However, they may also fail in some conditions (one negative example for OMP is discussed in [38]).

In practice, the exact constraint  $b = Ax$  is often relaxed to take into account the existence of measurement errors in the sensing process:

$$b = Ax + e. \quad (2)$$

Particularly, if the error term  $e$  is assumed to be white noise such that  $\|e\|_2 \leq \epsilon$ , the ground truth signal  $x_0$  can be well approximated by the so-called *basis pursuit denoising* (BPDN) [11], [10]:

$$(P_{1,2}): \quad \min \|x\|_1 \text{ subject to } \|b - Ax\|_2 \leq \epsilon. \quad (3)$$

Based on the nature of the measurement noise, the  $\ell_2$ -norm used in the penalty term can be replaced by other  $\ell_p$ -norms. For example, the following  $(P_{1,1})$  program has been considered in [43], [42], [45]:

$$(P_{1,1}): \quad \min \|x\|_1 \text{ subject to } \|b - Ax\|_1 \leq \epsilon. \quad (4)$$

As we will discuss further in Section IV, the problem  $(P_{1,1})$  assumes the measurement  $b$  may be corrupted by large and impulsive noise  $e$ , which itself may also be sparse.

In light of the high interest in finding more efficient algorithms to solve these problems, many new algorithms have been recently proposed. Although it is impossible to summarize all existing algorithms in the literature, in this paper, we provide a comprehensive review of five representative methods, namely, *gradient projection* [19], [29], *homotopy* [37],

[31], [17], *iterative shrinkage-thresholding* [13], [12], [24], [44], *proximal gradient* (also known as Nesterov's method) [34], [35], [5], [6], and *alternating direction* [45]. Unless stated otherwise, all algorithms are formulated to recover the approximate minimum  $\ell_1$ -norm solution ( $P_{1,2}$ ). It is easy to see that if  $\epsilon \rightarrow 0$  in (3), the solution becomes a good estimate of the basic BP problem ( $P_1$ ).

The paper intends to fill in a gap in the existing literature to systematically benchmark the performance of these algorithms using a fair and consistent experimental setting. Due to the attention given to compressive sensing and  $\ell_1$ -minimization in the community, there should be no surprise that other more advanced solutions will be conceived and studied in the near future, and we do not believe there exists an overall winner that could achieve the best performance in speed and accuracy for all applications. Therefore, in addition to extensive simulations on synthetic data, the experiment will be focused on an example application of robust face recognition [43], [42], where a sparse representation framework has recently been developed to recognize human identities from facial images, which may be affected by illumination, occlusion, and facial disguise. We have made the documentation of the  $\ell_1$ -min resources and benchmark scripts in MATLAB online at <http://www.eecs.berkeley.edu/~yang/software/11benchmark/>, which aims to provide useful references and guidelines to practitioners working in similar fields.

#### A. Notation

For a vector  $\mathbf{x} \in \mathbb{R}^n$ , we denote  $\mathbf{x}_+$  and  $\mathbf{x}_-$  that collect the positive and negative coefficients of  $\mathbf{x}$ , respectively:

$$\mathbf{x} = \mathbf{x}_+ - \mathbf{x}_-, \mathbf{x}_+ \geq 0, \mathbf{x}_- \geq 0. \quad (5)$$

We also denote

$$X = \text{diag}(x_1, x_2, \dots, x_n) \in \mathbb{R}^{n \times n} \quad (6)$$

as a square matrix with the coefficients of  $\mathbf{x}$  as its diagonal and zero otherwise. The concatenation of two (column) vectors will be written following the MATLAB convention:  $[\mathbf{x}_1; \mathbf{x}_2] \doteq \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}$ ;  $[\mathbf{x}_1, \mathbf{x}_2] \doteq [\mathbf{x}_1 \ \mathbf{x}_2]$ . The notation  $\mathbf{1}$  is a vector whose coefficients are all one with dimension defined within the context. In this paper, function  $\|\cdot\|$  without subscript represents the usual  $\ell_2$ -norm.

#### B. Primal-Dual Interior-Point Methods

We first discuss a classical solution to the  $\ell_1$ -min problem ( $P_1$ ), called the *primal-dual interior-point* method, which is usually attributed to the works of [20], [27], [32], [33], [30]. For the sake of simplicity, we assume here that the sparse solution  $\mathbf{x}$  is nonnegative.<sup>1</sup> Under this assumption, it is easy to see that ( $P_1$ ) can be converted to the standard primal and dual forms in linear programming (LP):

$$\begin{array}{ll} \text{Primal (P)} & \text{Dual (D)} \\ \min & \mathbf{c}^T \mathbf{x} \\ \text{s.t.} & A\mathbf{x} = \mathbf{b} \\ & \mathbf{x} \geq 0 \end{array} \quad \begin{array}{ll} \max & \mathbf{b}^T \mathbf{y} \\ \text{s.t.} & A^T \mathbf{y} + \mathbf{z} = \mathbf{c} \\ & \mathbf{z} \geq 0, \end{array} \quad (7)$$

<sup>1</sup>This constraint can be easily removed by considering the corresponding solution for another linear system  $\mathbf{b} = [A, -A][\mathbf{x}_+; \mathbf{x}_-]$ , where  $[\mathbf{x}_+; \mathbf{x}_-]$  is also nonnegative.

where for  $\ell_1$ -min,  $\mathbf{c} = \mathbf{1}$ . The algorithm simultaneously optimizes the primal-dual pair of the linear programming problems ( $P$ ) and ( $D$ ) in the domain:  $(\mathbf{x}, \mathbf{y}, \mathbf{z}) \in \mathbb{R}^n \times \mathbb{R}^d \times \mathbb{R}^n$ .

It was proposed in [20] that ( $P$ ) can be converted to a family of logarithmic barrier problems<sup>2</sup>:

$$(P_\mu) \quad \begin{array}{ll} \min & \mathbf{c}^T \mathbf{x} - \mu \sum_{i=1}^n \log x_i \\ \text{s.t.} & A\mathbf{x} = \mathbf{b}, \mathbf{x} > 0 \end{array} \quad (8)$$

Clearly, a feasible solution  $\mathbf{x}$  for ( $P_\mu$ ) cannot have zero coefficients. Therefore, we define the interiors of the solution domains for ( $P$ ) and ( $D$ ) as:

$$\begin{array}{ll} P_{+++} & = \{\mathbf{x} : A\mathbf{x} = \mathbf{b}, \mathbf{x} > 0\}, \\ D_{+++} & = \{(\mathbf{y}, \mathbf{z}) : A^T \mathbf{y} + \mathbf{z} = \mathbf{c}, \mathbf{z} > 0\}, \\ S_{+++} & = P_{+++} \times D_{+++}, \end{array} \quad (9)$$

and assume the sets are not empty.

Under these assumptions, one can show that problem ( $P_\mu$ ) has a unique global optimal solution  $\mathbf{x}(\mu)$  for all  $\mu > 0$ . As  $\mu \rightarrow 0$ ,  $\mathbf{x}(\mu)$  and  $(\mathbf{y}(\mu), \mathbf{z}(\mu))$  converge to optimal solutions of problems ( $P$ ) and ( $D$ ) respectively [32], [33].

The primal-dual interior-point algorithm seek the domain of the central trajectory for the problems ( $P$ ) and ( $D$ ) in  $S_{+++}$ , where the central trajectory is defined as the set  $S = \{(\mathbf{x}(\mu), \mathbf{y}(\mu), \mathbf{z}(\mu)) : \mu > 0\}$  of solutions to the following system of equations:

$$XZ\mathbf{1} = \mu\mathbf{1}, A\mathbf{x} = \mathbf{b}, A^T \mathbf{y} + \mathbf{z} = \mathbf{c}, \mathbf{x} \geq 0, \text{ and } \mathbf{z} \geq 0. \quad (10)$$

The condition (10) is also known as the Karush-Kuhn-Tucker (KKT) stationary condition for the convex function ( $P_\mu$ )[33], [30]

Hence, the update rule on the current value  $(\mathbf{x}^{(k)}, \mathbf{y}^{(k)}, \mathbf{z}^{(k)})$  is defined by the Newton direction  $(\Delta\mathbf{x}, \Delta\mathbf{y}, \Delta\mathbf{z})$ , which is evaluated from the following equations

$$\begin{cases} Z^{(k)} \Delta\mathbf{x} + X^{(k)} \Delta\mathbf{z} & = \hat{\mu}\mathbf{1} - X^{(k)} \mathbf{z}^{(k)} \\ A\Delta\mathbf{x} & = 0 \\ A^T \Delta\mathbf{y} + \Delta\mathbf{z} & = 0, \end{cases} \quad (11)$$

where  $\hat{\mu}$  is a free penalty parameter that generally is different from  $\mu$  in ( $P_\mu$ ).

In addition to the update rule (11), an algorithm also needs to specify the stopping criterion when the solution is close to the optimum. For  $\ell_1$ -min, some simple rules can be easily evaluated:

- 1) The relative change of the sparse support set becomes small;
- 2) The relative change (in the sense of the  $\ell_2$ -norm) of the update of the estimate becomes small;
- 3) The relative change of the objective function becomes small.

A more detailed discussion about choosing good stopping criteria in different applications is postponed to Section III.

Algorithm 1 summarizes the conceptual implementation of the interior-point methods.<sup>3</sup> For more details about how to choose the initial values  $(\mathbf{x}^{(0)}, \mathbf{y}^{(0)}, \mathbf{z}^{(0)})$  and the penalty parameter  $\hat{\mu}$ , the reader is referred to [30], [33].

<sup>2</sup>In general, any smooth function  $\Psi$  that satisfies  $\Psi(0^+) = -\infty$  is a valid barrier function [25].

<sup>3</sup>A MATLAB implementation of the primal-dual interior-point solver can be found in the SparseLab Toolbox at <http://sparselab.stanford.edu/>.

---

**Algorithm 1** Primal-Dual Interior-Point Algorithm (PDIPA)
 

---

**Input:** A full rank matrix  $A \in \mathbb{R}^{d \times n}$ ,  $d < n$ , a vector  $\mathbf{b} \in \mathbb{R}^d$ , initial guess  $(\mathbf{x}^{(0)}, \mathbf{y}^{(0)}, \mathbf{z}^{(0)})$ . Iteration  $k \leftarrow 0$ . Initial penalty  $\mu$  and a decreasing factor  $0 < \delta < \sqrt{n}$ .

1: **repeat**

2:  $k \leftarrow k + 1$ ,  $\mu \leftarrow \mu(1 - \delta/\sqrt{n})$ .

3: Solve (11) for  $(\Delta\mathbf{x}, \Delta\mathbf{y}, \Delta\mathbf{z})$ .

4:  $\mathbf{x}^{(k)} \leftarrow \mathbf{x}^{(k-1)} + \Delta\mathbf{x}$ ,  $\mathbf{y}^{(k)} \leftarrow \mathbf{y}^{(k-1)} + \Delta\mathbf{y}$ ,  $\mathbf{z}^{(k)} \leftarrow \mathbf{z}^{(k-1)} + \Delta\mathbf{z}$ .

5: **until** stopping criterion is satisfied.

**Output:**  $\mathbf{x}^* \leftarrow \mathbf{x}^{(k)}$ .

---

Algorithm 1 requires a total of  $O(\sqrt{n})$  iterations, and each iteration can be executed in  $O(n^3)$  operations for solving the linear system (11). In one simulation shown in Figure 1, the computational complexity of Algorithm 1 w.r.t. the sensing dimension  $d$  grows much faster than the other six algorithms in comparison. For example, at  $d = 1900$ , the fastest algorithm in this simulation, i.e., iterative thresholding, only takes about 0.2 sec to complete one trial, which is more than 4,000 times faster than PDIPA. Because of this reason, basic BP algorithms should only be used with caution in solving real-world applications.

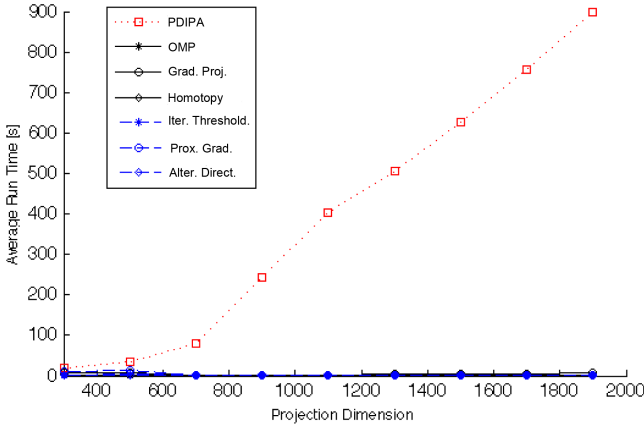


Fig. 1. Average run time of PDIPA in comparison to six other fast implementations under similar estimation accuracy. The simulation setup:  $n = 2000$ ,  $k = 200$ . The projection matrices are randomly generated based on the standard normal distribution with the dimension varies from 300 to 1900. The support of the ground truth  $\mathbf{x}_0$  is randomly selected at each trial, and the nonzero coefficients are sampled from the normal distribution. At each projection dimension, simulation repeats 50 trials.

Next, we will review the five fast  $\ell_1$ -min algorithms shown in Figure 1, namely, *gradient projection* in Section II-A, *homotopy* in Section II-B, *iterative shrinkage-thresholding* in Section II-C, *proximal gradient* in Section II-D, and *alternating direction* in Section II-E. Implementations of all the algorithms are available for download from their respective authors.

## II. FAST $\ell_1$ -MIN ALGORITHMS

### A. Gradient Projection Methods

We first discuss *gradient projection* (GP) methods that seek sparse representation  $\mathbf{x}$  along certain gradient direction,

which induces much faster convergence speed. The approach reformulates the  $\ell_1$ -min as a quadratic programming (QP) problem compared to the LP implementation in PDIPA.

We start with the  $\ell_1$ -min problem  $(P_{1,2})$ . It is equivalent to the so-called LASSO objective function [40]:

$$(P_q) : \min_{\mathbf{x}} \|\mathbf{b} - A\mathbf{x}\|_2^2 \quad \text{s.t.} \quad \|\mathbf{x}\|_1 \leq \sigma. \quad (12)$$

Using the Lagrangian method, the two problems  $(P_{1,2})$  and  $(P_q)$  can be both rewritten as an unconstrained optimization problem:

$$\mathbf{x}^* = \arg \min_{\mathbf{x}} F(\mathbf{x}) = \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{b} - A\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1, \quad (13)$$

where  $\lambda$  is the Lagrangian multiplier.

In the literature, there exist two slightly different methods to engage (13) as a quadratic programming problem, namely, *gradient projection sparse representation* (GPSR) [19] and *truncated Newton interior-point method* (TNIPM) [29].<sup>4</sup>

To formulate the GPSR algorithm, one can separate the positive coefficients  $\mathbf{x}_+$  and the negative coefficients  $\mathbf{x}_-$  in  $\mathbf{x}$ , and rewrite (13) as

$$\begin{aligned} \min \quad Q(\mathbf{x}) &= \frac{1}{2} \|\mathbf{b} - [A, -A][\mathbf{x}_+; \mathbf{x}_-]\|_2^2 + \lambda \mathbf{1}^T \mathbf{x}_+ + \lambda \mathbf{1}^T \mathbf{x}_- \\ \text{s.t.} \quad &\mathbf{x}_+ \geq 0, \mathbf{x}_- \geq 0. \end{aligned} \quad (14)$$

Problem (14) can be rewritten in the standard QP form as

$$\begin{aligned} \min \quad Q(\mathbf{z}) &\doteq \mathbf{c}^T \mathbf{z} + \frac{1}{2} \mathbf{z}^T B \mathbf{z} \\ \text{s.t.} \quad &\mathbf{z} \geq 0, \end{aligned} \quad (15)$$

where  $\mathbf{z} = [\mathbf{x}_+; \mathbf{x}_-]$ ,  $\mathbf{c} = \lambda \mathbf{1} + [-A^T \mathbf{b}; A^T \mathbf{b}]$ , and

$$B = \begin{bmatrix} A^T A & -A^T A \\ -A^T A & A^T A \end{bmatrix}. \quad (16)$$

Notice that the gradient of  $Q(\mathbf{z})$  is defined as

$$\nabla_{\mathbf{z}} Q(\mathbf{z}) = \mathbf{c} + B\mathbf{z}. \quad (17)$$

It leads to a basic algorithm that searches from each iterate  $\mathbf{z}^{(k)}$  along the negative gradient  $-\nabla Q(\mathbf{z})$ :

$$\mathbf{z}^{(k+1)} = \mathbf{z}^{(k)} - \alpha^{(k)} \nabla Q(\mathbf{z}^{(k)}), \quad (18)$$

where  $\alpha^{(k)}$  is the step size that remains undefined. This can be solved by the standard *line-search* process [26]. For example, in [19], a direction vector  $\mathbf{g}^{(k)}$  is defined as

$$g_i^{(k)} = \begin{cases} (\nabla Q(\mathbf{z}^{(k)}))_i, & \text{if } z_i^{(k)} > 0 \text{ or } (\nabla Q(\mathbf{z}^{(k)}))_i < 0 \\ 0, & \text{otherwise.} \end{cases} \quad (19)$$

Then the step size for the update is chosen to be

$$\alpha^{(k)} = \arg \min_{\alpha} Q(\mathbf{z}^{(k)} - \alpha \mathbf{g}^{(k)}), \quad (20)$$

which has a closed-form solution

$$\alpha^{(k)} = \frac{(\mathbf{g}^{(k)})^T \mathbf{g}^{(k)}}{(\mathbf{g}^{(k)})^T B \mathbf{g}^{(k)}}. \quad (21)$$

In terms of the computational complexity, the authors reported that the computational complexity and convergence of

<sup>4</sup>A MATLAB implementation of GPSR is available at <http://www.lx.it.pt/~mtf/GPSR>. A MATLAB Toolbox for TNIPM called LILS is available at [http://www.stanford.edu/~boyd/l1\\_ls/](http://www.stanford.edu/~boyd/l1_ls/).

GPSR is difficult to estimate exactly. Another issue is that the formulation of (15) doubles the dimension of the equations from (13). Therefore, the matrix operations involving  $B$  must take into account its special structure w.r.t.  $A$  and  $A^T$ .

The second GP algorithm, which our benchmark will be based on, is *truncated Newton interior-point method* (TNIPM) [29]. It transforms the same objective function (13) to a quadratic program but with inequality constraints:

$$\begin{aligned} \min \quad & \frac{1}{2} \|A\mathbf{x} - \mathbf{b}\|_2^2 + \lambda \sum_{i=1}^n u_i \\ \text{s.t.} \quad & -u_i \leq x_i \leq u_i, \quad i = 1, \dots, n. \end{aligned} \quad (22)$$

Then a *logarithmic barrier* for the constraints  $-u_i \leq x_i \leq u_i$  can be constructed [20]:

$$\Phi(\mathbf{x}, \mathbf{u}) = - \sum_i \log(u_i + x_i) - \sum_i \log(u_i - x_i). \quad (23)$$

Over the domain of  $(\mathbf{x}, \mathbf{u})$ , the central path consists of the unique minimizer  $(\mathbf{x}^*(t), \mathbf{u}^*(t))$  of the convex function

$$F_t(\mathbf{x}, \mathbf{u}) = t(\|A\mathbf{x} - \mathbf{b}\|_2^2 + \lambda \sum_{i=1}^n u_i) + \Phi(\mathbf{x}, \mathbf{u}), \quad (24)$$

where the parameter  $t \in [0, \infty)$ .

Using the primal barrier method that we discussed in Section I-B, the optimal search direction using Newton's method is computed by

$$\nabla^2 F_t(\mathbf{x}, \mathbf{u}) \cdot \begin{bmatrix} \Delta \mathbf{x} \\ \Delta \mathbf{u} \end{bmatrix} = -\nabla F_t(\mathbf{x}, \mathbf{u}) \in \mathbb{R}^{2n}. \quad (25)$$

Again, for large-scale problems, directly solving (25) is not computationally practical. Then in [29], the authors argued that it can be replaced by an approximate solution using the *pre-conditioned conjugate gradients* (PCG) algorithm. The reader is referred to [28], [36] for more details of the technique.

## B. Homotopy Methods

Both PDIPA and GP require the solution sequence to be close to a "central path", which is often difficult to satisfy and computationally expensive in practice. A natural question arises: *Are there any fast algorithms that are suitable for large-scale applications and yet can recover the sparsest solutions in similar conditions as  $\ell_1$ -min?*

In this section, we overview one such approach called *homotopy methods* [37], [31], [17]. The algorithm has intimate connection with two other greedy  $\ell_1$ -min approximations, namely, *least angle regression* (LARS) [18] and *polytope faces pursuit* (PFP) [38]. For instance, if a  $k$ -sparse signal is sufficiently sparse, all three algorithms can find it in  $k$  iterations. On the other hand, LARS would never remove indices from the current sparse support set, while the general homotopy and PFP include mechanisms to remove coefficients from the sparse support during the iteration. More importantly, the homotopy algorithm provably solves  $\ell_1$ -min ( $P_1$ ), while LARS and PFP are only approximate solutions. A more detailed discussion about homotopy, LARS, and PFP can be found in [17].

Recall that  $(P_{1,2})$  can be written as an unconstrained optimization problem:

$$\begin{aligned} \mathbf{x}^* = \arg \min_{\mathbf{x}} F(\mathbf{x}) &= \arg \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{b} - A\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1, \\ &\doteq \arg \min_{\mathbf{x}} f(\mathbf{x}) + \lambda g(\mathbf{x}) \end{aligned} \quad (26)$$

where  $\lambda$  is the Lagrangian multiplier. On one hand, w.r.t. a fixed  $\lambda$ , the optimal solution is achieved when  $\mathbf{0} \in \partial F(\mathbf{x})$ . On the other hand, similar to the interior-point algorithm, if we define

$$\mathcal{X} \doteq \{\mathbf{x}_\lambda^* : \lambda \in [0, \infty)\}, \quad (27)$$

$\mathcal{X}$  identifies a solution path that follows the change in  $\lambda$ : when  $\lambda \rightarrow \infty$ ,  $\mathbf{x}_\lambda^* = \mathbf{0}$ ; when  $\lambda \rightarrow 0$ ,  $\mathbf{x}_\lambda^*$  converges to the solution of  $(P_1)$ .

The homotopy methods exploit the fact that the objective function  $F(\mathbf{x})$  undergoes a homotopy from the  $\ell_2$  constraint to the  $\ell_1$  objective in (26) as  $\lambda$  decreases. One can further show that the solution path  $\mathcal{X}$  is piece-wise constant as a function of  $\lambda$  [37], [18], [17].<sup>5</sup> Therefore, in constructing a decreasing sequence of  $\lambda$ , it is only necessary to identify those "breakpoints" that lead to changes of the support set of  $\mathbf{x}_\lambda^*$ , namely, either a new nonzero coefficient added or a previous nonzero coefficient removed.

The major obstacle in computing  $\partial F(\mathbf{x})$  is that the  $\ell_1$ -norm term is not globally differentiable. Instead, one can consider the *subdifferential* of a convex function  $g$ , defined as [26]:

$$\partial g(\mathbf{x}) = \{\eta \in \mathbb{R}^n : g(\bar{\mathbf{x}}) - g(\mathbf{x}) \geq \eta^T(\bar{\mathbf{x}} - \mathbf{x}), \forall \bar{\mathbf{x}} \in \mathbb{R}^n\}. \quad (28)$$

If a line passes through  $\mathbf{x}$  and its gradient is within  $\partial g(\mathbf{x})$ , then all the points on the line are either touching or below the convex function  $g$ .

The first summand  $f$  in (26) is differentiable:  $\nabla f = A^T(A\mathbf{x} - \mathbf{b}) \doteq -\mathbf{c}(\mathbf{x})$ . The subdifferential of  $g(\mathbf{x}) = \|\mathbf{x}\|_1$  is the following set:

$$\mathbf{u}(\mathbf{x}) \doteq \partial \|\mathbf{x}\|_1 = \left\{ \mathbf{u} \in \mathbb{R}^n : \begin{array}{l} u_i = \text{sgn}(x_i), x_i \neq 0 \\ u_i \in [-1, 1], x_i = 0 \end{array} \right\}. \quad (29)$$

So the unconventional part of  $\mathbf{u}(\mathbf{x})$  is that when a coordinate  $x_i = 0$ ,  $u_i$  in its subdifferential is not a scalar but a set.

The algorithm operates in an iterative fashion with an initial value  $\mathbf{x}^{(0)} = \mathbf{0}$ . In each iteration w.r.t. a nonzero  $\lambda$ , once we assign  $\partial F(\mathbf{x}) = \mathbf{0}$ :

$$\mathbf{c}(\mathbf{x}) = A^T \mathbf{b} - A^T A \mathbf{x} = \lambda \mathbf{u}(\mathbf{x}). \quad (30)$$

Hence, by the definition (29), we maintain a sparse support set:

$$\mathcal{I} \doteq \{i : |c_i^{(l)}| = \lambda\}. \quad (31)$$

The algorithm computes the update for  $\mathbf{x}^{(k)}$  in terms of the positive/negative directions for its coefficients and the magnitude. Specifically, the update direction  $\mathbf{d}^{(k)}$  on the sparse support is the solution to the following system:

$$A_{\mathcal{I}}^T A_{\mathcal{I}} \mathbf{d}^{(k)}(\mathcal{I}) = \text{sgn}(\mathbf{c}^{(k)}(\mathcal{I})), \quad (32)$$

and the direction is manually set to zero on the coordinates not in  $\mathcal{I}$ . Along the direction indicated by  $\mathbf{d}^{(k)}$ , an update

<sup>5</sup>The reader who is familiar with LARS should see that LARS shares the similar property.

on  $\mathbf{x}$  may lead to a breakpoint where the condition (30) is violated. The first scenario occurs when an element of  $\mathbf{c}$  not in the support set would increase in magnitude beyond  $\lambda$ :

$$\gamma^+ = \min_{i \notin \mathcal{I}} \left\{ \frac{\lambda - c_i}{1 - \mathbf{a}_i^T A_{\mathcal{I}} \mathbf{d}^{(k)}(\mathcal{I})}, \frac{\lambda + c_i}{1 + \mathbf{a}_i^T A_{\mathcal{I}} \mathbf{d}^{(k)}(\mathcal{I})} \right\}. \quad (33)$$

The index that achieves  $\gamma^+$  is denoted as  $i^+$ . The second scenario occurs when an element of  $\mathbf{c}$  in the support set  $\mathcal{I}$  crosses zero, violating the sign agreement:

$$\gamma^- = \min_{i \in \mathcal{I}} \{-x_i/d_i\}. \quad (34)$$

The index that achieves  $\gamma^-$  is denoted as  $i^-$ . Hence, the homotopy algorithm marches to the next breakpoint, and updates the sparse support set by either appending  $\mathcal{I}$  with  $i^+$  or removing  $i^-$ :

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \min\{\gamma^+, \gamma^-\} \mathbf{d}^{(k)}. \quad (35)$$

The algorithm shall terminate when the update approaches to zero. Algorithm 2 summarizes the implementation of the homotopy methods.<sup>6</sup>

---

#### Algorithm 2 Homotopy

---

**Input:** A full rank matrix  $A = [\mathbf{v}_1, \dots, \mathbf{v}_n] \in \mathbb{R}^{d \times n}$ ,  $d < n$ , a vector  $\mathbf{b} \in \mathbb{R}^d$ , initial Lagrangian parameter  $\lambda = 2\|A^T \mathbf{b}\|_{\infty}$ .

- 1: Initialization:  $k \leftarrow 0$ . Find the first support index:  $i = \arg \max_{j=1}^n \|\mathbf{v}_j^T \mathbf{b}\|$ ,  $\mathcal{I} = \{i\}$ .
- 2: **repeat**
- 3:    $k \leftarrow k + 1$ .
- 4:   Solve for the update direction  $\mathbf{d}^{(k)}$  in (32).
- 5:   Compute the sparse support updates (33) and (34):  $\gamma^* \leftarrow \min\{\gamma^+, \gamma^-\}$ .
- 6:   Update  $\mathbf{x}^{(k)}$ ,  $\mathcal{I}$ , and  $\lambda \leftarrow \lambda - \gamma^*$ .
- 7: **until** stopping criterion is satisfied.

**Output:**  $\mathbf{x}^* \leftarrow \mathbf{x}^{(k)}$ .

---

Overall, solving (32) using a Cholesky factorization and the addition/removal of the sparse support elements dominate the computation. Since one can keep track of the rank-1 update of  $A_{\mathcal{I}}^T A_{\mathcal{I}}$  in solving (32) using  $O(d^2)$  operations in each iteration, the computational complexity of the homotopy algorithm is  $O(kd^2 + kdn)$ .

### C. Iterative Shrinkage-Thresholding Methods

The homotopy algorithm employs a more efficient iterative update rule that only involves operations on those submatrices of  $A$  corresponding to the nonzero support of the current vector  $\mathbf{x}$ . However, it may lose its computational competitiveness when the sparsity of  $\mathbf{x}$  grows proportionally with the observation dimension  $d$ . In such scenarios, the complexity may still approach the worst-case upper-bound  $O(n^3)$ . In this section, we discuss *iterative shrinkage-thresholding* (IST) methods [13], [12], [24], [44], whose implementation mainly

involves lightweight operations such as vector operations and matrix-vector multiplications. This is in contrast to most past methods that all involve expensive operations such as matrix factorization and solving linear least squares (LLE).

In a nutshell, IST considers solving  $(P_{1,2})$  as a special case of the following *composite objective function*:

$$\min_{\mathbf{x}} F(\mathbf{x}) \doteq f(\mathbf{x}) + \lambda g(\mathbf{x}), \quad (36)$$

where  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is a smooth and convex function, and  $g: \mathbb{R}^n \rightarrow \mathbb{R}$  as the regularization term is bounded from below but not necessarily smooth nor convex. For  $\ell_1$ -min in particular,  $g$  is also separable, that is,

$$g(\mathbf{x}) = \sum_{i=1}^n g_i(x_i). \quad (37)$$

Clearly, let  $f(\mathbf{x}) = \frac{1}{2}\|\mathbf{b} - A\mathbf{x}\|_2^2$  and  $g(\mathbf{x}) = \|\mathbf{x}\|_1$ , and the objective function (36) becomes the unconstrained BPDN problem.

The update rule to minimize (36) is computed using the linearized function of  $f$  [44], [5]:

$$\begin{aligned} \mathbf{x}^{(k+1)} &= \arg \min_{\mathbf{x}} \{f(\mathbf{x}^{(k)}) + (\mathbf{x} - \mathbf{x}^{(k)})^T \nabla f(\mathbf{x}^{(k)}) \\ &\quad + \frac{1}{2}\|\mathbf{x} - \mathbf{x}^{(k)}\|_2^2 \cdot \nabla^2 f(\mathbf{x}^{(k)}) + \lambda g(\mathbf{x})\} \\ &\approx \arg \min_{\mathbf{x}} \{(\mathbf{x} - \mathbf{x}^{(k)})^T \nabla f(\mathbf{x}^{(k)}) \\ &\quad + \frac{\alpha^{(k)}}{2}\|\mathbf{x} - \mathbf{x}^{(k)}\|_2^2 + \lambda g(\mathbf{x})\} \\ &= \arg \min_{\mathbf{x}} \left\{ \frac{1}{2}\|\mathbf{x} - \mathbf{u}^{(k)}\|_2^2 + \frac{\lambda}{\alpha^{(k)}} g(\mathbf{x}) \right\}, \\ &\doteq G_{\alpha^{(k)}}(\mathbf{x}^{(k)}), \end{aligned} \quad (38)$$

where

$$\mathbf{u}^{(k)} = \mathbf{x}^{(k)} - \frac{1}{\alpha^{(k)}} \nabla f(\mathbf{x}^{(k)}). \quad (39)$$

In (38), the hessian  $\nabla^2 f(\mathbf{x}^{(k)})$  is approximated by a diagonal matrix  $\alpha^{(k)} I$ .

Now if we replace  $g(\mathbf{x})$  in (38) by the  $\ell_1$ -norm  $\|\mathbf{x}\|_1$ , which is a separable function, then  $G(\mathbf{x}^{(k)}, \alpha^{(k)})$  has a closed-form solution w.r.t. each scalar coefficient:

$$x_i^{(k+1)} = \arg \min_{x_i} \left\{ \frac{(x_i - u_i^{(k)})^2}{2} + \frac{\lambda |x_i|}{\alpha^{(k)}} \right\} = \text{soft}(u_i^{(k)}, \frac{\lambda}{\alpha^{(k)}}), \quad (40)$$

where

$$\begin{aligned} \text{soft}(u, a) &\doteq \text{sgn}(u) \max\{|u| - a, 0\} \\ &= \begin{cases} \text{sgn}(u)(|u| - a) & \text{if } |u| > a \\ 0 & \text{otherwise} \end{cases} \end{aligned} \quad (41)$$

is the so-called *soft-thresholding* function [15].

There are two parameters remained to be determined in (38), that is, the regularizing coefficient  $\lambda$  and the coefficient  $\alpha^{(k)}$  that approximates the hessian matrix  $\nabla^2 f$ . Several IST algorithms differ in the strategies to pick the parameters in the iteration. For  $\alpha$ , since it is chosen such that  $\alpha I$  mimics the Hessian  $\nabla^2 f$ , we require that  $\alpha^{(k)}(\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}) \approx \nabla f(\mathbf{x}^{(k)}) - \nabla f(\mathbf{x}^{(k-1)})$  in the least-square sense. Hence,

$$\begin{aligned} \alpha^{(k+1)} &= \arg \min_{\alpha} \left\| \alpha(\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}) \right. \\ &\quad \left. - (\nabla f(\mathbf{x}^{(k)}) - \nabla f(\mathbf{x}^{(k-1)})) \right\|_2^2 \\ &= \frac{(\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)})^T (\nabla f(\mathbf{x}^{(k)}) - \nabla f(\mathbf{x}^{(k-1)}))}{(\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)})^T (\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)})}. \end{aligned} \quad (42)$$

This is the so-called *Barzilai-Borwein* equation [3], [39], [44].

<sup>6</sup>A MATLAB implementation [1] can be found at <http://users.ece.gatech.edu/~sasif/homotopy/>.

For choosing  $\lambda$ , instead of using a fixed value, several works have proposed a *continuation* procedure [24], [19], in which (38) is solved for a decreasing sequence of  $\lambda$ . Remember, as we mentioned in Section II-B, (38) recovers the optimal  $\ell_1$ -min solution when  $\lambda \rightarrow 0$ . However, it has been observed that the practical performance degrades by directly solving (38) for small values of  $\lambda$ , which has been dubbed as a ‘‘cold’’ starting point. Instead, *continuation* employs a warm-starting strategy by first solving (38) for a larger value of  $\lambda$ , then decreasing  $\lambda$  in steps towards its desired value.

The iterative shrinkage-thresholding algorithm (ISTA) is summarized in Algorithm 3.<sup>7</sup>

---

**Algorithm 3** Iterative Shrinkage-Thresholding Algorithm (ISTA)

---

**Input:** A full rank matrix  $A \in \mathbb{R}^{d \times n}$ ,  $d < n$ , a vector  $\mathbf{b} \in \mathbb{R}^d$ , Lagrangian  $\lambda_0$ , initial values for  $\mathbf{x}^{(0)}$  and  $\alpha^0$ ,  $k \leftarrow 0$ .

- 1: Generate a reducing sequence  $\lambda_0 > \lambda_1 > \dots > \lambda_N \rightarrow 0$ .
- 2: **for**  $i = 0, 1, \dots, N$  **do**
- 3:    $\lambda \leftarrow \lambda_i$
- 4:   **repeat**
- 5:      $k \leftarrow k + 1$ .
- 6:      $\mathbf{x}^{(k)} \leftarrow G(\mathbf{x}^{(k-1)})$ .
- 7:     Update  $\alpha^{(k)}$  using (42).
- 8:     **until** The objective function  $F(\mathbf{x}^{(k)})$  decreases.
- 9: **end for**

**Output:**  $\mathbf{x}^* \leftarrow \mathbf{x}^{(k)}$ .

---

#### D. Proximal Gradient Methods

Proximal gradient algorithms represent another class of algorithms that solve convex optimization problems in (36). Assume that  $f$  is a smooth convex function with Lipschitz continuous gradient, and  $g$  is a continuous convex function. The principle behind proximal gradient algorithms is to iteratively form quadratic approximations  $Q(\mathbf{x}, \mathbf{y})$  to  $F$  around a carefully chosen point  $\mathbf{y}$ , and to minimize  $Q(\mathbf{x}, \mathbf{y})$  rather than the original cost function  $F$ .

Again, we define  $g(\mathbf{x}) = \|\mathbf{x}\|_1$  and  $f(\mathbf{x}) = \frac{1}{2}\|A\mathbf{x} - \mathbf{b}\|_2^2$ . We note that  $\nabla f(\mathbf{x}) = A^T(A\mathbf{x} - \mathbf{b})$  is Lipschitz continuous with Lipschitz constant  $L_f \doteq \|A\|^2$ .<sup>8</sup> Define  $Q(\mathbf{x}, \mathbf{y})$  as:

$$Q(\mathbf{x}, \mathbf{y}) \doteq f(\mathbf{y}) + \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle + \frac{L_f}{2} \|\mathbf{x} - \mathbf{y}\|^2 + \lambda g(\mathbf{x}). \quad (43)$$

Thus, we have a slightly different problem whose solution gets closer to the solution set of (1) as  $\lambda \rightarrow 0$ .

One can show that  $F(\mathbf{x}) \leq Q(\mathbf{x}, \mathbf{y})$  for all  $\mathbf{y}$ , and

$$\arg \min_{\mathbf{x}} Q(\mathbf{x}, \mathbf{y}) = \arg \min_{\mathbf{x}} \left\{ \lambda g(\mathbf{x}) + \frac{L_f}{2} \|\mathbf{x} - \mathbf{u}\|^2 \right\}, \quad (44)$$

where  $\mathbf{u} = \mathbf{y} - \frac{1}{L_f} \nabla f(\mathbf{y})$  by the same trick used in (38). For the  $\ell_1$ -min problem, (44) has a closed-form solution given by

<sup>7</sup>A MATLAB implementation called *Sparse Reconstruction by Separable Approximation* (SpaRSA) [44] is available at <http://www.lx.it.pt/~mtf/SpaRSA/>.

<sup>8</sup> $\|A\|$  represents the spectral norm of the matrix  $A$ .

the soft-thresholding function:

$$\arg \min_{\mathbf{x}} Q(\mathbf{x}, \mathbf{y}) = \text{soft}\left(\mathbf{u}, \frac{\lambda}{L_f}\right). \quad (45)$$

However, unlike the iterative thresholding algorithm described earlier, we use a smoothed computation of the sequence  $\mathbf{y}_k$ . It has been shown that choosing

$$\mathbf{y}^{(k)} = \mathbf{x}^{(k)} + \frac{t_{k-1} - 1}{t_k} \left( \mathbf{x}^{(k)} - \mathbf{x}^{(k-1)} \right), \quad (46)$$

where  $\{t_k\}$  is a positive real sequence satisfying  $t_k^2 - t_k \leq t_{k-1}^2$ , achieves an accelerated non-asymptotic convergence rate of  $O(k^{-2})$  [34], [35], [5]. To further accelerate the convergence of the algorithm, one can make use of the continuation technique: rather than applying the proximal gradient algorithm directly to (36), we vary  $\lambda$ , starting from a large initial value  $\lambda_0$  and decreasing it with each iteration. Although the theoretical convergence rate is unchanged, it has been shown that this greatly reduces the number of iterations in practice.

Finally, for large problems, it is often computationally expensive to directly compute  $L_f = \|A\|^2$ .<sup>9</sup> A *backtracking* line-search strategy [5] can be used to generate a scalar sequence  $\{L_k\}$  that approximates  $L_f$ . We define

$$Q_L(\mathbf{x}, \mathbf{y}) \doteq f(\mathbf{y}) + (\mathbf{x} - \mathbf{y})^T \nabla f(\mathbf{y}) + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2 + \lambda g(\mathbf{x}). \quad (47)$$

Suppose that  $\eta > 1$  is a pre-defined constant. Then, given  $\mathbf{y}^{(k)}$  at the  $k$ th iteration, we set  $L_k = \eta^j L_{k-1}$ , where  $j$  is the smallest nonnegative integer such that the following inequality holds:

$$F(G_{L_k}(\mathbf{y}^{(k)})) \leq Q_{L_k}(G_{L_k}(\mathbf{y}^{(k)}), \mathbf{y}^{(k)}), \quad (48)$$

where  $G_L(\mathbf{y}) \doteq \arg \min_{\mathbf{x}} Q_L(\mathbf{x}, \mathbf{y}) = \text{soft}\left(\mathbf{u}, \frac{\lambda}{L}\right)$  for  $\mathbf{u} \doteq \mathbf{y} - \frac{1}{L} \nabla f(\mathbf{y})$ .

The algorithm, dubbed FISTA in [5], is summarized as Algorithm 4.<sup>10</sup> The convergence behavior of FISTA is given by

$$F(\mathbf{x}^{(k)}) - F(\mathbf{x}^*) \leq \frac{2L_f \|\mathbf{x}^{(0)} - \mathbf{x}^*\|^2}{(k+1)^2}, \quad \forall k. \quad (49)$$

The interested reader may refer to [35], [5], [6] for a proof of the above result.

#### E. Alternating Direction Methods

After the proximal gradient theory, many investigators believe that most relevant convex-optimization techniques have been exhausted in solving the  $\ell_1$ -min problem. After all, in the noiseless case, Homotopy provably solves  $(P_1)$  in  $k$  steps if the underlying sparse signal  $\mathbf{x}_0$  has only  $k$ -nonzeros; in the noisy case, proximal gradient algorithms provide a first-order solution that converges in the order of  $O(k^{-2})$ . Yet, as we are writing this paper, we become aware of a new development, called the *alternating direction method* (ADM) [45].

<sup>9</sup>This problem occurs in the IST algorithm as well.

<sup>10</sup>An implementation of FISTA can be download from the website of the paper: <http://www.eecs.berkeley.edu/~yang/software/11benchmark/>. Another Matlab toolbox called NESTA [6] is available at: <http://www.acm.caltech.edu/~nesta/>.



---

**Algorithm 4** Fast Iterative Shrinkage-Threshold Algorithm (FISTA)
 

---

**Input:**  $\mathbf{b} \in \mathbb{R}^m$ ,  $A \in \mathbb{R}^{m \times n}$ .
1: Set  $\mathbf{x}^{(0)} \leftarrow \mathbf{0}$ ,  $\mathbf{x}^{(1)} \leftarrow \mathbf{0}$ ,  $t_0 \leftarrow 1$ ,  $t_1 \leftarrow 1$ ,  $k \leftarrow 1$ .2: Initialize  $L_0$ ,  $\lambda_1$ ,  $\beta \in (0, 1)$ ,  $\bar{\lambda} > 0$ .3: **while** not converged **do**4:  $\mathbf{y}^{(k)} \leftarrow \mathbf{x}^{(k)} + \frac{t_{k-1}-1}{t_k} (\mathbf{x}^{(k)} - \mathbf{x}^{(k-1)})$ .5: Update  $L_k$  using (48) with  $\mathbf{y}^{(k)}$ .6:  $\mathbf{u}^{(k)} \leftarrow \mathbf{y}^{(k)} - \frac{1}{L_k} A^T (A\mathbf{y}^{(k)} - \mathbf{b})$ .7:  $\mathbf{x}^{(k+1)} \leftarrow \text{soft} \left( \mathbf{u}^{(k)}, \frac{\lambda_k}{L_k} \right)$ .8:  $t_{k+1} \leftarrow \frac{1 + \sqrt{4t_k^2 + 1}}{2}$ .9:  $\lambda_{k+1} \leftarrow \max(\beta \lambda_k, \bar{\lambda})$ .10:  $k \leftarrow k + 1$ .11: **end while****Output:**  $\mathbf{x}^* \leftarrow \mathbf{x}^{(k)}$ .

---

ADM can be attributed to the early works of [22], [21]. Its novelty is a procedure that alternates between optimizing the sparse signal  $\mathbf{x}$  and the residual term  $\mathbf{e}$ :

$$\min_{\mathbf{x}, \mathbf{e}} \|\mathbf{x}\|_1 + \frac{1}{2\mu} \|\mathbf{e}\|^2 \text{ subject to } \mathbf{b} = A\mathbf{x} + \mathbf{e}. \quad (50)$$

Using the Lagrangian method, (50) is converted to an unconstrained form with two additional variables  $\mathbf{y} \in \mathbb{R}^d$  and  $\lambda > 0$ :

$$\min_{\mathbf{x}, \mathbf{y}, \mathbf{e}} \left\{ \|\mathbf{x}\|_1 + \frac{1}{2\mu} \|\mathbf{e}\|^2 + \frac{1}{2\lambda} \|A\mathbf{x} + \mathbf{e} - \mathbf{b}\|^2 - \mathbf{y}^T (A\mathbf{x} + \mathbf{e} - \mathbf{b}) \right\}. \quad (51)$$

We solve (51) using an alternating minimization procedure w.r.t.  $\mathbf{e}^{(k)}$ ,  $\mathbf{x}^{(k)}$ , and  $\mathbf{y}^{(k)}$ , respectively. First, assume  $(\mathbf{x}^{(k)}, \mathbf{y}^{(k)})$  fixed, it is easy to show that the update rule for  $\mathbf{e}$  in (51) is given by:

$$\mathbf{e}^{(k+1)} = \frac{\mu}{\lambda + \mu} (\lambda \mathbf{y}^{(k)} - (A\mathbf{x}^{(k)} - \mathbf{b})). \quad (52)$$

Next, for  $(\mathbf{e}^{(k+1)}, \mathbf{y}^{(k)})$  fixed, the minimization of (51) w.r.t  $\mathbf{x}$  is equivalent to

$$\mathbf{x}^{(k+1)} = \arg \min_{\mathbf{x}} \{ \|\mathbf{x}\|_1 + h(\mathbf{x}) \}, \quad (53)$$

where  $h(\mathbf{x}) \doteq \frac{1}{2\lambda} \|A\mathbf{x} + \mathbf{e}^{(k+1)} - \mathbf{b} - \frac{1}{\lambda} \mathbf{y}^{(k)}\|^2$ . We already know that (53) has a closed-form solution given by the soft-thresholding function:

$$\mathbf{x}^{(k+1)} = \text{soft}(\mathbf{u}^{(k)}, \frac{1}{\alpha\lambda}), \quad (54)$$

where  $\mathbf{u}^{(k)} = \mathbf{x}^{(k)} - \frac{1}{\alpha} \nabla h(\mathbf{x}^{(k)})$  and the hessian  $\nabla^2 h(\mathbf{x}^{(k)})$  is approximated by a diagonal matrix  $\alpha I$ .

Finally, fixing  $(\mathbf{x}^{(k+1)}, \mathbf{e}^{(k+1)})$ , the update rule for the multiplier  $\mathbf{y}$  is

$$\mathbf{y}^{(k+1)} = \mathbf{y}^{(k)} - \gamma \frac{1}{\lambda} (A\mathbf{x}^{(k+1)} + \mathbf{e}^{(k+1)} - \mathbf{b}), \quad (55)$$

where  $\gamma > 0$  is a proper step size.

ADM appears to provide a versatile framework for solving different  $\ell_1$ -min variations. For starters, ADM easily applies to the approximate  $\ell_1$ -min problem  $(P_{1,2})$ . To enforce the  $\ell_2$ -penalty  $\|\mathbf{b} - A\mathbf{x}\| \leq \epsilon$ , we only need to change the update

rule of  $\mathbf{e}$  such that in each iteration the solution is projected onto the  $\ell_2$ -ball  $B_2^\epsilon$  of radius  $\epsilon$ :

$$\mathbf{e}^{(k+1)} = \mathcal{P}_{B_2^\epsilon} (\lambda \mathbf{y}^{(k)} - (A\mathbf{x}^{(k)} - \mathbf{b})), \quad (56)$$

where  $\mathcal{P}$  denotes the projection operator. Furthermore, ADM can be also applied to the dual problems of  $\ell_1$ -min, which we briefly introduced in Section I-B (see [45] for more details). An implementation of the algorithm, called YALL1, iterates in both the primal and dual spaces to converge to the optimal solutions  $(\mathbf{x}^*, \mathbf{r}^*, \mathbf{y}^*)$ .<sup>11</sup>

Finally, the experiments shown in [45] posed an interesting question regarding choosing a proper  $\ell_p$ -penalty in the BPDN model. Recall in Section I that the common  $\ell_2$ -penalty used in (3) can be replaced by other  $\ell_p$ -norms. In particular, the  $(P_{1,1})$  program

$$(P_{1,1}) : \min \|\mathbf{x}\|_1 \text{ subject to } \|\mathbf{b} - A\mathbf{x}\|_1 \leq \epsilon \quad (57)$$

is a good criterion if the measurement error  $\mathbf{e} = \mathbf{b} - A\mathbf{x}$  is also assumed to be sparse. The  $(P_{1,1})$  program can be reformulated as a standard  $(P_1)$  program:

$$\min_{\hat{\mathbf{x}}} \|\hat{\mathbf{x}}\|_1 \text{ subject to } \mathbf{b} = \hat{A}\hat{\mathbf{x}}, \quad (58)$$

where  $\hat{A} = [A, I]$  and  $\hat{\mathbf{x}} = [\mathbf{x}; \mathbf{e}]$ , in addition to a normalization step to balance the rows or columns of the new matrix  $\hat{A}$  if necessary. Then all the previous  $\ell_1$ -min techniques naturally apply to *simultaneously* recover the underlying signal  $\mathbf{x}_0$  and the error  $\mathbf{e}$ .

The authors in [45] claims that the  $\ell_1$ -penalty improves the estimation over the  $\ell_2$ -penalty when the data may contain large, impulsive data noise. In addition, even without impulsive noise, enforcing the  $\ell_1$ -penalty does not seem to harm the solution quality as long as the data do not contain a large amount of white noise (in which case the  $\ell_2$ -penalty is superior). The observations are intriguing partly because in most practical, real-world applications, the measurement error  $\mathbf{e}$  rarely satisfies a white noise model, and some of its coefficients contain large, impulsive values. In addition,  $(P_{1,1})$  has a connection to a recent work about robust face recognition [43], where the concept of sparse representation is utilized to recognize human identities from facial images. In Section IV, we will discuss and rank the performance of the pervious  $\ell_1$ -min algorithm based on the framework of robust face recognition.

### III. SIMULATION: RANDOM SPARSE SIGNALS

In this section, we present two sets of experiments to benchmark the performance of the five fast  $\ell_1$ -min algorithms on random sparse signals, namely, TNIPM/L1LS, Homotopy, SpaRSA, FISTA, and YALL1, together with the classical OMP algorithm [14], [41], [17], [8] as the baseline. It is important to note however that OMP as a greedy algorithm does not solve the  $\ell_1$ -min problem  $(P_1)$ . The benchmark of PDIPA is not reported in the paper, as its performance markedly lags behind the rest of the fast algorithms (as shown in Figure 1).

<sup>11</sup>A MATLAB package for YALL1 is available at <http://www.caam.rice.edu/~yzhang/YALL1/>.

One factor that we should pay special attention to is the stopping criteria used in benchmarking these algorithms. As we first mentioned in Section I-B, choosing a good stopping criterion is important to properly exit an iteration when the estimate becomes close to a local or global optimum. On one hand, in general, straightforward rules do exist, such as the relative change of the objective function:

$$\frac{\|F(\mathbf{x}^{(k+1)}) - F(\mathbf{x}^{(k)})\|}{\|F(\mathbf{x}^{(k)})\|}, \quad (59)$$

or the relative change of the estimate:

$$\frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|}{\|\mathbf{x}^{(k)}\|}. \quad (60)$$

However, their efficacy depends on a proper step size of the update rule: If the step size is poorly chosen, the algorithm may terminate prematurely when the solution is still far away from the optimum. On the other hand, certain special criteria are more effective to some algorithms than the others. For example, for PDIPA, it is natural to use the (relative) duality gap between the primal and dual solutions; for Homotopy, it is easy to measure the relative change of the sparse support as the stopping criterion, as in each iteration a certain number of coefficients will be added or removed from the sparse support set.

In order to design a fair comparison for the six algorithms, in this section, we will take advantage of the available groundtruth sparse signal  $\mathbf{x}_0$ : if the  $\ell_2$ -norm difference between the  $\ell_1$ -min estimate  $\mathbf{x}^*$  and  $\mathbf{x}_0$  is smaller than a threshold, the iteration should exit.<sup>12</sup> In addition, we set the maximal iteration for all algorithms equal to 1,000. If an algorithm fails to converge to the ground truth, it will quit after it reaches the maximal iteration. All experiments are performed in MATLAB on a Dell PowerEdge 1900 workstation with dual quad-core 2.66GHz Xeon processors and 8GB of memory.

#### A. $\rho$ - $\delta$ Plot in the Noise-Free Case

The first experiment is designed to measure how accurate the algorithms recover exact sparse signals in the noise-free case ( $P_1$ ). A good performance metric is the so-called  $\rho$ - $\delta$  plot, where the sparsity rate  $\rho = k/n \in (0, 1]$  and the sampling rate  $\delta = d/n \in (0, 1]$ . At each  $\delta$ , the percentages of successes that an  $\ell_1$ -min algorithm finds the ground-truth solution  $\mathbf{x}_0$  (with a very small tolerance threshold) are measured over different  $\rho$ 's. Then a fixed success rate, say of 95%, over all  $\delta$ 's can be interpolated as a curve in the  $\rho$ - $\delta$  plot. In general, the higher the success rates, the better an algorithm can recover dense signals in the ( $P_1$ ) problem.

Figure 2 shows the 95% success-rate curves for the six algorithms. In the simulation, the ambient dimension  $d = 1000$  is fixed. To generate the ground-truth signal, a random subset of nonzero coefficients is chosen, and their values are drawn from the standard normal distribution and normalized to have

<sup>12</sup>Note that the MATLAB Toolbox for YALL1 is copyright protected, which by default uses (60) as the stopping criterion. Although we were not able to modify the source code, we have tuned the stopping parameter for the rest of the algorithms that roughly align with the average accuracy of YALL1, and then compare the computational complexity after this normalization step.

unit length. The projection matrix  $A$  is a Gaussian dictionary, whose coefficients are randomly generated from the standard normal distribution. We sample the average success rates on a grid of  $(\rho, \delta)$  pairs for each of the  $\ell_1$ -min algorithms, and the coordinates of the 95% rate are interpolated from the grid values.

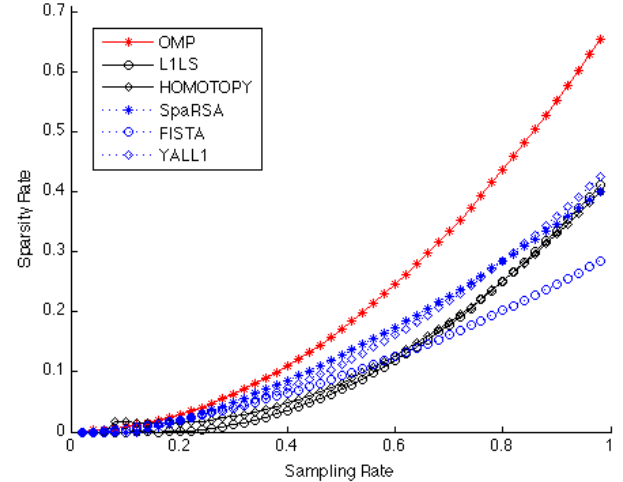


Fig. 2. The  $\rho$ - $\delta$  plot (in color) that shows the 95% success-rate curves for the six fast  $\ell_1$ -min algorithms.

The observations of the experiment are summarized below:

- 1) On average, OMP achieves the highest success rates. It shows OMP can be an excellent choice in the ideal scenario where the data noise is low. However, OMP is also one of the slowest algorithms, as we will show in the next experiment.
- 2) The success rates of SpaRSA and YALL1 are comparable over all sampling rates, and they also outperform the rest of the algorithms.
- 3) The success rates of L1LS and Homotopy are comparable over all sampling rates. In particular, their accuracy approaches that of SpaRSA and YALL1 in the high sampling-rate regime. In the low sampling-rate regime, Homotopy is slightly better than L1LS.
- 4) The success rates of FISTA are only higher than L1LS and Homotopy in the low sampling-rate regime. They are the lowest among the six algorithms with the increase of the sampling rate and the sparsity rate.

#### B. Performance with Moderate Data Noise

We are more interested in comparing the  $\ell_1$ -min algorithms when the measurement contains moderate amounts of data noise. In the second experiment, we rank the six algorithms under two scenarios: First, we measure the performance in the low-sparsity regime, where the ambient dimension  $n = 2000$  and the sparsity rate  $\rho = k/n = 0.1$  are fixed, and the dimension of the Gaussian random projection varies  $d = 300 - 1900$ . Second, we measure the performance when  $\mathbf{x}$  becomes dense w.r.t. a fixed sampling rate, where  $n = 2000$  and  $d = 1500$  are fixed, and the sparsity ratio  $\rho = k/n$  varies from 0.1 to 0.5. The results are shown in Figure 3 and 4. In both experiments,

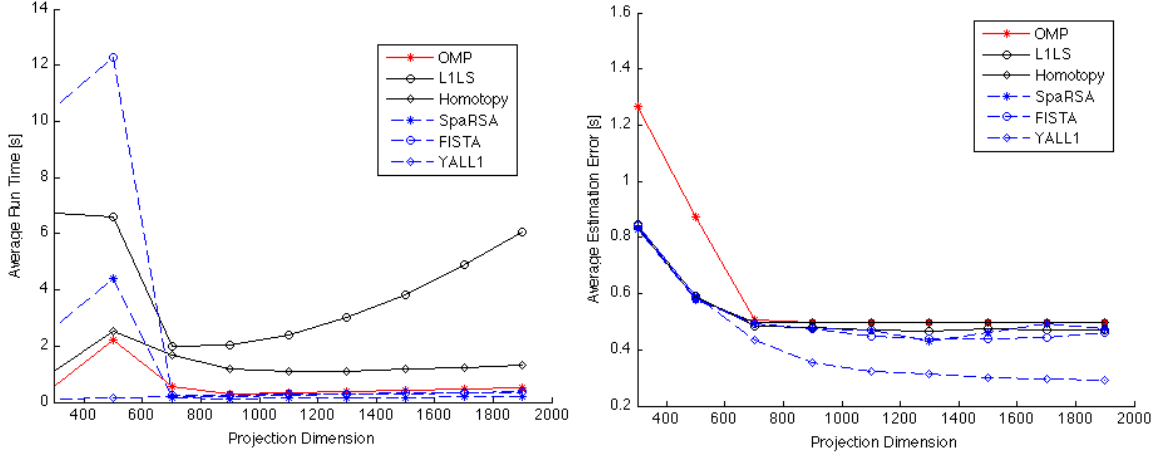


Fig. 3. Comparison of the six fast  $\ell_1$ -min algorithms w.r.t. a fixed sparsity ratio ( $n = 2000, k = 200$ ), and varying projection dimensions  $d = 300 - 1900$ . **Left:** Average run time. **Right:** Average  $\ell_2$ -norm error.

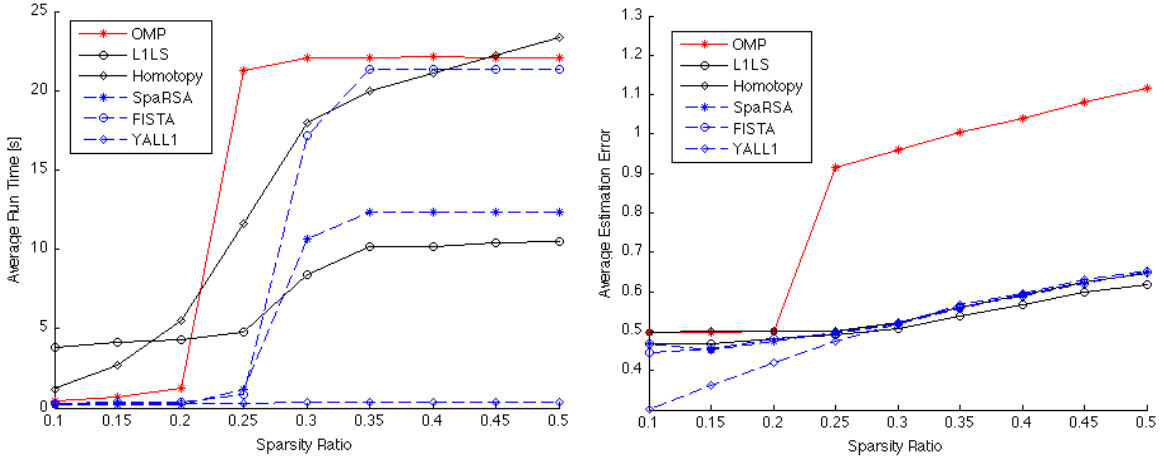


Fig. 4. Comparison of the six fast  $\ell_1$ -min algorithms w.r.t. a fixed sampling ratio ( $n = 2000, d = 1500$ ), and varying sparsity ratio  $k/n = 0.1 - 0.5$ . **Left:** Average run time. **Right:** Average  $\ell_2$ -norm error.

we corrupt the measurement vector  $\mathbf{b}$  with  $\mathbf{e}$ , an additive white noise term whose entries are i.i.d. distributed as  $N(0, 0.01)$ .

From the results, we draw the following observations. First, when a low sparsity ratio of  $\rho = 0.1$  is fixed in Figure 3,  $\ell_1$ -min becomes better conditioned as the projection dimension increase, and all algorithms converge to good approximate solutions when  $d > 750$  as indicated in Figure 3 Right. We then compare the speed of the six algorithm in Figure 3 Left:

- 1) When the projection dimension is small (e.g.,  $d < 750$ ) as  $\ell_1$ -min fails to converge to the global optimum, FISTA, L1LS and SpaRSA take a much longer time to exit than Homotopy, OMP, and YALL1.
- 2) When  $d > 750$ , the average run time of L1LS grows superlinearly with the projection dimension, while the run time of the rest algorithms largely remains constant.
- 3) The average run time of YALL1 is one of the lowest over all projection dimensions, which makes it the best algorithm in this comparison.

Second, when the projection dimension  $d = 1500$  is fixed in Figure 4, we compare both the average run time and the average estimation error when the sparsity varies:

- 1) The average estimation error of OMP quickly blows up when the sparsity ratio increases in Figure 4 Right. It shows that OMP is not stable when the data are noisy.
- 2) In the high-sparsity regime, The average run time of OMP, Homotopy, and FISTA is significantly higher than the other algorithms. It shows that the algorithms are not as effective when the signal becomes dense.
- 3) In the low-sparsity regime, L1LS is the slowest algorithm. However, its computational cost only increases modestly in the high-sparsity regime, and outperforms other approximate algorithms such as SpaRSA and FISTA.
- 4) YALL1 is among the fastest in both the low-sparsity regime and high-sparsity regime, and its run time remains almost constant while the sparsity ratio increases. This makes YALL1 the best algorithm in this comparison.

#### IV. EXPERIMENT: ROBUST FACE RECOGNITION

In this section, we benchmark the performance of the six algorithms in robust face recognition. The experiment is set

up to estimate sparse representation of real face images based on a so-called *cross-and-bouquet* (CAB) model [42].

More specifically, It has been known in face recognition that a well-aligned frontal face image  $\mathbf{b}$  under different lighting and expression lies close to a special low-dimensional linear subspace spanned by the training samples from the same subject, called a face subspace [7], [4]:

$$A_i = [\mathbf{v}_{i,1}, \mathbf{v}_{i,2}, \dots, \mathbf{v}_{i,n_i}] \in \mathbb{R}^{d \times n_i}, \quad (61)$$

where  $\mathbf{v}_{i,j}$  represents the  $j$ -th training image from the  $i$ -th subject stacked in the vector form. Given  $C$  subjects and a new test image  $\mathbf{b}$  (also in the vector form), we seek the sparsest linear representation of the sample with respect to all training examples:

$$\mathbf{b} = [A_1, A_2, \dots, A_C][\mathbf{x}_1; \mathbf{x}_2; \dots; \mathbf{x}_C] = A\mathbf{x}, \quad (62)$$

where  $A \in \mathbb{R}^{d \times n}$  collects all the training images.

Clearly, if  $\mathbf{b}$  is a valid test image, it must be associated with one of the  $C$  subjects. Therefore, the corresponding representation in (62) has a sparse representation  $\mathbf{x} = [\dots; \mathbf{0}; \mathbf{x}_i; \mathbf{0}; \dots]$ : on average only a fraction of  $\frac{1}{C}$  coefficients are nonzero, and the dominant nonzero coefficients in sparse representation  $\mathbf{x}$  reveal the true subject class.

In addition, we consider the situation where the query image  $\mathbf{b}$  may be severely occluded or corrupted. The problem is modeled by a corrupted set of linear equations  $\mathbf{b} = A\mathbf{x} + \mathbf{e}$ , where  $\mathbf{e} \in \mathbb{R}^d$  is an unknown vector whose nonzero entries correspond to the corrupted pixels. In [43], the authors proposed to estimate  $\mathbf{w} \doteq [\mathbf{x}; \mathbf{e}]$  together as the sparsest solution to the extended equation:

$$\min \|\mathbf{w}\|_1 \text{ subject to } \mathbf{b} = [A, I]\mathbf{w}. \quad (63)$$

The new dictionary  $[A, I]$  was dubbed a cross-and-bouquet model in the following sense. The columns of  $A$  are highly correlated, as the convex hull spanned by all face images of all subjects occupies an extremely tiny portion of the ambient space. These vectors are tightly bundled together as a “bouquet,” whereas the vectors associated with the identity matrix and its negative  $\pm I$  form a standard “cross” in  $\mathbb{R}^d$ , as shown in Figure 5. Finally, a quite surprising result was shown in [42] that accurate recover of sparse signals  $\mathbf{x}$  is possible and computationally feasible even when the nonzero corruption in  $\mathbf{e}$  approaches 100%.

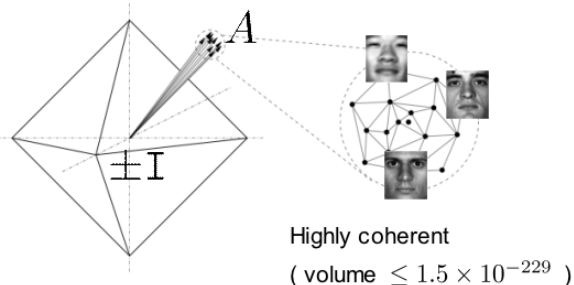


Fig. 5. The cross-and-bouquet model for face recognition. The raw images of human faces expressed as columns of  $A$  are clustered with very small variance. (Courtesy of John Wright [42])

The performance of the six  $\ell_1$ -min algorithms using the CAB model is benchmarked on the CMU Multi-PIE face database [23]. A subset of 249 subjects from the data set (Session 1) are used, each of which is captured in 20 frontal images under a fixed set of illumination settings. The images are then manually aligned and cropped, and down-sampled to  $40 \times 30$  pixels. Out of the 20 images for each subject, images  $\{0, 1, 7, 13, 14, 16, 18\}$  with extreme illumination conditions are chosen as the training images, and the rest 13 images are designated for testing. Finally, a certain number of image pixels are randomly corrupted with the corruption percentage from 0% to 80%, as shown in Figure 6.



Fig. 6. An aligned face image of Subject 1 in Multi-PIE, Session 1, under the ambient lighting condition (No. 0) is shown on the left. On the right, 20%, 40%, 60% and 80% of image pixels are randomly selected and corrupted with values ranges in  $[0, 255]$ , respectively.

We measure the performance of the algorithms in terms of the final recognition accuracy and the speed. In choosing a proper stopping criterion, the stopping threshold is individually tuned for each algorithm to achieve the highest recognition rate. Our priority is to achieve the highest accuracy in face recognition (e.g., must exceed 99% accuracy for any real-world scenarios), and the computational cost is only a secondary metric. The results are shown in Tables I and II.

TABLE I  
AVERAGE RECOGNITION ACCURACY (IN PERCENTAGE) ON THE MULTI-PIE DATABASE.

Corruption	0%	20%	40%	60%	80%
OMP	99.91	90.41	43.42	11.16	2.10
L1LS	100	100	100	88.49	15.93
Homotopy	100	100	99.55	92.82	34.09
SpaRSA	99.79	99.82	98.8	63.63	9.42
FISTA	100	92.95	61.98	21.46	20.52
YALL1	99.82	90.94	62.03	19.37	2.32

TABLE II  
AVERAGE RUN TIME (IN SECOND) ON THE MULTI-PIE DATABASE.

Corruption	0%	20%	40%	60%	80%
OMP	1.99	1.95	1.66	3.33	17.65
L1LS	28.75	27.77	25.08	14.55	8.13
Homotopy	3.56	8.40	15.44	28.23	23.81
SpaRSA	85.89	87.00	100.49	100.05	98.32
FISTA	151.3	17.06	16.47	18.20	3.48
YALL1	2.47	2.62	2.84	2.81	3.00

In Table I, clearly the Homotopy method achieves the best overall performance in recognition accuracy. For instance, with 60% of the pixels randomly occluded, its recognition rate based on the CAB model is about 93%. The worst performer is OMP, which corroborates that OMP does not perform well in the presence of practical data noise.

Among the three approximate  $\ell_1$ -min solutions, namely, SpaRSA, FISTA, and YALL1, the best performer is SpaRSA.

For instance, SpaRSA achieves 99% accuracy with 40% pixels randomly corrupted, while the accuracy of the rest two algorithms is only around 60%. However, note that the improved accuracy of SpaRSA carries a heavy penalty that the speed of SpaRSA is much slower than L1LS and Homotopy. For YALL1, although its overall speed shown in Table II is the fastest, its accuracy quickly drops with increase in the number of corrupted pixels.

Finally, it is more interesting to compare the difference in accuracy between L1LS and Homotopy, which provably solve the  $(P_1)$  problem, and SpaRSA, FISTA, and YALL1, which essentially rely on the soft thresholding function and approximation of the gradients of the objective function. In robust face recognition, we observe that the exact solutions as a whole significantly outperform the approximate solutions.

## V. CONCLUSION AND DISCUSSION

The paper has provided a comprehensive review of the five state-of-the-art fast  $\ell_1$ -min methods, i.e., *gradient projection*, *homotopy*, *soft shrinkage-thresholding*, *proximal gradient*, and *alternating direction*. The extensive experiment has shown that, under a wide range of data conditions, there is no clear winner that always achieves the best performance. For perfect, noise-free data, on average OMP is more effective than the rest of the algorithms, albeit at a much lower speed. Under random Gaussian dictionaries, approximate  $\ell_1$ -min solutions (i.e., SpaRSA, FISTA, and YALL1) are efficient to estimate sparse signals in both low-sparsity and high-sparsity regimes. In the application of robust face recognition, a special CAB model was constructed based on real training images representing a large set of human subjects. Homotopy and L1LS in turn achieve the highest recognition rate, and their computational cost is comparable to that of the other fast  $\ell_1$ -min algorithms. To aid peer evaluation, all the experimental scripts and data have been made available on our website.

## ACKNOWLEDGEMENTS

The authors would like to thank Dr. John Wright at Microsoft Research Asia and Zihan Zhou at the University of Illinois for their valuable comments. Yang also appreciates the hospitality of the Visual Computing Group at Microsoft Research Asia during his visit there in 2009.

## REFERENCES

- [1] M. Asif. Primal dual pursuit: A homotopy based algorithm for the dantzig selector. M. S. Thesis, Georgia Institute of Technology, 2008.
- [2] D. Baron, M. Wakin, M. Duarte, S. Sarvotham, and R. Baraniuk. Distributed compressed sensing. *preprint*, 2005.
- [3] J. Barzilai and J. Borwein. Two point step size gradient methods. *IMA Journal of Numerical Analysis*, 8:141–148, 1988.
- [4] R. Basri and D. Jacobs. Lambertian reflectance and linear subspaces. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 25(2):218–233, 2003.
- [5] A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM Journal on Imaging Sciences*, 2(1):183–202, 2009.
- [6] S. Becker, J. Bobin, and E. Candes. NESTA: a fast and accurate first-order method for sparse recovery. *preprint*, 2009.
- [7] P. Belhumeur, J. Hespanda, and D. Kriegman. Eigenfaces vs. Fisherfaces: recognition using class specific linear projection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 19(7):711–720, 1997.
- [8] A. Bruckstein, D. Donoho, and M. Elad. From sparse solutions of systems of equations to sparse modeling of signals and images. (*in press*) *SIAM Review*, 2007.
- [9] E. Candès. Compressive sampling. In *Proceedings of the International Congress of Mathematicians*, 2006.
- [10] E. Candès, J. Romberg, and T. Tao. Stable signal recovery from incomplete and inaccurate measurements. *Communications on Pure and Applied Math*, 59(8):1207–1223, 2006.
- [11] S. Chen, D. Donoho, and M. Saunders. Atomic decomposition by basis pursuit. *SIAM Review*, 43(1):129–159, 2001.
- [12] P. Combettes and V. Wajs. Signal recovery by proximal forward-backward splitting. *SIAM Multiscale Modeling and Simulation*, 4:1168–1200, 2005.
- [13] I. Daubechies, M. Defrise, and C. Mol. An iterative thresholding algorithm for linear inverse problems with a sparsity constraint. *Communications on Pure and Applied Math*, 57:1413–1457, 2004.
- [14] G. Davis, S. Mallat, and M. Avellaneda. Adaptive greedy approximations. *Journal of Constructive Approximation*, 13:57–98, 1997.
- [15] D. Donoho. De-noising by soft-thresholding. *IEEE Transactions on Information Theory*, 41:613–627, 1995.
- [16] D. Donoho. For most large underdetermined systems of linear equations the minimal  $\ell^1$ -norm near solution approximates the sparsest solution. *preprint*, 2004.
- [17] D. Donoho and Y. Tsaig. Fast solution of  $\ell^1$ -norm minimization problems when the solution may be sparse. *preprint*, <http://www.stanford.edu/tsaig/research.html>, 2006.
- [18] B. Efron, T. Hastie, I. Johnstone, and R. Tibshirani. Least angle regression. *The Annals of Statistics*, 32(2):407–499, 2004.
- [19] M. Figueiredo, R. Nowak, and S. Wright. Gradient projection for sparse reconstruction: Application to compressed sensing and other inverse problems. *IEEE Journal of Selected Topics in Signal Processing*, 1(4):586–597, 2007.
- [20] K. Frisch. The logarithmic potential method of convex programming. Technical report, University Institute of Economics (Oslo, Norway), 1955.
- [21] D. Gabay and B. Mercier. A dual algorithm for the solution of nonlinear variational problems via finite-element approximations. *Computers and Mathematics with Applications*, 2:17–40, 1976.
- [22] R. Glowinski and A. Marrocco. Sur l’approximation par éléments finis d’ordre un, et la résolution, par pénalisation-dualité d’une classe de problèmes de dirichlet nonlinéaires. *Revue Française d’Automatique, Informatique, Recherche Opérationnelle*, 9(2):41–76, 1975.
- [23] R. Gross, I. Mathews, J. Cohn, T. Kanade, and S. Baker. Multiple. In *IEEE International Conference on Automatic Face and Gesture Recognition*, 2006.
- [24] E. Hale, W. Yin, and Y. Zhang. A fixed-point continuation method for  $\ell^1$ -regularized minimization with applications to compressed sensing. Technical Report CAAM Technical Report TR07-07, Rice University, Houston, TX, 2007.
- [25] D. Hertog, C. Roos, and T. Terlaky. On the classical logarithmic barrier function method for a class of smooth convex programming problems. *Journal of Optimization Theory and Applications*, 73(1):1–25, 1992.
- [26] J. Hiriart-Urruty and C. Lemaréchal. *Convex analysis and minimization algorithms*. Springer-Verlag, 1996.
- [27] N. Karmarkar. A new polynomial time algorithm for linear programming. *Combinatorica*, 4:373–395, 1984.
- [28] C. Kelley. *Iterative methods for linear and nonlinear equations*. SIAM, Philadelphia, 1995.
- [29] S. Kim, K. Koh, M. Lustig, S. Boyd, and D. Gorinevsky. An interior-point method for large-scale  $\ell_1$ -regularized least squares. *IEEE Journal of Selected Topics in Signal Processing*, 1(4):606–617, 2007.
- [30] M. Kojima, N. Megiddo, and S. Mizuno. Theoretical convergence of large-step primal-dual interior point algorithms for linear programming. *Mathematical Programming*, 59:1–21, 1993.
- [31] D. Malioutov, M. Cetin, and A. Willsky. Homotopy continuation for sparse signal representation. In *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing*, 2005.
- [32] N. Megiddo. Pathways to the optimal set in linear programming. In *Progress in Mathematical Programming: Interior-Point and Related Methods*, pages 131–158, 1989.
- [33] R. Monteiro and I. Adler. Interior path following primal-dual algorithms. Part I: Linear programming. *Mathematical Programming*, 44:27–41, 1989.
- [34] Y. Nesterov. A method of solving a convex programming problem with convergence rate  $O(1/k^2)$ . *Soviet Mathematics Doklady*, 27(2):372–376, 1983.
- [35] Y. Nesterov. Gradient methods for minimizing composite objective function. *ECORE Discussion Paper*, 2007.
- [36] J. Nocedal and S. Wright. *Numerical Optimization*. Springer, New York, 2nd edition, 2006.

- [37] M. Osborne, B. Presnell, and B. Turlach. A new approach to variable selection in least squares problems. *IMA Journal of Numerical Analysis*, 20:389–404, 2000.
- [38] M. Plumbley. Recovery of sparse representations by polytope faces pursuit. In *Proceedings of International Conference on Independent Component Analysis and Blind Source Separation*, pages 206–213, 2006.
- [39] T. Serafini, G. Zanghirati, and L. Zanni. Gradient projection methods for large quadratic programs and applications in training support vector machines. *Optimization Methods and Software*, 20(2–3):353–378, 2004.
- [40] R. Tibshirani. Regression shrinkage and selection via the LASSO. *Journal of the Royal Statistical Society B*, 58(1):267–288, 1996.
- [41] J. Tropp and A. Gilbert. Signal recovery from partial information by orthogonal matching pursuit. *Preprint*, 2005.
- [42] J. Wright and Y. Ma. Dense error correction via  $\ell^1$ -minimization. (accepted) *IEEE Transactions on Information Theory*, 2010.
- [43] J. Wright, A. Yang, A. Ganesh, S. Sastry, and Y. Ma. Robust face recognition via sparse representation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 31(2):210 – 227, 2009.
- [44] S. Wright, R. Nowak, and M. Figueiredo. Sparse reconstruction by separable approximation. In *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2008.
- [45] J. Yang and Y. Zhang. Alternating direction algorithms for  $\ell_1$ -problems in compressive sensing. (preprint) *arXiv:0912.1185*, 2009.