

Minimizing Nonconvex Functions for Sparse Vector Reconstruction

Nasser Mourad and James P. Reilly

Department of Electrical and Computer Engineering
McMaster University, Hamilton, Ont. CANADA L8S 4K1

Abstract—In this paper we develop a novel methodology for minimizing a class of non-convex (concave on the non-negative orthant) functions for solving under-determined linear system of equations $As = x$ when the solution vector s is known a priori to be sparse. The proposed technique is based on locally replacing the original objective function by a quadratic convex function which is easily minimized. The resulting algorithm is iterative and is absolutely converging to a fixed point of the original objective function. For a certain selection of convex objective functions, the class of algorithms called Iterative Reweighted Least Squares (IRLS) are shown to be a special case of the proposed methodology. Thus the proposed algorithms are a generalization and unification of the previous methods. In addition, we also propose a new class of algorithms with better convergence properties compared to the regular IRLS algorithms and hence can be considered as enhancements to these algorithms. Since the original objective functions are non-convex, the proposed algorithm is susceptible to convergence to a local minimum. To alleviate this difficulty, we propose a random perturbation technique that enhances the performance of the proposed algorithm. The numerical results show that the proposed algorithms outperform some of the well known algorithms that are usually utilized for solving the same problem.

I. INTRODUCTION

There are many situations in science and technology that seek solutions to under-determined systems of equations, i.e. systems of linear equations with fewer equations than unknowns. Such problems are extremely common for a variety of reasons. For instance, the number of sensors may be small due to physical limitations as in breast cancer imaging, or the sensing process may be slow so that one can only measure the object a few times, as in MRI. Many other examples in inverse problems, array signal processing, and biomagnetic imaging all come to mind. This problem can be expressed mathematically as finding a unique solution to a linear under-determined systems of equations $As = x$, where $x \in \mathbb{R}^m$ is the vector of measurements, $A \in \mathbb{R}^{m \times n}$ is the sensing (measuring) matrix, and $m < n$. However, it is known that a system of linear equations with fewer equations than unknowns has infinitely many solutions, and thus it is necessary to impose constraints on the candidate solution to identify which of these candidate solutions is the desired one.

A powerful constraint that can be used in this regard is the “sparsity” of the solution vector. Sparsity means that the solution vector has few nonzero elements. The problem can be readily extended to the case of non-sparse solution vectors that are piecewise constant (hence its gradient is sparse [1]), or can be sparsely represented in some basis (e.g., Fourier basis and wavelet related basis [2]). The problem of reconstruction of sparse signals from a limited number of linear measurements is known in the literature as “compressed sensing” or “compressive sampling” (CS). Compressed sensing has a wide range of applications and it has been studied in the literature under many different names such as subset selection [3], sparse coding [4], sparse component analysis (SCA) with application to blind source separation of more sources than sensors [5], [6]. Other applications are biomagnetic inverse problems [7], [8], dictionary learning [9], [10], and image restoration [1], [11]. See [12] and the references therein for a list of more applications. In all these applications, the underlying linear inverse problem is the same and can be stated as follows: Represent a signal of interest using the minimum number of “atoms” (vectors) from an overcomplete dictionary.

Restricting the solution vector to be sparse converts the underlying problem from being impossible into being a tractable, but nevertheless still difficult, problem. There are many algorithms that have been developed in the literature for solving the under-determined linear system of equations under the constraint that the candidate solution vector is sparse. Examples include greedy algorithms [13], [14], the basis pursuit [15], [16], iterative-thresholding algorithms [17], and iterative reweighted norm algorithms [1], [7], [12], [18]–[20]. A brief description of some of these approaches is presented in Section II.

In this paper a novel methodology is developed to minimize a class of non-convex (concave on the non-negative orthant) functions for solving the afore mentioned problem. The proposed technique is based on locally replacing the original objective function by a quadratic convex function which is easily minimized. The resulting algorithm is iterative, and it will be shown that the penalty imposed by the convex objective function at any iteration depends on the gradient of the original objective function evaluated at the previous solution vector. Accordingly, if the objective function is carefully chosen, the penalty imposed by the convex function on the small entries of the solution vector can increase as the algorithm gets closer to a local minimum of the original objective function. As will be shown in this paper, for a certain selection of the convex

objective function, the class of algorithms called Iterative Re-weighted Least Squares (IRLS) [1], [7], [18]–[20] can be derived from the proposed methodology. Thus the proposed algorithms are a generalization and unification of the previous methods.

In this paper we propose a straightforward technique for selecting a convex function such that, for any starting solution vector s_0 , the algorithm generates a sequence $\{s_k\}_{k=1}^{\infty}$ that converges to a fixed point of the original objective function. Since the original objective functions are non-convex, the proposed algorithm is susceptible to convergence to a local minimum. To alleviate this difficulty, we propose a random perturbation technique that enhances the performance of the proposed algorithm. The proposed algorithm is called MCCR, as an abbreviation of **Minimizing a Concave function via a Convex function Replacement**.

Our contributions in this paper can be summarized as follows. (1) we propose a technique that provides a deeper understanding of the IRLS class of algorithms and provides a natural mechanism for deriving IRLS algorithms starting from suitably chosen *diversity* (antisparsity) measures, where the diversity of a vector is defined as the number of its nonzero entries, (2) we propose a new class of algorithms with better convergence properties compared to the regular IRLS algorithms and, hence, can be considered as enhancements to these algorithms, and (3) we suggest some novel techniques for improving the performance of IRLS methods.

The paper is organized as follows. A revision of some of the previous approaches is presented in Section II. The proposed convex objective function is presented in Section III. The MCCR algorithm is derived in Section IV. Some techniques that enhance the performance of the derived algorithms are presented in Section V. Section VI presents some computer simulations for assessing the performance of the proposed algorithms. Finally, conclusion remarks are given in Section VII.

II. REVIEW OF PREVIOUS APPROACHES

The problem of finding a *sparse* solution to an under-determined linear system of equations can be reformulated as one of basis selection from an over-complete dictionary matrix. Let \mathbf{A} denote an $(m \times n)$ dictionary matrix whose columns comprise an over-complete set of basis vectors, i.e. $m < n$, where it is assumed that $\text{rank}(\mathbf{A}) = m$. It is required to represent a vector $\mathbf{x} \in \mathbb{R}^m$ by the *smallest* possible number of columns of the dictionary matrix, that is

$$\mathbf{x} = \mathbf{A}\mathbf{s}^*, \quad (1)$$

where $\mathbf{s}^* \in \mathbb{R}^n$ is an l -sparse vector, i.e., a vector with $l \ll n$ nonzero entries. It is also assumed that $l < m$.

The goal of a sparse-signal recovery algorithm is to obtain an estimate of \mathbf{s}^* given only \mathbf{x} and \mathbf{A} . This problem is non-trivial since \mathbf{A} is overcomplete. Since the original vector \mathbf{s}^* is sparse, the problem of finding the desired solution vector can be phrased as an optimization problem where the objective is to maximize (minimize) an appropriate measure of sparsity (diversity) of the solution vector while simultaneously

satisfying the constraints defined by (1). This can be expressed mathematically as

$$\hat{\mathbf{s}} = \arg \min_{\mathbf{s}} g(\mathbf{s}) \quad \text{subject to} \quad \mathbf{x} = \mathbf{A}\mathbf{s}, \quad (2)$$

where $g(\cdot)$ is an objective function to be minimized that encourages sparsity in the solution vector. In this section we consider this function to be of the form

$$g_p(\mathbf{s}) = \|\mathbf{s}\|_p^p = \sum_i |s[i]|^p, \quad (3)$$

where $p \geq 0$, and $s[i]$ is the i -th element of \mathbf{s} . Eq. (3) expresses the p -th norm of \mathbf{s} (although it is not strictly a valid norm for $0 \leq p < 1$). We now briefly discuss issues relating to solving (3) for various values of $0 \leq p \leq 1$.

The ℓ_0 -norm counts the number of non-zero entries in \mathbf{s} . Therefore, the ℓ_0 -norm of an l -sparse solution vector \mathbf{s}^* is expressed as $\|\mathbf{s}^*\|_{\ell_0} = l$. An l -sparse solution vector can be obtained by searching over the $\binom{n}{l}$ possible ways in which the basis sets can be chosen to find the best solution. Unfortunately, the cost of such combinatorial searches is prohibitive, requiring investigation of other approaches.

For $p = 1$, (2) is usually called basis pursuit (BP) [15]. Since $p = 1$ is the smallest value of p for which $g_p(\mathbf{s})$ is convex, ℓ_1 -minimization has been utilized in the context of sparse solutions for many years. See [12], [20] and the references therein for the history of ℓ_1 -minimization and its applications.

It was shown in [24]–[26] that minimizing ℓ_q -norm, for values of $0 < q < 1$ performs better than minimizing the ℓ_1 -norm in the sense that a smaller number of measurements are needed for exact reconstruction of the sparse solution vector.

Another class of algorithms that can be utilized for finding a sparse solution vector is called Iterative Re-weighted Least Squares (IRLS) [7], [18]–[20], [27], [28]. IRLS algorithms combine the simplicity of the least squares solution while imposing sparsity on the solution vector. IRLS has been successfully utilized for minimizing a wide class of objective functions, e.g., the ℓ_p -norm with $0 < p < 2$, the logarithmic function of the form $\sum_i \log(|s[i]|)$ [7], and the Shannon entropy [18]. See [7], [20], [29], [30] for the history of IRLS and their convergence properties. IRLS algorithms have the form

$$(P_{w\ell_2}) \quad \min_{\mathbf{s}} \|\mathbf{W}^{-1}\mathbf{s}\|_2^2 \quad \text{subject to} \quad \mathbf{x} = \mathbf{A}\mathbf{s}, \quad (4)$$

where \mathbf{W} is a diagonal weighting matrix that reflects our prior knowledge about the solution vector \mathbf{s} . The resulting algorithm is iterative, and the estimated solution at the k th iteration can be expressed as

$$\mathbf{s}_k = \mathbf{W}_k(\mathbf{A}\mathbf{W}_k)^\dagger \mathbf{x}, \quad (5)$$

where † is the Moore-Penrose inverse [31]. The difference between different IRLS algorithms resides in the way that the diagonal matrix is defined. In [7] \mathbf{W}_k was selected as $\mathbf{W}_k = \text{diag}(|\mathbf{s}_{k-1}|)$ which was shown in [18] to be equivalent to minimizing $g(\mathbf{s}) = \sum_i \log(|s[i]|)$, while it was selected as $\mathbf{W}_k = \text{diag}(|\mathbf{s}_{k-1}[i]|^{1-0.5q})$ in [18] for minimizing ℓ_q -norm¹ by following the affine scaling methodology. The local rate

¹unless explicitly stated, $0 < q < 1$.

of convergence varies according to the expression of \mathbf{W} ; for instance it was shown to be quadratic in [7], while in [20] it was either linear or super-linear for the algorithm approximating ℓ_1 -norm or ℓ_q -norm, respectively.

Another approach for weighted norm minimization is the one proposed in [12], where $g(s)$ in (2) is replaced by $g_{w\ell_1}(s) = \|\mathbf{W}_k^{-1}s\|_{\ell_1}$, where $\mathbf{W}_k = \text{diag}(|s_{k-1}|)$ is also a diagonal weighting matrix. It was shown in [12] that this algorithm performs much better than the ℓ_1 -norm minimization and converges in a few iterations. However each iteration is computationally expensive compared with an IRLS iteration.

III. PROPOSED CONVEX FUNCTION

In the absence of noise, a sparse solution vector of the under-determined system of equations, $\mathbf{A}s = \mathbf{x}$, can be obtained by solving the following optimization problem

$$\hat{s} = \arg \min_s g(s) \quad \text{subject to} \quad \mathbf{x} = \mathbf{A}s, \quad (6)$$

where $\mathbf{A} \in \mathbb{R}^{m \times n}$, $m < n$, and $g(\cdot)$ is an objective function to be minimized that encourages sparsity in the solution. In this paper we restrict our attention to the class of non-convex objective functions that have the following two properties [9]:

- (P1) $g(s)$ is separable, i.e., $g(s) = \sum_i g_c(s[i])$, where $g_c(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$ is a scalar function, and $s[i]$ is the i th entry of s . Therefore, $g(s)$ is also permutation invariant, i.e., $g(s) = g(\mathbf{P}s)$ for any permutation matrix \mathbf{P} .
- (P2) The function $g_c(s) : \mathbb{R} \rightarrow \mathbb{R}$ referred to in P1 is sign invariant and concave-and-monotonically increasing on the nonnegative orthant \mathcal{O}_1 .

Examples of $g_c(s)$ that were extensively used in the literature are $g_q(s) = |s|^q$ and $g_{\log}(s) = \log(|s|)$ [18], [19], [24], [25]. More objective functions that obey property P2 are suggested in [32].

Since $g(s)$ is sign invariant and concave on \mathcal{O}_1 , it is concave on each of the other orthants \mathcal{O}_k , $1 \leq k \leq 2^n$. However, this does not imply that $g(s)$ is concave on \mathbb{R}^n .

Consider the following theorem.

Theorem 1 ([9] Theorem 8): Let $g(s) : \mathbb{R}^n \rightarrow \mathbb{R}$ be permutation invariant, sign invariant, and concave on the nonnegative orthant \mathcal{O}_1 . Then the global minimum of the optimization problem

$$\min_s g(s) \quad \text{subject to} \quad \mathbf{x} = \mathbf{A}s.$$

has at most m nonzero entries, where $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $m < n$.

Proof: See [9]. \square

Theorem 1 implies that if the purpose of solving (6) is to achieve a sparse solution vector, then the permutation invariant and sign invariant concave functions can be used as diversity measure functions.

The objective of this paper is to iteratively replace $g(s)$ by a convex smooth function $f(s)$, which can be easily minimized in an iterative fashion. The advantages of the proposed approach is that, by replacing the nonconvex function $g(s)$ by $f(s)$, each iteration becomes very tractable. However, since the original function $g(s)$ is nonconvex, the proposed algorithm may still converge to a local minimum.

For the proposed approach to be useful, the function $f(s)$ must be selected such that its global minimizer is guaranteed to reduce the original function $g(s)$. Towards that end, and to simplify the exposition, we consider replacing the one dimensional function $g_c(s)$ by a convex function $f(s)$ first. We then generalize to the n dimensional case. Consider the following theorem

Theorem 2: Given a differentiable one dimensional function $g_c(s) : \mathbb{R} \rightarrow \mathbb{R}$ that obeys Property P2, a one dimensional convex function $f(s) : \mathbb{R} \rightarrow \mathbb{R}$ with a global minimizer s^o , and an initial point s_0 , then sufficient conditions for $g_c(s)$ to be reduced at s^o , i.e., $g_c(s^o) \leq g_c(s_0)$, are: (1) $\text{sign}(f'(s_0)) = \text{sign}(g'_c(s_0))$, and (2) $|s^o| \leq |s_0|$, where $g'(s_0)$ is the gradient of $g(s)$ at $s = s_0$.

Proof: Since the negative of the gradient of any objective function provides a valid direction for minimizing this objective function, the first condition in Theorem 2 insures that the minimizing direction of $f(s)$ is a valid minimizing direction for the original function $g(s)$. On the other hand, from P2, it is known that $g(s)$ is sign invariant and concave-and-monotonically increasing on the nonnegative orthant. Accordingly, the second condition $|s^o| \leq |s_0|$ in Theorem 2 insures that that $g(s^o) \leq g(s_0)$. \square

Remark 1: For the special case where we select $f(s)$ such that $f(s) \geq g_c(s)$ and $f(s_0) = g_c(s_0)$, the convex objective function $f(s)$ majorizes $g(s)$ and the proposed technique is reduced to the Majorization-Minimization (MM) method [33]. However, in this paper we restrict $f(s)$ to satisfy the two conditions of Theorem 2 rather than the two conditions of the MM method. As will be shown in Section IV, the convex objective function proposed in Section IV-A satisfies the two conditions of the MM method, while the function proposed in Section IV-B violates these conditions.

There are many ways of selecting convex objective functions that satisfy the two conditions of Theorem 2. We propose a straightforward technique for doing this task. The proposed function has the following form

$$f(s) = g_c(s_0) + g'_c(s_0)(s - s_0) + \beta_0(s - s_0)^2, \quad (7)$$

where $g'_c(s_0)$ is the gradient of $g_c(s)$ at s_0 , and $\beta_0 \geq 0$. Clearly the first condition of Theorem 2 is always satisfied, since $f'(s_0) = g'_c(s_0)$, and the second condition can be met by properly selecting the value of β_0 .

Remark 2: If β_0 in (7) is selected as $\beta_0 = -0.5g''_c(s_0)$, where $g''_c(s_0)$ is the Hessian of $g_c(s)$ at $s = s_0$, then the function $f(s)$ in (7) is equivalent to the second-order Taylor series expansion of $g_c(s)$, with the sign of the quadratic term reversed. This results in $f(s)$ being a convex function. However it does not guarantee that the second condition of Theorem 2 is satisfied.

Assuming that the value of β_0 was chosen such that $f(s)$, defined in (7), satisfies the second condition of Theorem 2, then the point $s_1 = \arg \min_s f(s)$ is the global minimizer to $f(s)$ but not to $g(s)$. Accordingly, to reach the global minima of $g(s)$, the procedure is iterated, i.e., set $s_0 \leftarrow s_1$ and select β_0 such that (7) satisfies the conditions of Theorem 2. This

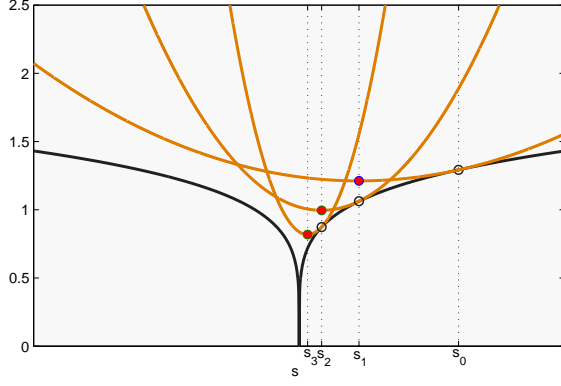


Fig. 1. The function $f(s)$ (orange curves) in successive iterations of the proposed algorithm for minimizing the nonconvex function $g_c(s)$ (black curve). The hollow circles represent the starting points, while the solid ones represent the global minima of the locally-convex functions.

procedure is repeated until $g(s)$ is minimized. The function $f(s)$ in successive iterations is shown in Fig. 1.

Since $g(s)$ is separable, the one dimensional function (7) can be readily extended to the n dimensional case by constructing the following separable convex function

$$\begin{aligned} f(s) &= \sum_{i=1}^n f(s[i]) \\ &= \sum_{i=1}^n g_c(s_0[i]) + g'_c(s_0[i])(s[i] - s_0[i]) + \beta_0[i](s[i] - s_0[i])^2 \\ &= g(s_0) + (s - s_0)^T d_0 + (s - s_0)^T B_0 (s - s_0), \end{aligned}$$

where $d_0 = \nabla g(s_0)$ is the gradient of $g(s)$ at $s = s_0$, and B_0 is a diagonal positive semidefinite matrix. Instead of using (8) in deriving the proposed algorithms, the following simplified function will be used

$$\tilde{f}(s) = s^T (d_0 - 2B_0 s_0) + s^T B_0 s, \quad (9)$$

which is equivalent to $f(s)$ after dropping the constant terms. Although the constant term in (8) does not affect the value of the solution vector of $f(s)$, it is important in proving the convergence of the derived algorithms.

IV. THE MCCR ALGORITHMS

In this section we derive algorithms for finding a sparse solution vector of the noise-free linear underdetermined system of equations, based on replacing the nonconvex objective function $g(s)$ in (6) by the convex objective function $f(s)$ defined in (8), or equivalently $\tilde{f}(s)$ defined in (9). Since $f(s)$ satisfies the first condition of Theorem 2, the remaining part of this section considers the problem of selecting B_0 such that, for a given starting point s_0 , the original objective function $g(s)$ satisfies $g(s_1) \leq g(s_0)$, where

$$\begin{aligned} s_1 &= \arg \min_s s^T (d_0 - 2B_0 s_0) + s^T B_0 s \\ &\text{subject to } x = As. \end{aligned} \quad (10)$$

The optimization problem (10) can be solved by following the standard method of Lagrangian multipliers, which results in the following general expression for s_1 [25]

$$s_1 = s_0 - 0.5B_0^{-1}d_0 + 0.5B_0^{-1}A^T(AB_0^{-1}A^T)^{-1}AB_0^{-1}d_0. \quad (11)$$

Eq. (11) represents the general expression of the proposed MCCR algorithm. In the next subsections we propose two different choices for B_0 that guarantee $g(s)$ is reduced at s_1 , i.e., $g(s_1) \leq g(s_0)$.

A. First approach: Select B_0 such that $\tilde{f}(s) = \|W^{-1}s\|_2^2$

In this subsection we select B_0 such that the global minima of $\tilde{f}(s)$, and hence $f(s)$, occurs at $s^o = 0$.² From (8), or equivalently (9), the global minima of $f(s)$ occurs at s^o , which can be calculated from the following equation

$$\nabla_s f(s)|_{s^o} = d_0 + 2B_0(s^o - s_0) = 0,$$

which leads to

$$s^o = s_0 - 0.5B_0^{-1}d_0. \quad (12)$$

Accordingly, setting $s^o = 0$ is equivalent to

$$s_0 = 0.5B_0^{-1}d_0. \quad (13)$$

Substituting (13) into (11), and utilizing the fact that s_0 is feasible, i.e., $As_0 = x$, the expression of s_1 is reduced to

$$s_1 = B_0^{-1}A^T(AB_0^{-1}A^T)^{-1}x. \quad (14)$$

Note that s_1 is the global minima of the convex problem (10) but not of the original problem (6). However, since $g(s_1) \leq g(s_0)$ from Theorem 3, a stable point of (6) can be obtained by iteratively repeating (14), i.e., for $k \geq 0$

$$s_{k+1} = B_k^{-1}A^T(AB_k^{-1}A^T)^{-1}x, \quad (15)$$

where B_k is a diagonal matrix which depends on s_k and $d_k = \nabla g(s_k)$, the gradient of $g(s)$ at s_k . The diagonal entries of B_k can be easily calculated from (13) as follows

$$B_k[i, i] = \frac{d_k[i]}{2s_k[i]}, \quad i = 1, \dots, n. \quad (16)$$

Note that $B_k[i, i]$ calculated from (16) is nonnegative, i.e., B_k is positive semidefinite, which is a necessary condition for $f(s)$ to be convex. This follows readily from the sign-invariant property of $g(s)$ and its monotonically increasing property on the nonnegative orthant \mathcal{O}_1 . Note that, substituting (13) into (9), the expression of $\tilde{f}(s)$ is reduced to $\tilde{f}(s) = \|W^{-1}s\|_2^2$ as desired, where $W = B_0^{-1/2}$.

Since the original objective function $g(s)$ is nonconvex, the convergence of the proposed algorithm (15) into the “global” minima of the original problem (6) is not guaranteed. However, the proposed algorithm always converges to a fixed point of (6) and does not have undesired properties such as divergence or oscillation. The proof of convergence of (15) to a fixed point of (6) is shown by the following theorem.

²It should be noted that, even though the global minimizer of $\tilde{f}(s)$ is $s^o = 0$, the solutions to the constrained problem (10) is given by (11) and is generally not $s^o = 0$.

Theorem 3: Given a starting solution vector $\mathbf{s}_k \in \mathbb{R}^n$ and a separable function $g(\mathbf{s}) = \sum_i g_c(s[i]) : \mathbb{R}^n \rightarrow \mathbb{R}$, where $g_c(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$ is differentiable and obeys P2. Let \mathbf{d}_k denote the gradient of $g(\mathbf{s})$ at \mathbf{s}_k . If \mathbf{B}_k in (8) is calculated using (16) and a new solution vector \mathbf{s}_{k+1} is obtained using (15), then $g(\mathbf{s}_{k+1}) \leq g(\mathbf{s}_k)$ with the equality holds only if \mathbf{s}_k is a fixed point of the original problem (6).

Proof: Since the global minima of $f(\mathbf{s})$ occurs at $\mathbf{s}^o = 0$, it is straightforward to show that $f(\mathbf{s})$ is sign invariant, i.e., $f(\mathbf{s}) = f(|\mathbf{s}|)$. Since both $g(\mathbf{s})$ and $f(\mathbf{s})$ are sign invariant, it will be assumed without loss of generality that $\mathbf{s}_k \in \mathcal{O}_1$, the nonnegative orthant. From (8) it is clear that both $f(\mathbf{s})$ and $g(\mathbf{s})$ share a common tangent at \mathbf{s}_k , which is given by the first two terms of (8). From the convexity of $f(\mathbf{s})$ and the concavity of $g(\mathbf{s})$ on \mathcal{O}_1 , a common tangent to $f(\mathbf{s})$ and $g(\mathbf{s})$ on \mathcal{O}_1 implies that $g(\mathbf{s})$ is upper bounded by $f(\mathbf{s})$ on \mathcal{O}_1 . This is also applicable to the other $2^n - 1$ orthants. Therefore, $g(\mathbf{s}) \leq f(\mathbf{s})$, where the equality holds only at the 2^n tangent points, i.e., for a given starting point $\mathbf{s}_k \in \mathcal{O}_1$ we have $g(\mathbf{s}_k) = f(\mathbf{s}_k)$. However, from (10), and since $f(\mathbf{s})$ is a quadratic convex function, we have $f(\mathbf{s}_{k+1}) \leq f(\mathbf{s}_k)$ where the equality holds if and only if $\mathbf{s}_{k+1} = \mathbf{s}_k$, i.e., \mathbf{s}_{k+1} can not be one of the other $2^n - 1$ tangent points. In conclusion, if \mathbf{s}_{k+1} is different from \mathbf{s}_k , then $g(\mathbf{s}_{k+1}) < f(\mathbf{s}_{k+1}) < f(\mathbf{s}_k) = g(\mathbf{s}_k)$, while $g(\mathbf{s}_{k+1}) = f(\mathbf{s}_{k+1}) = f(\mathbf{s}_k) = g(\mathbf{s}_k)$ if and only if $\mathbf{s}_{k+1} = \mathbf{s}_k$. In this last case, \mathbf{s}_k is a fixed point of (6). \square

Thus, starting from any feasible solution vector $\mathbf{s}_0 \in \mathbb{R}^n$, the algorithm (15) generates a sequence $\{\mathbf{s}_k\}_{k=1}^\infty$ that converges at least to a local minimum of the original problem (6). At this point, \mathbf{s} does not vary from one iteration to the next.

Theorem 3 shows that any algorithm of the form (15) with \mathbf{B}_k calculated using (16) is guaranteed to converge to a fixed point of the original optimization problem (6) as long as $g(\mathbf{s}) \in \mathcal{G}$, where \mathcal{G} is the set of all objective functions that obey P1–P2. The exact expression of the objective function $g(\mathbf{s})$ affects only the rate of convergence of (15). The rate of convergence of some functions $g(\mathbf{s}) \in \mathcal{G}$ are derived in [32].

Relation with IRLS algorithms: Define $\mathbf{W}_k = \mathbf{B}_k^{-\frac{1}{2}}$, then (15) can be written as

$$\mathbf{s}_{k+1} = \mathbf{W}_k(\mathbf{A}\mathbf{W}_k)^\dagger \mathbf{x}, \quad (17)$$

where † is the Moore-Penrose inverse [31]. Comparing (17) with (5) we readily find that the derived algorithm has the form of the IRLS algorithms. The convergence of an IRLS algorithm to a sparse solution vector depends on the relation between \mathbf{W}_k and the previous solution vector \mathbf{s}_k . Some IRLS algorithms may not converge to a sparse solution vector, e.g., the algorithm derived in [18] for minimizing the Shannon entropy. However, the proposed technique provides a general methodology for deriving IRLS algorithms that converge to sparse solution vectors and enjoy fixed point convergence.

The performance of the proposed algorithm in finding a sparse solution is shown in Fig. 2. The difference between Fig. 2 and Fig. 1 is that the global minimum of each convex function in Fig. 2 is at the origin. Note that this is the global minimum to the unconstrained minimization of $f(\mathbf{s})$ but not to the constrained problem (10), which has a global minima in

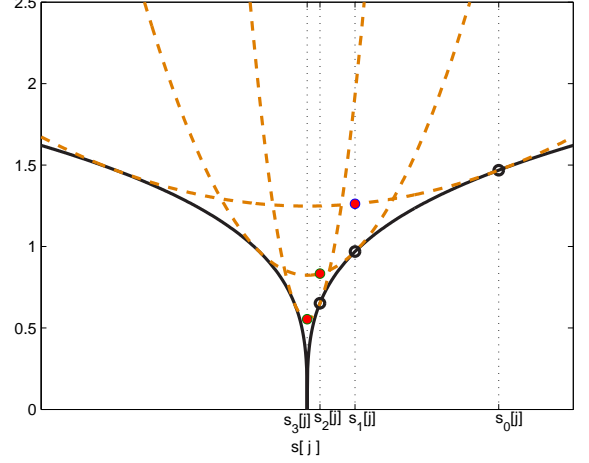


Fig. 2. Successive iterations of the proposed algorithm for minimizing the nonconvex function $g_c(s[j])$ (black curve). The hollow circles represent the starting points, while the solid ones represent the global minima of the locally-convex functions.

this case given by (15). Fig. 2 shows the successive iterations associated with reducing a zero entry of the solution vector \mathbf{s} towards its optimum value of zero. As shown in this figure, the algorithm takes large steps when the value of $s[j]$ is large, while the steps get smaller and smaller as $s_k[j]$ approaches zero. This performance is due to the weighting matrix \mathbf{B}_k whose i th diagonal element is given by (16). From (16) we can see that as $s_k[i] \rightarrow 0$ the weights become large due to the fact that the function $g_c(s)$ is concave and monotonically increasing on \mathcal{O}_1 , and hence its gradient increases as the value of s decreases. Large weights associated with small values of s are required for a sparse solution. Note that the quadratic replacement function $f(\mathbf{s})$ on the other hand imposes small penalties on small values of \mathbf{s} , a fact which explains why the MCCR approach can not achieve a sparse solution in one step.

B. Second Approach: Select \mathbf{B}_0 such that $\tilde{f}(\mathbf{s}) = \|\mathbf{W}^{-1}(\mathbf{s} - \theta\mathbf{s}_0)\|_2^2$

In Section IV-A we proved that selecting \mathbf{B}_k as in (16) reduces (8) to $f(\mathbf{s}) = \|\mathbf{W}_k^{-1}\mathbf{s}\|_2^2$, where $\mathbf{W}_k = \mathbf{B}_k^{-\frac{1}{2}}$. Consequently, the resulting algorithm (15), or equivalently (17), has the general form of the IRLS algorithms. Moreover, it was shown by Theorem 3 that, starting from any bounded solution vector, the proposed algorithm is guaranteed to converge to a stable solution vector of the original problem (6). However, the rate of convergence of the proposed algorithm is limited by the nature of the weighted ℓ_2 -ball, as shown in Fig. 3. In contrast to the ℓ_1 -ball, the ℓ_2 -ball, and hence ℓ_{w2} -ball, the weighted ℓ_2 -ball, have rounded tips. As a result, the proposed algorithm, and all other IRLS algorithms, may require large number of iterations to converge to a solution vector, especially when m and n are large.

In this figure, the exact and the estimated solution vectors are represented by the blue circle at $\mathbf{s}^* = [0 \ 1 \ 0]^T$ and the green circle, respectively, while \mathcal{H} , the set of all points $\mathbf{s} \in \mathbb{R}^3$ obeying $\mathbf{A}\mathbf{s} = \mathbf{A}\mathbf{s}^*$, is represented by the red line passing

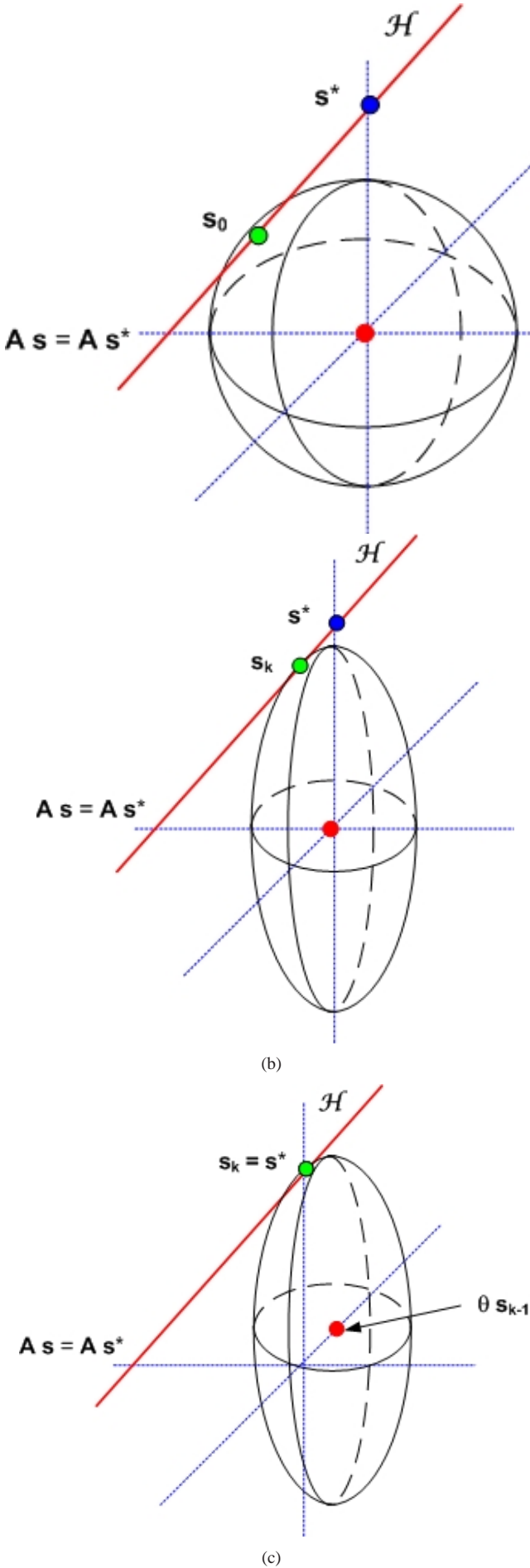


Fig. 3. Geometric interpretation of the solution vector obtained by: (a) $\|s\|_2^2$ minimization, (b) $\|W_{k-1}^{-1} s\|_2^2$ minimization, and (c) $\|W_{k-1}^{-1}(s - \theta s_{k-1})\|_2^2$ minimization.

through s^* . The ℓ_2 minimizer of (6) is the point on \mathcal{H} closest to the origin. This point can be found by blowing up the ℓ_2 -ball, represented by the hypersphere in Fig. 3(a), until it contacts \mathcal{H} . Due to the randomness of the entries of A , \mathcal{H} is oriented at random angle. Accordingly, with high probability, the closest point s_0 will live away from the coordinate axes and hence will be neither sparse nor close to the correct solution vector s^* [34].

Geometrically, incorporating a diagonal weighting matrix into the ℓ_2 -norm causes the ℓ_2 -ball to elongate along certain direction and no longer be symmetric. Assuming that the proposed algorithm (15) is converging to the exact solution vector s^* , then at the k th iteration the ℓ_{w2} -ball will contact \mathcal{H} at s_k , which is closer to the exact solution vector s^* as shown in Fig. 3(b). At the next iteration, and due to the new weighting matrix, the ℓ_{w2} -ball will be further squeezed and contact \mathcal{H} at a point s_{k+1} closer to s^* than the previous solution point s_k . The procedure continues until the ℓ_{w2} -ball contacts \mathcal{H} at s^* .

Since the sparse solution vector s^* is restricted to live at the tip of the ℓ_{w2} -ball, algorithm (15) may take a large number of iterations to converge to that vector, i.e., for the ℓ_{w2} -ball to contact \mathcal{H} at s^* . However, the number of iterations can be significantly reduced if the sparse solution vector s^* is allowed to live on the surface of the ℓ_{w2} -ball. This can be achieved by changing the origin of the ℓ_{w2} -ball in each iteration of (15), where the origin at the k th iteration depends on the previous solution vector. The proposed modification is shown in Fig. 3(c), where the origin of the ℓ_{w2} -ball is shifted to θs_{k-1} , where θ is a parameter to be determined. As shown in this figure, although \mathcal{H} contacts the ℓ_{w2} -ball at a point away from its tip, θ can be chosen so that this point is closer to the exact solution vector s^* than the solution given by the IRLS solution.

Based on this geometric interpretation, we propose selecting B_0 in (9) such that $\tilde{f}(s) = \|\mathbf{W}^{-1}(s - \theta s_0)\|_2^2$, where \mathbf{W} is a diagonal weighting matrix and s_0 is a given solution vector. Since \mathbf{W} is diagonal, $\tilde{f}(s)$ can be written as

$$\tilde{f}(s) = \|\mathbf{W}^{-1}(s - \theta s_0)\|_2^2 = -2\theta s^T \mathbf{W}^{-2} s_0 + s^T \mathbf{W}^{-2} s + C, \quad (18)$$

where C is a constant term that does not depend on s . Equating (9) and (18) we get

$$\begin{aligned} \mathbf{W} &= \mathbf{B}_0^{-\frac{1}{2}}, \\ B_0[i, i] &= \frac{d_0[i]}{2(1 - \theta)s_0[i]}, \quad i = 1, \dots, n, \end{aligned} \quad (19)$$

where the second equation results from the comparison of the first terms in (9) and (18), and utilizing the fact that \mathbf{B}_0 is diagonal.

As described in Section IV-A, the quantity $d_0[i]/s_0[i]$ is always nonnegative. Accordingly, for \mathbf{B}_0 in (19) to be non-negative definite, the value of θ must be less than 1. As will be shown later, this is also a necessary condition for (18) to produce a sequence of solution vectors that minimize the original objective function, i.e., $g(s_{k+1}) < g(s_k)$, for all $k \geq 0$.

Based on this formulation, the optimization problem (10) is reduced to

$$\mathbf{s}_1 = \arg \min_{\mathbf{s}} \quad \|\mathbf{W}_0^{-1}(\mathbf{s} - \theta \mathbf{s}_0)\|_2^2 \quad \text{subject to} \quad \mathbf{x} = \mathbf{A}\mathbf{s}, \quad (20)$$

where \mathbf{W}_0 is a diagonal weighting matrix calculated from (19).

Remark 3: It is worth mentioning that the subtraction strategy in (20) has somewhat the flavour of recent Bregman methods [35], [36] for improving convergence of such constrained problems. It is not difficult to show that the objective function (20) is equivalent to the Bregman distance (based on the objective function $f(\mathbf{s}) = \|\mathbf{W}_0^{-1}\mathbf{s}\|_2^2$) between \mathbf{s} and $\theta \mathbf{s}_0$, i.e., $f_2(\mathbf{s}) = \|\mathbf{W}_0^{-1}(\mathbf{s} - \theta \mathbf{s}_0)\|_2^2 = D_f^{\mathbf{p}}(\mathbf{s}, \theta \mathbf{s}_0)$, where $D_f^{\mathbf{p}}(\mathbf{s}, \theta \mathbf{s}_0)$ is the Bregman distance between \mathbf{s} and $\theta \mathbf{s}_0$ based on the objective function $f(\mathbf{s})$, and $\mathbf{p} \in \partial f(\mathbf{s}_0)$. A topic for future research is to consider utilizing the Bregman distance based on the original objective function $g(\mathbf{s})$, rather than $f(\mathbf{s})$, to derive new algorithms for solving an under-determined linear system of equations.

Equation (20) is not in the standard form of a weighted minimum norm problem. Accordingly, for solving this problem, we follow the following steps. First, assuming that \mathbf{W}_0 is invertible, let \mathbf{y} be defined as

$$\mathbf{y} \triangleq \mathbf{W}_0^{-1}(\mathbf{s} - \theta \mathbf{s}_0).$$

Then, \mathbf{s} is expressed as

$$\mathbf{s} = \mathbf{W}_0 \mathbf{y} + \theta \mathbf{s}_0. \quad (21)$$

Finally, substituting (21) into (20) we get

$$\tilde{\mathbf{y}} = \arg \min_{\mathbf{y}} \quad \|\mathbf{y}\|_2^2 \quad \text{subject to} \quad \bar{\mathbf{x}} = \bar{\mathbf{A}}\mathbf{y}, \quad (22)$$

where $\bar{\mathbf{A}} = \mathbf{A}\mathbf{W}_0$ and $\bar{\mathbf{x}} = \mathbf{x} - \theta \mathbf{A}\mathbf{s}_0 = (1 - \theta)\mathbf{x}$. Equation (22) is a minimum norm problem in the standard form and has the following solution

$$\tilde{\mathbf{y}} = (\bar{\mathbf{A}})^\dagger \bar{\mathbf{x}} = (1 - \theta)(\mathbf{A}\mathbf{W}_0)^\dagger \mathbf{x}. \quad (23)$$

Substituting (23) into (21), the solution vector at the $(k+1)$ th iteration is expressed as

$$\begin{aligned} \mathbf{s}_{k+1} &= \theta \mathbf{s}_k + (1 - \theta) \mathbf{W}_k (\mathbf{A}\mathbf{W}_k)^\dagger \mathbf{x} \\ &= \theta \mathbf{s}_k + (1 - \theta) \mathbf{B}_k^{-1} \mathbf{A}^T (\mathbf{A}\mathbf{B}_k^{-1} \mathbf{A}^T)^{-1} \mathbf{x}, \end{aligned} \quad (24)$$

where \mathbf{B}_k is calculated using (19) and \mathbf{s}_0 and \mathbf{d}_0 have been replaced by \mathbf{s}_k and \mathbf{d}_k , respectively, where $\mathbf{d}_k = \nabla g(\mathbf{s}_k)$. In the remaining part of this paper, the acronyms MCCR and IRLS will be used to refer to (24) and (15), respectively.

Selecting the value of θ : The performance of the MCCR algorithm (24) depends on the value of θ . From (24) we note that multiplying \mathbf{B}_k by any scaling parameter does not affect the value of \mathbf{s}_{k+1} . Accordingly, the constant term in the expression of $B_k[i, i]$ in (19) can be dropped without affecting the value of \mathbf{s}_{k+1} . As a result, the dependency of \mathbf{B}_k on the unknown parameter θ can be dropped by calculating $B_k[i, i]$, or equivalently $B_k^{-1}[i, i]$, using the following equation

$$B_k^{-1}[i, i] = \frac{s_k[i]}{d_k[i]}, \quad i = 1, \dots, n. \quad (25)$$

TABLE I

SOME OBJECTIVE FUNCTIONS THAT CAN BE USED WITH MCCR.

Function	$d[i]$	$B^{-1}[i, i]$
$g(\mathbf{s}) = \sum_i s[i] ^q, 0 < q < 1$	$qs[i] s[i] ^{q-2}$	$ s[i] ^{2-q}$
$g(\mathbf{s}) = \sum_i \log(s[i])$	$\frac{s[i]}{ s[i] ^2}$	$ s[i] ^2$
$g(\mathbf{s}) = \sum_i \log(1 + s[i] /\delta)$	$\frac{s[i]}{ s[i] (\delta + s[i])}$	$ s[i] (\delta + s[i])$
$g(\mathbf{s}) = \sum_i \text{atan}(s[i] /\delta)$	$\frac{\delta s[i]}{ s[i] (\delta^2 + s^2[i])}$	$ s[i] (\delta^2 + s^2[i])$
$g(\mathbf{s}) = \sum_i \frac{ s[i] }{ s[i] + \delta}$	$\frac{\delta s[i]}{ s[i] (\delta + s[i])^2}$	$ s[i] (\delta + s[i])^2$

The expression of $B^{-1}[i, i]$ derived from some objective functions of the form $g(\mathbf{s}) = \sum_i g_c(s[i])$, where $g_c(s[i])$ obeys P2, is presented in Table I. Note that all the constant terms are dropped from the expression of $B^{-1}[i, i]$ in the third column of Table I. The free parameter $\delta > 0$ in the last three objective functions of Table I can be selected to control convergence. Many approaches for selecting the value of δ are suggested in [32]. In this paper the value of δ at the k th iteration is calculated as $\delta_k = 0.5 \text{ mean}(|\mathbf{s}_{k-1}|)$.

After dropping the dependency of $B_k^{-1}[i, i]$ on the unknown parameter θ , Eq. (24) can be written as

$$\mathbf{s}_{k+1} = \theta \mathbf{s}_k + (1 - \theta) \tilde{\mathbf{s}}_{k+1}, \quad (26)$$

where $\tilde{\mathbf{s}}_{k+1} = \mathbf{B}_k^{-1} \mathbf{A}^T (\mathbf{A}\mathbf{B}_k^{-1} \mathbf{A}^T)^{-1} \mathbf{x}$ is the IRLS solution vector obtained using (15). Eq. (26) shows that the MCCR solution at the $(k+1)$ th iteration is an *affine* combination of the solution vector at the k th iteration and the IRLS solution at the $(k+1)$ th vector. Geometrically this means that the solution vector \mathbf{s}_{k+1} exists at a point determined by the parameter θ , somewhere along the line connecting \mathbf{s}_k and $\tilde{\mathbf{s}}_{k+1}$. Note that \mathbf{s}_{k+1} calculated using (26) is always feasible, i.e., $\mathbf{A}\mathbf{s}_{k+1} = \mathbf{x}$. Accordingly, we will refer to the set of all affine combinations of \mathbf{s}_k and $\tilde{\mathbf{s}}_{k+1}$ for various values of θ as the *feasible*-line.

In Section IV-A we proved that the IRLS algorithms derived in this paper are fixed point converging, i.e., for a given feasible solution vector \mathbf{s}_k , the IRLS solution vector $\tilde{\mathbf{s}}_{k+1}$ in (26) satisfies $g(\tilde{\mathbf{s}}_{k+1}) \leq g(\mathbf{s}_k)$. This implies that the direction from \mathbf{s}_k to $\tilde{\mathbf{s}}_{k+1}$ on the *feasible*-line provides a descending direction for the original objective function $g(\mathbf{s})$. This in turn implies that θ must satisfy the condition ($\theta \leq 1$) in order that $g(\mathbf{s}_{k+1}) \leq g(\mathbf{s}_k)$. The optimum value of θ , denoted θ_o , is the one that minimizes $g(\mathbf{s})$ for all \mathbf{s} belonging to the *feasible*-line. Accordingly, θ_o can be determined by solving the following optimization problem

$$\theta_o = \arg \min_{\theta} \quad g(\theta \mathbf{s}_k + (1 - \theta) \tilde{\mathbf{s}}_{k+1}) \quad \text{subject to} \quad \theta_{\min} < \theta < 1, \quad (27)$$

where θ_{\min} is a lower bound on θ . This formulation is readily implemented using, e.g., the Golden-section search method, and can be solved using the Matlab function "fminbnd.m". In the simulation results to be shown we empirically set $\theta_{\min} = -2$.

Note that if no point along the feasible line decreases $g(\mathbf{s})$, then the optimum value of θ is zero. In this case, the MCCR solution given by (24) coincides with the IRLS solution given by (15). Furthermore, if the IRLS solution is a fixed point, then the MCCR solution is also a fixed point.

A further advantage of the MCCR solution (24) over the IRLS solution (15) is that (24) can be easily adapted to solve the following problem

$$\hat{\mathbf{s}} = \arg \min_{\mathbf{s}} g(\mathbf{s}) \quad \text{subject to} \quad \mathbf{s} \geq 0, \quad \mathbf{x} = \mathbf{A}\mathbf{s}. \quad (28)$$

This can be readily accomplished by selecting a starting feasible solution vector and properly selecting the value of $\theta_{min} \leq 1$ in (27) to ensure that $\theta \mathbf{s}_k + (1 - \theta) \tilde{\mathbf{s}}_{k+1} \geq 0$ for all $\theta_{min} \leq \theta \leq 1$.

V. PERFORMANCE ENHANCEMENT

In this section the performance of the derived algorithms is discussed. Since the MCCR solution (24) is an affine combination of the previous solution and the IRLS solution (15), we will consider only the factors that affect the performance of the IRLS solution (15). The discussion presented in this section is general and is applicable to all of the objective functions in Table I. There are two issues that can affect the performance of the updating equation (15). These issues, as well as their mitigating interventions, are now discussed.

A. Inversion Operation Associated with (15)

The first issue is the inversion operation associated with (15). Define $\mathbf{C}_k = \mathbf{A}\mathbf{B}_k^{-1}\mathbf{A}^T$. Then for the solution of (15) to exist, \mathbf{C}_k must be invertible. Recall from (16) that ³

$$\mathbf{B}_k^{-1}[i, i] = \frac{s_k[i]}{d_k[i]}, \quad i = 1, \dots, n. \quad (29)$$

Accordingly, \mathbf{C}_k can be expressed as

$$\mathbf{C}_k = \sum_i^n \frac{s_k[i]}{d_k[i]} \mathbf{a}_i \mathbf{a}_i^T, \quad (30)$$

where \mathbf{a}_i is the i th column of \mathbf{A} . Thus \mathbf{C}_k is a summation of n rank-1 ($m \times m$) matrices. Accordingly, for \mathbf{C}_k to be invertible, at least m elements of \mathbf{s}_k must be significantly different from zero, which contradicts the assumption that \mathbf{s} is sparse. To overcome this difficulty, we redefine $\mathbf{B}_k^{-1}[i, i]$ as follows

$$\mathbf{B}_k^{-1}[i, i] = \begin{cases} \frac{s_k[i]}{d_k[i]} & \text{if } |s_k[i]| \geq \epsilon \\ \frac{s_k[i]}{d_k[i]} \big|_{s_k[i] \leftarrow \epsilon}, & \text{otherwise,} \end{cases} \quad (31)$$

where ϵ is a small positive number. The second line in (31) means that, if the condition $|s_k[i]| \geq \epsilon$ is not satisfied, then $\mathbf{B}_k^{-1}[i, i]$ is calculated by evaluating $s_k[i]/d_k[i]$ first then replacing each $s_k[i]$ by ϵ . For example, for $g_c(s) = |s|^q$ we have $d = qs|s|^{q-2}$, and if $|s| < \epsilon$ then the corresponding entry of \mathbf{B}^{-1} will equal ϵ^{2-q}/q .

Unfortunately, it was shown in [32] that the performance of MCCR depends on the value of ϵ . To overcome this difficulty, we follow the procedure suggested in [19]. Here, ϵ is initiated to a relatively large value, e.g. $\epsilon = 1$. Its value is then reduced by a factor of 10 at iterations where the condition $\|\mathbf{s}_k - \mathbf{s}_{k-1}\|/\|\mathbf{s}_k\| < \sqrt{\epsilon}/100$ is satisfied. The MCCR algorithm is summarized in Table II.

³In (29) the constant associated with $\mathbf{B}_k^{-1}[i, i]$ is neglected.

TABLE II
THE MCCR ALGORITHM

Algorithm 1: The MCCR Algorithm

Given an $(m \times n)$ matrix \mathbf{A} of basis vectors, and a vector $\mathbf{x} \in \mathbb{R}^m$, select one of the value for ϵ , e.g. $\epsilon = 1$, an empirically-selected value for θ_{min} , e.g. $\theta_{min} = -2$, a small \mathbf{s}_0 . This point can be selected as the least squares solution, i.e. $\mathbf{s}_0 = \mathbf{A}^T(\mathbf{A}\mathbf{A}^T)^{-1}\mathbf{x}$. Steps:

- 1) Repeat until convergence:
 - Calculate \mathbf{B}_k^{-1} using (31).
 - Calculate $\tilde{\mathbf{s}}_{k+1} = \mathbf{B}_k^{-1}\mathbf{A}^T(\mathbf{A}\mathbf{B}_k^{-1}\mathbf{A}^T)^{-1}\mathbf{x}$.
 - Calculate θ using (27).
 - Set $\mathbf{s}_{k+1} = \theta \mathbf{s}_k + (1 - \theta) \tilde{\mathbf{s}}_{k+1}$.
 - if $\|\mathbf{s}_{k+1} - \mathbf{s}_k\|_2 / \|\mathbf{s}_{k+1}\|_2 < \sqrt{\epsilon}/100$, set $\epsilon = \epsilon/10$ end
 - Set $k = k + 1$,
- End
- 2) Output \mathbf{s}_{k+1} as the solution.

B. Non-convexity of $g(\mathbf{s})$

The second issue that affects the performance of (15), and hence (24), is the non-convexity of the objective functions considered in this paper. As a result, MCCR may converge to a local minima. This problem may be partially alleviated by the *perturbed* MCCR algorithm (PMCCR) described in Table III. At each iteration, the solution vector \mathbf{s}_j^0 is perturbed by a random noise vector \mathbf{v} , which is constrained to be in the null space of the mixing matrix; i.e. $\mathbf{v} = \mathbf{F}\mathbf{u}$ where \mathbf{F} is any matrix whose range is the null space of \mathbf{A} and $\mathbf{u} \in \mathbb{R}^{n-m}$ is a random noise vector. The columns of \mathbf{F} can be chosen by first calculating the singular value decomposition (SVD) of $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$, and then selecting the columns of \mathbf{F} as the last $(n - m)$ columns of \mathbf{V} . In the simulation results, the elements of \mathbf{u} are sampled from a zero mean uniformly distributed random variable between $\pm \alpha \bar{s}_{max}$, where \bar{s}_{max} is the maximum absolute element in the MCCR solution vector $\bar{\mathbf{s}}$, and α is an empirically determined non-negative number in the range e.g., $1 \leq \alpha \leq 2$. Note that the perturbation noise is constrained to be in the null space of the mixing matrix to insure the feasibility of the new perturbed vector. A stopping criterion could be a maximum number of iterations, or a pre-specified value of the cardinality of the solution. Note that, in each step, the new solution vector is accepted only if its cardinality is less than the cardinality of the previous solution vector. Accordingly, by following this strategy, it is guaranteed that no performance degradation occurs.

VI. SIMULATION RESULTS

In this section, a set of examples are presented in order to examine the effect of θ and the perturbation procedure on the performance of the MCCR algorithm, and to provide a comparison between the MCCR algorithm and other well known algorithms. In the first Example, we present a comparison between the proposed MCCR algorithms and the counterpart IRLS algorithms. This example reflects the effect of the parameter θ . A comparison between MCCR and PMCCR is presented by the second example. Finally, the third example presents a comparison between the MCCR algorithm and some algorithms that are usually utilized for estimating sparse vectors. In all these examples, the parameter ϵ , which

TABLE III
THE PMCCR ALGORITHM

Algorithm 2: The Perturbed MCCR Algorithm (PMCCA)

Initialization: Select one of the objective functions in Table I; select a feasible solution \mathbf{s}_0 to obtain \mathbf{s}_j^0 .

- 1) For the selected objective function, execute MCCR with initial value \mathbf{s}_0 to obtain \mathbf{s}_j^0 .
- 2) Repeat until convergence:
 - evaluate perturbed feasible solution $\mathbf{s}_j^p = \mathbf{s}_j^0 + \mathbf{F}\mathbf{u}$. Refer to the text for details.
 - execute MCCR with initial value \mathbf{s}_j^p to obtain \mathbf{s}_j^1 .
 - if $\|\mathbf{s}_j^1\|_{\ell_0} \leq \|\mathbf{s}_j^0\|_{\ell_0}$, $\mathbf{s}_{j+1}^0 = \mathbf{s}_j^1$.
 - else $\mathbf{s}_{j+1}^0 = \mathbf{s}_j^0$.
 - $j \leftarrow j + 1$

End

- 3) Output \mathbf{s}_{j+1}^0 as the solution.

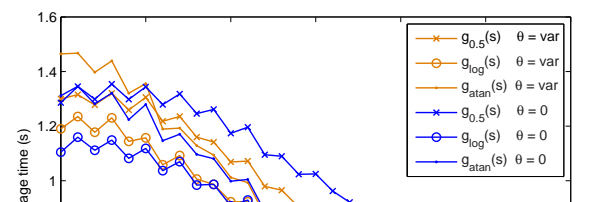
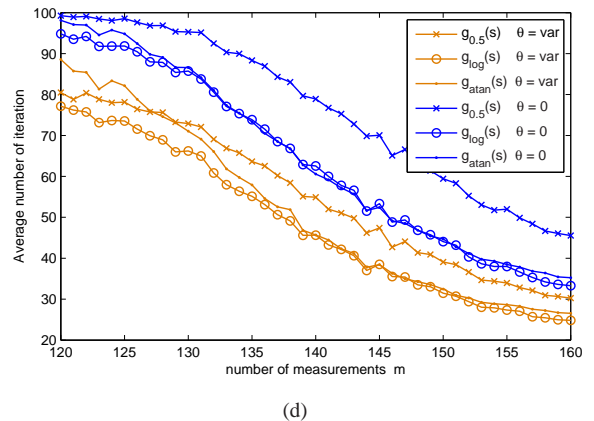
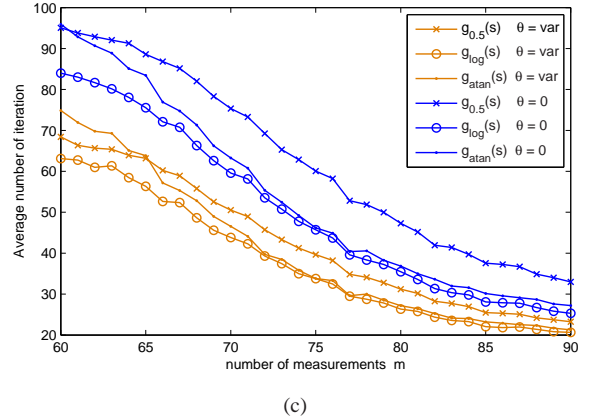
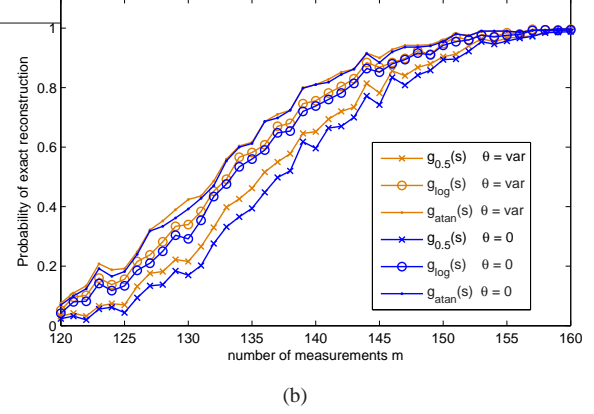
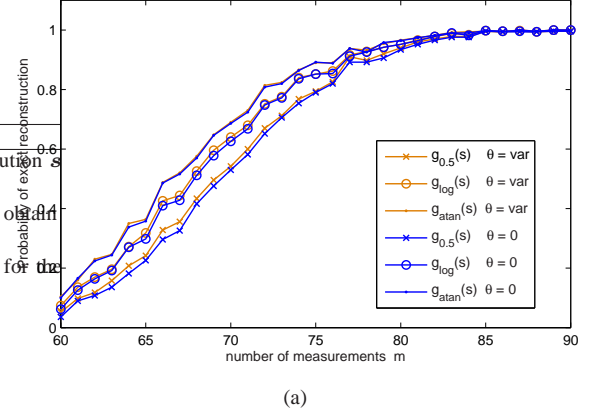
is used with the IRLS, MCCR and PMCCR algorithms, is selected as in Table II, i.e., each algorithm is started with $\epsilon = 1$. Then the value of ϵ is changed to $\epsilon \leftarrow \epsilon/10$ whenever the condition $\|\mathbf{s}_{k+1} - \mathbf{s}_k\|_2 / \|\mathbf{s}_{k+1}\|_2 < \sqrt{\epsilon}/100$ is met.

In all these examples, two parameters are used as measures of performance. The first parameter is the probability of exact reconstruction (PER) of the solution vector. The PER is defined as the ratio between the number of runs at which the algorithm successfully estimates the sparse solution vector, to the total number of runs. The second parameter is the average number of iterations taken by each algorithm to converge to a solution vector. Each one of these two parameters is plotted as a function of the number of measurements (m).

Example 1. Effect of θ : In this Example we compare the effect of θ on the performance of MCCR. Recall that the difference between the MCCR algorithm (24) and the IRLS algorithm (15) is the parameter θ . Three different objective functions are incorporated in this comparison. The functions are $g_{0.5}(s) = |s|^{0.5}$, $g_{\log}(s)$, and $g_{atan}(s)$. Since the computation of θ requires solving a nonlinear optimization problem (27), it is important to compare not just the average number of iterations but also the average amount of time taken by each algorithm to converge to a solution vector. For doing this, the results are calculated for two different conditions. In the first condition, the low dimensional case, the parameters are $l = 30$, $n = 256$, and m increases from 60 to 90, while for the second condition, the high dimensional case, the parameters are $l = 60$, $n = 512$, and m increases from 120 to 160, where l is the number of nonzero entries in the exact solution vector.

For the two cases, the measuring matrix \mathbf{A} and the sparse vector \mathbf{s}^* are generated as follows. For each value of m , a random $(m \times n)$ matrix \mathbf{A} is created whose entries are each Gaussian random variables with zero mean and unit variance. A sparse vector $\mathbf{s}^* \in \mathbb{R}^n$ with l nonzero entries is then created and the corresponding \mathbf{x} is generated as $\mathbf{x} = \mathbf{A}\mathbf{s}^*$. The indices of these l entries of \mathbf{s}^* are randomly selected, and their values are chosen randomly from a zero mean Gaussian random variable with variance = 4. The total number of runs for each case is 500, and the results are shown in Fig. 4. The results for MCCR are presented in faint brown while the results for IRLS ($\theta = 0$) are presented in blue.

As shown in Fig. 4(a) and Fig. 4(b) for the two conditions



and for all the objective functions, MCCR and IRLS have almost similar PERs, meaning that they almost always converge to the same solution vectors. However, comparing the average number of iterations taken by MCCR and IRLS algorithms in Fig. 4(c) and Fig. 4(d) it is clear that MCCR takes fewer iterations than IRLS to converge to the solution vectors. For example, in Fig. 4(d), and for the objective function $g_{0.5}(s)$ and the case $m = 140$ we find that, on average, IRLS converges after 80 trials while MCCR converges after 54 trials only. However, since the computation of θ requires extra computational time, it is important to check the average time taken by each algorithm to converge to a solution vector. This comparison is shown in Fig. 4(e) and Fig. 4(f).

For the low dimensional case presented in Fig. 4(e) it is clear that, except for $g_{0.5}(s)$, the two algorithms take on average the same amount of time to converge to a solution vector, which implies that there is no advantage of replacing the IRLS algorithm by the MCCR algorithm for low dimensional problems. However, the average computational time could be significantly reduced if a faster algorithm is used for calculating θ , or if the computational cost of θ is low compared with the inversion operation in (24). This later case is presented in Fig. 4(f). From this figure it is clear that MCCR in this case converges in a shorter time than IRLS. Accordingly, we conclude this example by recommending the MCCR algorithm to solve large scale problems.

Example 2. Comparison between MCCR and PMCCR :

In this example we examine the effect of the perturbation approach on the performance of the PMCCR. Specifically, we compare the performance of the MCCR and the PMCCR algorithm for the case $g(s) = \|s\|_{\ell_q}^q$ for the following values of q ; 0.1, 0.5, and 0.9. The comparison is made in terms of the probability of exact reconstruction of the sparse vector, and the average number of iterations required for convergence to that solution vector, as functions of the number of observations m , when both n and l are fixed at 40 and 3, respectively. For each value of m , A and s^* are generated as in Example 1. For the two algorithms, the values of ϵ and θ are selected as in Table II. The results are shown in Fig. 5.

As shown in Fig. 5(a) the probability of the exact reconstruction of the PMCCR algorithm is significantly better than that of the MCCR algorithm. For a low value of m , e.g., $m = 10$, the probability of exact reconstruction for PMCCR is 100% for $q = 0.1$ and 0.5. The number of iterations taken by each algorithm is presented in Fig. 5(b). The number of iterations of the PMCCR algorithm presented in Fig. 5(b) is the number of times the MCCR algorithm was called by the PMCCR algorithm, i.e., the final value of j in Table III. It is obvious that the PMCCR algorithm takes large number of iterations compared with the MCCR algorithm, especially when the number of measurements is small, e.g., $m < 10$ in this case. However, for larger m , e.g., $m \geq 15$, the performance of PMCCR in terms of successful reconstruction is 100%, while the number of outer iterations required is roughly less than or equal to 3. Therefore, PMCCR can be seen as a high-performance, more costly method for difficult cases (i.e., low values of m), whereas it is a higher-performance method at roughly the same cost, for cases which work well

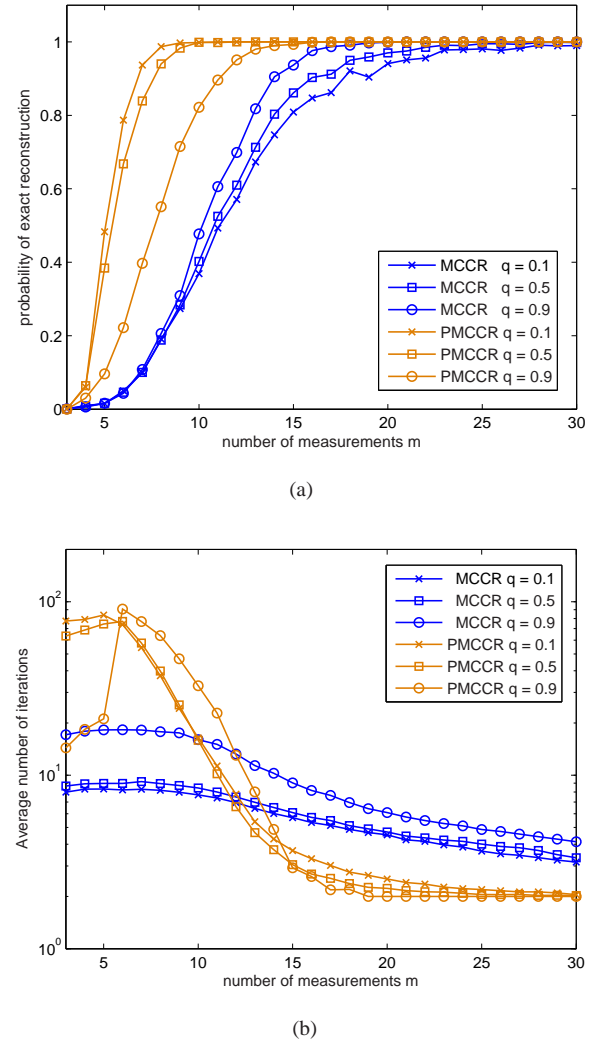


Fig. 5. Comparison between MCCR and PMCCR. (a) The probability of exact reconstruction of the original sparse vector; (b) The average number of iterations required to get a sparse solution vector. The number of iterations of the PMCCR algorithm represents the value of j in Table III.

with other methods.

Example 3. Comparison between MCCR and other algorithms: We conclude this paper by a comparison between MCCR and four different algorithms. The four algorithms are: the ℓ_1 -norm; the weighted ℓ_1 -norm [12], denoted as ℓ_{w1} -norm; the $\ell_{0.5}$ -norm implemented using the algorithm presented in [19], which is equivalent to MCCR when $g_c(s) = |s|^{0.5}$ and $\theta = 0$; and the recently developed algorithm called smooth ℓ^0 -norm algorithm [37], denoted as Sl_0 in Fig. 6. The objective function used with MCCR is $g_{atan}(s)$. In this example, $n = 256$, $l = 40$, and m increases from 80 to 140. For each value of m , A and s^* are generated as in Example 1. The total number of runs in this example is 200, and the results are presented in Fig. 6.

As shown in Fig. 6(a) MCCR has the best performance in terms of the PER, while the smooth ℓ^0 -norm algorithm has the worst performance among the compared algorithms. Also it is clear from Fig. 6(a) that the PER of $\ell_{0.5}$ -norm algorithm is very close to that of MCCR and the ℓ_{w1} -norm performs better than the regular ℓ_1 -norm. However, in terms of the average

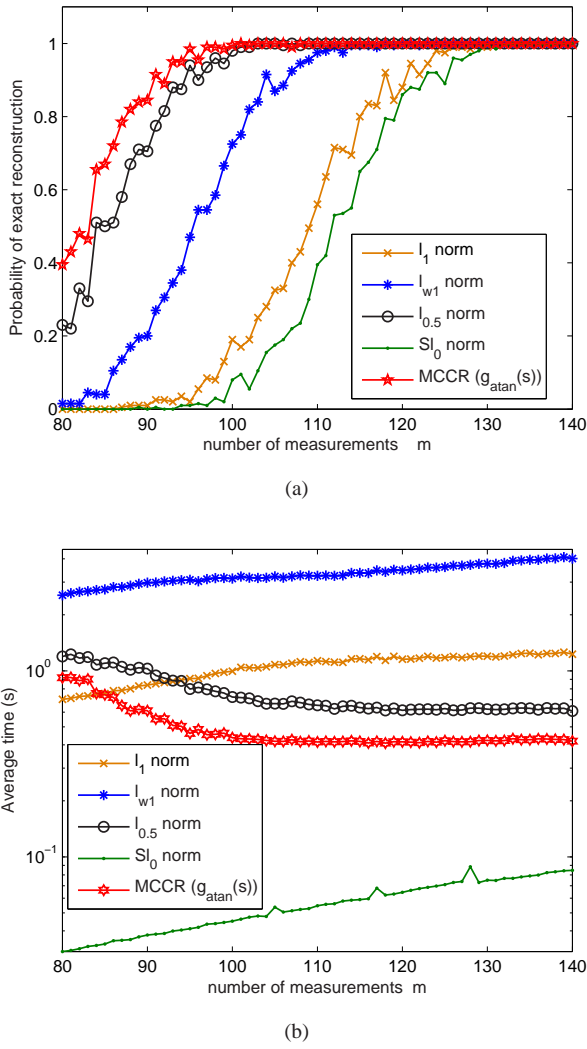


Fig. 6. Comparison between MCCR and four different algorithms. (a) The probability of exact reconstruction of the original sparse vector as a function of the number of measurements (m); (b) The average number of iterations required to get a sparse solution vector.

amount of time shown in Fig. 6(b) it is clear that, aside from the smooth ℓ^0 -norm, the MCCR is the fastest algorithm, while the ℓ_{w1} -norm is the most computationally expensive one. Although both MCCR and $\ell_{0.5}$ -norm algorithms estimated all the solution vectors exactly for all $m > 100$, the MCCR converged in a significantly shorter time. For example, for the case when $m = 100$, MCCR converged in 0.43 seconds while the $\ell_{0.5}$ -norm algorithm converged in 0.72 seconds.

VII. CONCLUSION

In this paper a novel methodology was developed and employed to minimize a class of non-convex (concave on the non-negative orthant) functions for solving an under-determined system of linear equations for the case of sparse solution vector. The proposed technique is based on locally replacing the original objective function by a quadratic convex function which is easily minimized. It was shown in this paper, for a certain selection of the convex objective function, the Iterative Re-weighted Least Squares (IRLS) class of algorithms are a special case of the proposed methodology. Thus the proposed algorithms are a generalization and unification

of the previous methods. In this paper we also proposed a convex objective function that produces an algorithm that can converge to the sparse solution vector in a significantly fewer number of iterations than the IRLS algorithms. Other selections of the convex objective function may produce algorithms with convergence properties better than the IRLS algorithms.

In this paper we also proposed a straightforward technique for selecting a convex function such that, for any starting solution vector s_0 , the algorithm generates a sequence $\{s_k\}_{k=1}^{\infty}$ that converges to a fixed point of the original objective function. Since the original objective functions are non-convex, the proposed algorithm is susceptible to convergence to a local minimum. To alleviate this difficulty, we proposed a random perturbation technique that enhances the performance of the proposed algorithm. Simulation results were presented to examine the performance of the proposed algorithm and to compare its performance with some well known algorithms that are usually utilized for solving the same problem. The simulation results show that the proposed algorithm outperforms the existing algorithms in terms of the execution time and the accuracy of reconstructing a sparse solution vector.

REFERENCES

- [1] P. Rodríguez and B. Wohlberg, "An iterative reweighted norm algorithm for total variation regularization," *IEEE Signal Processing Letters*, vol. 14, no. 12, pp. 948–951, Dec 2007.
- [2] S. G. Mallat and Z. Zhang, "Matching pursuits with time-frequency dictionaries," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. 41, pp. 33973415, Dec. 1993.
- [3] G. H. Golub and C. F. Van Loan, "Matrix Computations." Baltimore: The Johns Hopkins University Press, 1989.
- [4] S. Lesage, R. Gribonval, F. Bimbot, and L. Benaroya, "Learning unions of orthonormal bases with thresholded singular value decomposition," in *Proc. IEEE Int. Conf. Acoust. Speech Signal Process.*, Philadelphia, PA, Mar. 1823, 2005, vol. 5, pp. v/293v/296.
- [5] Y. Li, A. Cichocki, and S. Amari, "Sparse Component Analysis for Blind Source Separation With Less Sensors Than Sources," in *Proc. 4th Int. Symp. Independent Component Analysis Blind Source Separation (ICA BSS)*, 2003, pp. 89–94.
- [6] Y. Li, A. Cichocki, S. Amari, S. Shishkin, J. Cao, and F. Gu, "Sparse Representation and Its Applications in Blind Source Separation," in *7th Annual Conference on Neural Information Processing Systems (NIPS-2003)*, (Vancouver), Dec. 2003.
- [7] I. F. Gorodnitsky and B. D. Rao, "Sparse signal reconstructions from limited data using FOCUSS: A re-weighted minimum norm algorithm," *IEEE Trans. Signal Processing*, vol. 45, pp. 600–616, Mar. 1997.
- [8] P. Xu, Y. Tian, H. Chen, and D. Yao, "LP norm iterative sparse solution for EEG source localization," *IEEE Trans. Biomed. Eng.*, vol. 54, no. 3, pp. 400409, Mar. 2007.
- [9] K. Kreutz-Delgado and B. D. Rao, "FOCUSS-based dictionary learning algorithms," in *Proc. Wavelet Appl. Signal Image Process.*, Bellingham, WA, Jul.-Aug. 2000, vol. 4119, pp. 459473.
- [10] M. Aharon, M. Elad, and A. Bruckstein, "The K-SVD: An algorithm for designing overcomplete dictionaries for sparse representation," *IEEE Trans. Signal Process.*, vol. 54, no. 11, pp. 43114322, Nov. 2006.
- [11] B. Jeffs and M. Gonsay, "Restoration of blurred star field images by maximally sparse optimization," *IEEE Trans. Image Process.*, vol. 2, no. 2, pp. 202211, Mar. 1993.
- [12] E. J. Candès, M. B. Wakin, and S. P. Boyd, "Enhancing sparsity by reweighted ℓ_1 minimization," *Journal of Fourier Analysis and Applications*, vol. 14, no. 5, pp. 877–905, special issue on sparsity, December 2008.
- [13] S. Mallat and Z. Zhang, "Matching pursuits with time-frequency dictionaries," *IEEE Trans. Signal Process.*, vol. 41, no. 12, pp. 33973415, Dec. 1993.
- [14] J. A. Tropp, A. C. Gilbert, and M. J. Strauss, "Algorithms for simultaneous sparse approximation, Part I: Greedy pursuit," *Signal Process.*, vol. 86, no. 3, pp. 572588, 2006.

- [15] S. Chen and D. Donoho, "Basis pursuit," in Proc. Twenty-Eighth Asilomar Conf. Signals, Syst. Comput., Monterey, CA, Nov. 1994, vol. I, pp. 41-44.
- [16] S. S. Chen, D. L. Donoho, and M. A. Saund, "Atomic decomposition by basis pursuit," SIAM J. Sci. Comput., vol. 20, no. 1, pp. 3361, 1998.
- [17] T. Blumensath, M. E. Davies, "Iterative Thresholding for Sparse Approximations," The Journal of Fourier Analysis and Applications, vol. 14, no. 5, pp. 629-654, Dec. 2008
- [18] B. D. Rao and K. Kreutz-Delgado, "An affine scaling methodology for best basis selection," IEEE Trans. Signal Process., vol. 47, no. 1, Jan. 1999.
- [19] R. Chartrand and W. Yin, "Iteratively Reweighted Algorithms for Compressive Sensing," in Proc. Acoustics, Speech and Signal Processing, ICASSP 2008, pp. 3869-3872.
- [20] I. Daubechies, R. DeVore, M. Fornasier, and S. Gunturk, " Iteratively re-weighted least squares minimization for sparse recovery," to appear in Commun. Pure Appl. Math. 2009.
- [21] D. Baron, M. F. Duarte, M. B. Wakin, S. Sarvotham, and R. G. Baraniuk, "Distributed Compressive Sensing," arXiv:0901.3403v1, Jan. 2009.
- [22] E. Candes and T. Tao, "Near optimal signal recovery from random projections: universal, encoding strategies?" IEEE Trans. Inf. Theory, vol. 52, pp. 5406-5425, 2006.
- [23] D. L. Donoho and J. Tanner, "Thresholds for the recovery of sparse solutions via L1 minimization," 40th Annual Conference on Information Sciences and Systems, 2006, pp 202-206.
- [24] R. Chartrand, "Exact reconstruction of sparse signals via nonconvex minimization," IEEE Signal Process. Lett., vol. 14, no. 10, Oct. 2007.
- [25] N. Mourad and J. Reilly, " ℓ_p Minimization for Sparse Vector Reconstruction", ICASSP2009, Taipei, Taiwan, Apr. 19-24, 2009.
- [26] R. Chartrand and V. Staneva, "Restricted isometry properties and non-convex compressive sensing", Inverse Problems, vol. 24, no. 035020, pp. 1-14, 2008.
- [27] Y. Li, "A globally convergent method for ℓ_p problems," SIAM J. Optim. vol. 3, no. 3, pp. 609-629, 1993.
- [28] D. L. Donoho and Y. Tsaig, "Fast solution of ℓ_1 -norm minimization problems when the solution may be sparse," preprint (2006).
- [29] R. Wolke and H. Schwetlick, "Iteratively reweighted least squares: algorithms, convergence analysis, and numerical comparisons," SIAM J. Sci. Stat. Comput., vol. 9, no. 5, pp. 907-921, 1988.
- [30] R. Wolke, "Iteratively reweighted least squares: A comparison of several single step algorithms for linear models," BIT Numerical Mathematics, vol. 32, no. 3, pp. 506-524, 1992.
- [31] S. L. Campbell and C. D. Meyer, "Generalized inverse of linear transformations," London, U. K.: Pitman, 1979.
- [32] Nasser Mourad, "Advances in sparse signal analysis with applications to blind source separation and EEG/MEG signal processing," Ph.D. dissertation, McMaster University, Hamilton, On., Canada, 2009.
- [33] D. R. Hunter and K. Lange, "A tutorial on MM algorithms," The American Statistician, vol. 58, no. 1, pp. 30-37, 2004.
- [34] R. G. Baraniuk, "Compressive Sensing," IEEE Trans. Signal Proc., vol. 24, no. 4, pp. 118 - 121, July 2007.
- [35] S. Osher, Y. Mao, B. Dong, and A. W. Yin, "Fast linearized bregman iteration for compressive sensing and sparse denoising," Available at <ftp://ftp.math.ucla.edu/pub/camreport/cam08-37.pdf>
- [36] W. Yin, S. Osher, D. Goldfarb, and J. Darbon, "Bregman iterative algorithms for ℓ_1 -minimization with applications to compressed sensing," SIAM J. Imaging Science, vol. 1, no. 1, pp. 143-168, 2008.
- [37] G. H. Mohimani, M. Babaie-Zadah, and C. Jutten, "A Fast approach for overcomplete sparse decomposition based on smoothed ℓ^0 norm," IEEE Trans. Signal Processing, vol. 57, no. 1, pp. 289-301, Jan. 2009.