

# Efficient Sparse Recovery Pursuits with Least Squares

Guy Leibovitz\*, Raja Giryes†

\*guyleibovitz@mail.tau.ac.il, †raja@tauex.tau.ac.il  
School of Electrical Engineering, Tel-Aviv University  
Tel-Aviv, Israel

## Abstract

We present a new greedy strategy, with an efficient implementation technique, that enjoys similar computational complexity and stopping criteria like OMP. Moreover, its recovery performance in the noise free and the Gaussian noise cases is comparable and in many cases better than other existing sparse recovery algorithms both with respect to the theoretical and empirical reconstruction ability. Our framework has other appealing properties. Convergence is always guaranteed by design also in the case that the recovery conditions are violated. In addition, our implementation method is useful for improving the computational costs of other greedy pursuits such as orthogonal least squares (OLS).

## I. INTRODUCTION

In this paper we consider the problem of finding a sparse representation  $\mathbf{x} \in \mathbb{R}^n$  in a dictionary  $\mathbf{D} \in \mathbb{R}^{m \times n}$  for the vector  $\mathbf{y} \in \mathbb{R}^m$  such that  $\mathbf{y} = \mathbf{D}\mathbf{x} + \mathbf{w}$ , where  $\mathbf{w}$  is an added Gaussian noise vector, and  $m \leq n$ . The sparsity  $k$  (number of non-zero elements) of  $\mathbf{x}$  is its  $\ell_0$  pseudo-norm  $\|\mathbf{x}\|_0 = k$ . This problem received a lot of attention in the signal processing and statistics communities. It is prevalent in areas such as regression [22], signal and image restoration [3], and compressed sensing [10]. Finding the best  $k$ -sparse approximation of the vector  $\mathbf{y}$  is proven to be NP-hard, and thus only approximate solutions can be obtained in polynomial time, unless the dictionary satisfies some regularity condition. A popular condition is the Restricted Isometry Property (RIP) [5] that defines a condition on subsets of columns of  $\mathbf{D}$  as follows: A matrix  $\mathbf{D}$  is said to satisfy the RIP of order  $s$  (an integer) with a constant  $\delta_s \in (0, 1)$  if

$$(1 - \delta_s)\|\mathbf{v}\|^2 \leq \|\mathbf{D}_T\mathbf{v}\|^2 \leq (1 + \delta_s)\|\mathbf{v}\|^2, \quad \forall |T| \leq s, \mathbf{v} \in \mathbb{R}^s, \quad (\text{I.1})$$

where the matrix  $\mathbf{D}_T$  is comprised of the columns of  $\mathbf{D}$  indexed by the set of indices  $T$ . One of the common approaches for sparse signal estimation relies on  $\ell_1$ -based convex optimization [7, 22]. This strategy has been proven to be very potent in finding the best sparse approximation under the RIP condition [4], but  $\ell_1$  methods generally require more computations than another widespread methodology, the so-called greedy pursuit approach that includes the popular Orthogonal Matching Pursuit (OMP) and Orthogonal least squares (OLS)<sup>1</sup>. OMP and OLS find the sparse representation by greedily selecting one atom of the dictionary  $\mathbf{D}$  at a time. OMP picks the most correlated atom of  $\mathbf{D}$  to the residual error between  $\mathbf{y}$  and its temporal least-squares estimate with the currently selected columns. OLS on the other hand, improves over OMP by picking the atom that yields the smallest approximation error. However, this comes with an additional computational cost (order of  $k$  times the OMP complexity). Optimized OMP (OOMP) [20] is an acceleration of OLS, but requires a significant amount of added memory. In Section S1 of the supplementary material we propose a more efficient method for accelerating OLS based on the tools we develop in this paper. Both OLS and OMP use a one-off strategy, where an atom never leaves the selected support after it enters. One option for improving one-off methods is called back-tracing, the re-consideration of atoms in the selected support. Pursuit algorithms that use this approach include: CoSaMP and Subspace Pursuit (SP) [8, 12, 19, 21], which received a lot of attention for their high reconstruction guarantees ( $\delta_{4k} \leq 0.3843$  and  $\delta_{3k} \leq 0.4859$  respectively, [11, 21]) Another option is OMP with Replacement (OMPR) [14]. The guarantees for OMPR are also based on the RIP but they impose an additional constraint on the mutual coherence of  $\mathbf{D}$  ( $\delta_2$ ) and on a parameter of this algorithm that requires tuning.

In this study we propose new greedy methods that resemble OMP and OLS. Both of our proposed pursuits are designed as such that there is no parameter affecting their convergence rate or accuracy, other than an input of either target sparsity (in both of them) or target residual (in one method) similar to the way OMP and OLS work. Our techniques are guaranteed to converge to the true solution under better RIP conditions compared to other greedy strategies. We validate the supremacy of our suggested programs also in simulations.

This paper is organized as follows: Section II briefly describes the OLS algorithm and presents some preliminary lemmas that are used later with the proposed algorithms in Section III. Section IV contains the theoretical performance guarantees and their proofs. Section V follows with numerical examples, and Section VI concludes the paper.

<sup>1</sup> The OLS algorithm entertains a plethora of designations in the literature being rediscovered under different names several times. For example, in statistics it is known as forward stepwise regression. Other names include Least Squares Orthogonal Matching Pursuit (LSOMP), forward-regression, and Order Recursive Matching Pursuit (ORMP) to name a few. See [2] for discussion and literature survey.

## II. PRELIMINARIES

We use the following notation in this paper:  $\mathbf{a}$  is a vector,  $\vec{\mathbf{a}} = \mathbf{a}/\|\mathbf{a}\|$  is a unit vector,  $a$  is a scalar, and  $\mathbf{A}$  is a matrix.  $\mathbf{a}(i)$  is the  $i$ -th entry of  $\mathbf{a}$ , whereas  $\mathbf{a}_{\setminus i}$  is the vector  $\mathbf{a}$  without entry  $i$ , similarly  $\mathbf{A}_{\setminus i}$  is  $\mathbf{A}$  without column  $i$ , and the  $i$ -th column and row are deleted from  $\mathbf{A}$  to get  $\mathbf{A}_{\setminus i, \setminus i}$ . Unless otherwise noted in the subscript  $\|\cdot\|$  is the usual  $\ell_2$  norm.  $\mathbf{D} \in \mathbb{R}^{m \times n}$  is the dictionary,  $\mathbf{d}_{(i)}$  designates the  $i$ -th atom (or column) of  $\mathbf{D}$  ( $\|\mathbf{d}_{(i)}\|_2 = 1 \quad \forall i$ ).  $\mathbf{x} \in \mathbb{R}^n$  is the unknown sparse representation with sparsity  $k = \|\mathbf{x}\|_0$ , and  $\mathbf{y} = \mathbf{D}\mathbf{x} + \mathbf{w}$  is the signal we have. Uppercase non-bold letters are sets of indices (e.g.  $T$ ). We denote by  $\mathbf{A} = \mathbf{D}_T$  the sub-matrix of  $\mathbf{D}$  comprised of the columns indexed by  $T$ , and interchange the notations at convenience. The least squares estimate of  $\mathbf{y}$  using the atoms in  $T$  is denoted as  $\hat{\mathbf{x}}_{\mathbf{A}} = \hat{\mathbf{x}}_T = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}$ ; the orthogonal projection onto the column space of  $\mathbf{A}$  as  $P_{\mathbf{A}} = P_T = \mathbf{A}(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$ ; and the orthogonal complement  $R_{\mathbf{A}} = R_T = I - P_{\mathbf{A}}$ . Finally,  $\odot$  represents an element-wise multiplication.

The following are preliminary lemmas that will aid us throughout this paper in the derivation of two new greedy techniques together with their theoretical recovery guarantees. We start with two variants of the Sherman-Morrison matrix inversion lemma for a column addition and deletion [15, 16] that follow from a simple modification of the original lemma.

**Lemma II.1** (Sherman-Morrison matrix inversion lemma for column addition). *Let  $\mathbf{B} = (\mathbf{A}^T \mathbf{A})^{-1}$  and  $\tilde{\mathbf{A}} = [\mathbf{A} \ \mathbf{a}]$ . Then we may calculate  $\tilde{\mathbf{B}} = (\tilde{\mathbf{A}}^T \tilde{\mathbf{A}})^{-1}$  as follows:*

$$\tilde{\mathbf{B}} = ([\mathbf{A} \ \mathbf{a}]^T [\mathbf{A} \ \mathbf{a}])^{-1} = \begin{bmatrix} \mathbf{B} & \mathbf{0} \\ \mathbf{0} & 0 \end{bmatrix} + \frac{1}{r} \begin{bmatrix} \hat{\mathbf{e}} \\ -1 \end{bmatrix} [\hat{\mathbf{e}}^T \ -1], \quad (\text{II.1})$$

where  $\hat{\mathbf{e}} = \mathbf{A}^\dagger \mathbf{a}$ , and  $r = \|R_{\mathbf{A}} \mathbf{a}\|^2$ .

A straight forward consequence of Lemma II.1 is the following update for column removal:

**Lemma II.2** (Sherman-Morrison matrix inversion lemma for column removal). *Let  $\tilde{\mathbf{B}} = (\tilde{\mathbf{A}}^T \tilde{\mathbf{A}})^{-1}$  with  $\tilde{\mathbf{A}} = [\mathbf{A} \ \mathbf{a}]$ . Then we may calculate  $\mathbf{B} = (\mathbf{A}^T \mathbf{A})^{-1}$  as follows:*

$$\hat{\mathbf{e}} = -\tilde{\mathbf{B}}_{\setminus i, \setminus i}(:, i), \quad r = \tilde{\mathbf{B}}(i, i)^{-1}, \quad \mathbf{B} = \tilde{\mathbf{B}}_{\setminus i, \setminus i} - r \hat{\mathbf{e}} \hat{\mathbf{e}}^T. \quad (\text{II.2})$$

Note that in order to use Lemma II.1 for column insertion at a general location in  $\mathbf{A}$ , simply insert  $\mathbf{a}$  at the last index and permute  $\tilde{\mathbf{A}}$  and  $\tilde{\mathbf{B}}$  afterwards.

**Lemma II.3** (The change in error). *Let the mean squared error of estimating  $\mathbf{y}$  using  $\mathbf{A}$  be  $\|\mathbf{y} - \hat{\mathbf{y}}\|^2 = \|\mathbf{y} - P_{\mathbf{A}} \mathbf{y}\|^2 = \|R_{\mathbf{A}} \mathbf{y}\|^2$ . Then, the change in  $\|R_{\tilde{\mathbf{A}}} \mathbf{y}\|^2 = \|R_{[\mathbf{A} \ \mathbf{a}]} \mathbf{y}\|^2$  after a vector addition, in terms of  $\|R_{\mathbf{A}} \mathbf{y}\|^2$  can be written as:*

$$\langle \overrightarrow{R_{\tilde{\mathbf{A}}} \mathbf{a}}, \mathbf{y} \rangle^2 = (\mathbf{y}^T \overrightarrow{R_{\tilde{\mathbf{A}}} \mathbf{a}})^2 = \|R_{\mathbf{A}} \mathbf{y}\|^2 - \|R_{\tilde{\mathbf{A}}} \mathbf{y}\|^2 = \|P_{\tilde{\mathbf{A}}} \mathbf{y}\|^2 - \|P_{\mathbf{A}} \mathbf{y}\|^2. \quad (\text{II.3})$$

*Proof.* By utilizing lemma II.1, we can write

$$\|R_{\tilde{\mathbf{A}}} \mathbf{y}\|^2 = \|\mathbf{y} - P_{\tilde{\mathbf{A}}} \mathbf{y}\|^2 \stackrel{(a)}{=} \|\mathbf{y}\|^2 - \mathbf{y}^T \mathbf{A} \mathbf{B} \mathbf{A}^T \mathbf{y} - \frac{1}{r} \left( \mathbf{y}^T [\mathbf{A} \ \mathbf{a}] \begin{bmatrix} \hat{\mathbf{e}} \\ -1 \end{bmatrix} \right)^2, \quad (\text{II.4})$$

where (a) follows from (II.1) and some simple algebraic steps. Since  $\|R_{\mathbf{A}} \mathbf{y}\|^2 = \|\mathbf{y}\|^2 - \|P_{\mathbf{A}} \mathbf{y}\|^2 = \|\mathbf{y}\|^2 - \mathbf{y}^T \mathbf{A} \mathbf{B} \mathbf{A}^T \mathbf{y}$ , we have,

$$\|R_{\tilde{\mathbf{A}}} \mathbf{y}\|^2 = \|R_{\mathbf{A}} \mathbf{y}\|^2 - \frac{1}{r} (\mathbf{y}^T (\mathbf{A} \hat{\mathbf{e}} - \mathbf{a}))^2 = \|R_{\mathbf{A}} \mathbf{y}\|^2 - (\mathbf{y}^T \overrightarrow{R_{\tilde{\mathbf{A}}} \mathbf{a}})^2. \quad (\text{II.5})$$

Reordering (II.5) together with the fact that  $\|P_{\mathbf{A}} \mathbf{y}\|^2 = \|\mathbf{y}\|^2 - \|R_{\mathbf{A}} \mathbf{y}\|^2$  leads to (II.3).  $\square$

**Lemma II.4** (The value of  $\hat{\mathbf{x}}_{\mathbf{A}}(i)$ ). *Let  $\hat{\mathbf{x}}_{\mathbf{A}} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}$ , the least squares estimate of  $\mathbf{x}$ . It is in effect the representation of  $\mathbf{y}$  in a bi-orthogonal basis for the space spanned by the columns of  $\mathbf{A}$ , i.e., its  $i$ -th entry is*

$$\hat{\mathbf{x}}_{\mathbf{A}}(i) = \frac{1}{\|R_{\mathbf{A}_{\setminus i}} \mathbf{d}_{(i)}\|^2} \langle R_{\mathbf{A}_{\setminus i}} \mathbf{d}_{(i)}, \mathbf{y} \rangle. \quad (\text{II.6})$$

*Proof.* Without loss of generality, we prove this formula for the last entry of  $\hat{\mathbf{x}}_{\mathbf{A}}$  (i.e.  $i = k$ ). Let  $\tilde{\mathbf{A}} = [\mathbf{A} \ \mathbf{a}]$ , with  $k$  columns, consider the  $k$ -th entry of the expression  $\hat{\mathbf{x}}_{\tilde{\mathbf{A}}} = (\tilde{\mathbf{A}}^T \tilde{\mathbf{A}})^{-1} \tilde{\mathbf{A}}^T \mathbf{y} = \tilde{\mathbf{B}} \tilde{\mathbf{A}}^T \mathbf{y}$ . Using Lemma II.1 together with the notation therein we may rewrite  $\hat{\mathbf{x}}_{\tilde{\mathbf{A}}}$  as (recall  $\mathbf{A} = \tilde{\mathbf{A}}_{\setminus k}$ ):

$$\hat{\mathbf{x}}_{\tilde{\mathbf{A}}}(k) = \frac{1}{\|R_{\mathbf{A}} \mathbf{a}\|^2} (\mathbf{a}^T - \mathbf{a}^T \mathbf{A} \mathbf{B} \mathbf{A}^T) \mathbf{y},$$

which equals (II.6) for  $i = k$ . By permuting the entries of  $\hat{\mathbf{x}}$  and  $\tilde{\mathbf{A}}$ , the same can be proved for all  $i \in [1, k]$ .  $\square$

To conclude this section, the following lemma is an interesting consequence of the results above.

**Lemma II.5.** Let  $\hat{\mathbf{x}} = \tilde{\mathbf{B}}\tilde{\mathbf{A}}^T\mathbf{y}$ , with  $\tilde{\mathbf{B}} = (\tilde{\mathbf{A}}^T\tilde{\mathbf{A}})^{-1}$ , then the least contributing column  $i$  of  $\tilde{\mathbf{A}}$  for the least squares estimate of  $\mathbf{x}$  from  $\mathbf{y}$  is the one whose index corresponds to

$$\arg \min_i \hat{\mathbf{x}}(i)^2 / \tilde{\mathbf{B}}(i, i). \quad (\text{II.7})$$

*Proof.* From Lemma (II.6) we have that

$$\hat{\mathbf{x}}_{\mathbf{A}}(i)^2 = \frac{1}{\|R_{\mathbf{A}_{\setminus i}}\mathbf{a}\|^4} \langle \overrightarrow{R_{\mathbf{A}}\mathbf{a}}, \mathbf{y} \rangle^2.$$

Combining this with the expression for  $r$  in (II.2) that specifies the value of the diagonal of  $\tilde{\mathbf{B}}$ , we get that  $\langle \overrightarrow{R_{\mathbf{A}}\mathbf{a}}, \mathbf{y} \rangle = \hat{\mathbf{x}}(i)^2 / \tilde{\mathbf{B}}(i, i)$ . From Lemma (II.3) we have  $\|R_{\mathbf{A}}\mathbf{y}\|^2 - \|R_{\tilde{\mathbf{A}}}\mathbf{y}\|^2 = \langle \overrightarrow{R_{\mathbf{A}}\mathbf{a}}, \mathbf{y} \rangle$  which validates the claim.  $\square$

#### A. Orthogonal Least Squares (OLS)

Lemma II.3 implies that given a set of columns  $\mathbf{A}$  used for the estimation of a signal, adding to it a column that satisfies  $\mathbf{d} = \arg \max_{\mathbf{d} \in \mathcal{D}} (\mathbf{y}^T \overrightarrow{R_{\mathbf{A}}\mathbf{d}})^2$  yields the smallest residual among all atoms in the dictionary. This notion is the basis for OLS. Notice that in OMP (and other methods such as SP, CoSaMP and OMPR ) the selection criterion is based on  $\arg \max_{\mathbf{d} \in \mathcal{D}} (\mathbf{y}^T R_{\mathbf{A}}\mathbf{d})^2$ , lacking the normalization by  $\|R_{\mathbf{A}}\mathbf{d}\|^2$ . OLS enlarges its selected support iteratively, where at each step a new atom that satisfies Lemma II.3 is selected. We relegate to the supplementary material the OLS pseudo-code and an acceleration we propose for it.

### III. PROPOSED ALGORITHMS

In this section we introduce two new algorithms that utilize a similar approach to OLS, by selecting atoms for inclusion/exclusion based on  $\langle \overrightarrow{R_T\mathbf{d}_{(j)}}, \mathbf{y} \rangle$  as the metric.

#### A. Orthogonal least squares with replacement (OLSR)

The first method we introduce is called OLS with replacement since it starts by producing a support estimation of size  $k + 1$ , similarly to OLS, followed by subsequent attempts to improve the residual by replacing atoms in the support according to Lemma II.3, until convergence occurs. At each step OLSR extracts an atom from the selected support (of size  $k + 1$ ), then, if an atom exists that lowers the residual compared to the beginning of the iteration it is inserted to the support. The algorithm stops when there are no such atoms in the dictionary, ending up with a support of size  $k$ .

#### B. Iterative orthogonal least squares with replacement (IOLSR)

The second improvement to OLS (in the same computational cost) introduced in this paper is the IOLSR algorithm. It resembles Efroymsen's algorithm ([18], section 3.3) for stepwise regression. At each step a new atom enters the support according to the regular OLS selection rule (II.3), and a test is performed to see if taking out one of the other atoms in the support will lower the residual compared to the beginning of the iteration (i.e. if the least-contributing column (II.7) is different than the one just added). If yes, a column is removed from the selected support, otherwise, the selection set length is enlarged by 1.

#### C. IOLSR and OLSR properties

Both methods require a single  $\mathcal{O}(\text{Mul})$  operation at each loop iteration, similar to OMP, where  $\mathcal{O}(\text{Mul})$  is the complexity of multiplication by  $\mathbf{D}$ . In general  $\mathcal{O}(\text{Mul}) = \mathcal{O}(mn)$  but for specific dictionaries it might be lower, e.g., for DCT it is  $\mathcal{O}(n \log n)$ . To facilitate this, we introduce two length  $n$  vectors with an efficient update scheme that requires a single dictionary times vector multiplication. To formulate this update scheme recall Lemma II.3. We calculate  $\langle \overrightarrow{R_T\mathbf{d}_{(j)}}, \mathbf{y} \rangle^2$  for each of the atoms in the dictionary by dividing the square of the elements of  $\mathbf{c}(j) = \langle R_T\mathbf{d}_{(j)}, \mathbf{y} \rangle$  by the elements of  $\boldsymbol{\rho}(j) = \|R_T\mathbf{d}_{(j)}\|^2$ . Denote  $\mathbf{v} = R_T\mathbf{d}_{(i)}$ ,  $\tilde{\boldsymbol{\rho}} = \mathbf{D}^T\mathbf{v}$ , where  $i$  is an index of the atom that was currently added to  $T$ , and  $T' = T \cup i$ . Then upon the addition of  $\mathbf{d}_{(i)}$  to the support,  $\mathbf{c}$  and  $\boldsymbol{\rho}$  are updated as follows:

$$\begin{aligned} \mathbf{c}_n(j) &= \langle R_{T'}\mathbf{d}_{(j)}, \mathbf{y} \rangle = \langle (R_T - \tilde{\mathbf{v}}\tilde{\mathbf{v}}^T)\mathbf{d}_{(j)}, \mathbf{y} \rangle = \langle R_T\mathbf{d}_{(j)}, \mathbf{y} \rangle - \langle \tilde{\mathbf{v}}, \mathbf{d}_{(j)} \rangle \langle \tilde{\mathbf{v}}, \mathbf{y} \rangle \\ &= \mathbf{c}_{n-1}(j) - \tilde{\boldsymbol{\rho}}(j)\mathbf{c}_{n-1}(i) / \|\mathbf{v}\|^2 \end{aligned} \quad (\text{III.1})$$

$$\begin{aligned} \boldsymbol{\rho}_n(j) &= \|R_{T'}\mathbf{d}_{(j)}\|^2 = \langle (R_T - \tilde{\mathbf{v}}\tilde{\mathbf{v}}^T)\mathbf{d}_{(j)}, \mathbf{d}_{(j)} \rangle = \langle R_T\mathbf{d}_{(j)}, \mathbf{d}_{(j)} \rangle - \langle \tilde{\mathbf{v}}, \mathbf{d}_{(j)} \rangle^2 \\ &= \boldsymbol{\rho}_{n-1}(j) - \tilde{\boldsymbol{\rho}}(j)^2 / \|\mathbf{v}\|^2. \end{aligned} \quad (\text{III.2})$$

Removing a column works in a similar way (with addition instead of subtraction in (III.1) and (III.2) as the residual is enlarged as the subspace of  $\mathbf{D}_T$  is reduced). The following two functions perform the above procedure efficiently based on (III.1) and (III.2) and lemmas II.4 and II.5 using the updated value of  $\mathbf{B}$  and  $\hat{\mathbf{x}}$  from the previous iteration (denote  $\tau = |T|$ , we always insert to the last index):

TABLE I  
THE OLSR AND IOLSR ALGORITHMS

OLS with replacement (OLSR)	Iterative OLS with replacement (IOLSR)
Input: dictionary $\mathbf{D}$ , measurement $\mathbf{y}$ , cardinality $k$ (or target residual $\epsilon_t$ -IOLSR only) Output: $\hat{\mathbf{x}}$ with $k$ elements (or $\epsilon_t$ residual -IOLSR only), and $T$ its support	
$\{T, \mathbf{c}, \boldsymbol{\rho}, \epsilon_0\} \leftarrow \text{OLS}(\mathbf{D}, \mathbf{y}, k+1)$ <b>if</b> $\epsilon_0 \approx 0$ , <b>return</b> $T, \hat{\mathbf{x}} = \mathbf{D}_T^\dagger \mathbf{y}$ , and <b>exit</b> $\mathbf{c}^0 \leftarrow \mathbf{D}^T \mathbf{y}$ $\mathbf{B} = (\mathbf{D}_T^T \mathbf{D}_T)^{-1}$ $\hat{\mathbf{x}} \leftarrow \mathbf{B} \mathbf{c}_T^0$ <b>loop</b> $j \leftarrow \arg \min_{j \in T} \{\hat{\mathbf{x}}(j)^2 / \mathbf{B}(j, j)\}$ $\{\mathbf{c}, \boldsymbol{\rho}\} \leftarrow \text{updRem}(\mathbf{D}, \mathbf{B}, T, \hat{\mathbf{x}}, \mathbf{c}, \boldsymbol{\rho}, j)$ remove $\mathbf{d}_{(j)}$ from $\mathbf{B}$ using (II.2) $T \leftarrow T \setminus \{j\}$ $i \leftarrow \arg \max_{i \notin T} \{\mathbf{c}(i)^2 / \boldsymbol{\rho}(i)\}$ <b>exit loop if</b> $\mathbf{c}(i)^2 / \boldsymbol{\rho}(i) \geq \hat{\mathbf{x}}(j)^2 / \mathbf{B}(j, j)$ update $\mathbf{B}$ with $\mathbf{d}_{(i)}$ using (II.1) $T \leftarrow T \cup \{i\}$ $\hat{\mathbf{x}} \leftarrow \mathbf{B} \mathbf{c}_T^0$ $\{\mathbf{c}, \boldsymbol{\rho}\} \leftarrow \text{updAdd}(\mathbf{D}, \mathbf{B}, T, \hat{\mathbf{x}}, \mathbf{c}, \boldsymbol{\rho})$ <b>end loop</b> <b>return</b> $T, \hat{\mathbf{x}}$	$T \leftarrow \{\}, \epsilon_0 \leftarrow \ \mathbf{y}\ ^2, \boldsymbol{\rho} \leftarrow \mathbf{1}_{N \times 1}$ $\mathbf{c}^0 \leftarrow \mathbf{D}^T \mathbf{y}, \mathbf{c} \leftarrow \mathbf{c}^0$ <b>while</b> $( T  < k+1)$ or $(\epsilon_0 < \epsilon_t \text{ and } j =  T )$ $i \leftarrow \arg \max_{i \notin T} \{\mathbf{c}(i)^2 / \boldsymbol{\rho}(i)\}$ update $\mathbf{B}$ with $\mathbf{d}_{(i)}$ according to (II.1) $T \leftarrow T \cup \{i\}$ $\hat{\mathbf{x}} \leftarrow \mathbf{B} \mathbf{c}_T^0$ $\epsilon_0 \leftarrow \epsilon_0 - \hat{\mathbf{x}}( T )^2 / \mathbf{B}( T ,  T )$ $\{\mathbf{c}, \boldsymbol{\rho}\} \leftarrow \text{updAdd}(\mathbf{D}, \mathbf{B}, T, \hat{\mathbf{x}}, \mathbf{c}, \boldsymbol{\rho})$ $j \leftarrow \arg \min_j \{\hat{\mathbf{x}}(j)^2 / \mathbf{B}(j, j)\}$ <b>if</b> $j \neq  T $ $\{\mathbf{c}, \boldsymbol{\rho}\} \leftarrow \text{updRem}(\mathbf{D}, \mathbf{B}, T, \hat{\mathbf{x}}, \mathbf{c}, \boldsymbol{\rho}, j)$ $\epsilon_0 \leftarrow \epsilon_0 + \hat{\mathbf{x}}(j)^2 / \mathbf{B}(j, j)$ remove $\mathbf{d}_{(j)}$ from $\mathbf{B}$ using (II.2) $T \leftarrow T \setminus \{j\}$ <b>end if</b> <b>end while</b> <b>if</b> $( T  > k)$ or $(\epsilon_0 + \hat{\mathbf{x}}(j)^2 / \mathbf{B}(j, j) < \epsilon_t)$ perform: $T \leftarrow T \setminus \{j\}$ <b>return</b> $T, \hat{\mathbf{x}} = \mathbf{D}_T^\dagger \mathbf{y}$

$$\begin{aligned} \{\mathbf{c}, \boldsymbol{\rho}\} &= \text{updAdd}(\mathbf{D}, \mathbf{B}, T, \hat{\mathbf{x}}, \mathbf{c}, \boldsymbol{\rho}) \\ \mathbf{v} &\leftarrow \mathbf{D}_T \mathbf{B}(:, \tau) \\ \tilde{\boldsymbol{\rho}} &\leftarrow \mathbf{D}^T \mathbf{v} \\ \mathbf{c} &\leftarrow \mathbf{c} - (\hat{\mathbf{x}}(k) / \mathbf{B}(\tau, \tau)) \tilde{\boldsymbol{\rho}} \\ \boldsymbol{\rho} &\leftarrow \boldsymbol{\rho} - (1 / \mathbf{B}(\tau, \tau)) \tilde{\boldsymbol{\rho}} \odot \tilde{\boldsymbol{\rho}} \end{aligned}$$

$$\begin{aligned} \{\mathbf{c}, \boldsymbol{\rho}\} &= \text{updRem}(\mathbf{D}, \mathbf{B}, T, \hat{\mathbf{x}}, \mathbf{c}, \boldsymbol{\rho}, j) \\ \mathbf{v} &\leftarrow \mathbf{D}_T \mathbf{B}(:, j) \\ \tilde{\boldsymbol{\rho}} &\leftarrow \mathbf{D}^T \mathbf{v} \\ \mathbf{c} &\leftarrow \mathbf{c} + (\hat{\mathbf{x}}(j) / \mathbf{B}(j, j)) \tilde{\boldsymbol{\rho}} \\ \boldsymbol{\rho} &\leftarrow \boldsymbol{\rho} + (1 / \mathbf{B}(j, j)) \tilde{\boldsymbol{\rho}} \odot \tilde{\boldsymbol{\rho}} \end{aligned}$$

IOLSR and OLSR share the same performance guarantees in theorems IV.6 and IV.4, and we provide convergence speed analysis for OLSR in Theorem IV.5. Note that in simulations IOLSR demonstrates linear dependency on  $k$  with a similar slope as OMP and OLSR. Another difference worth noting between IOLSR and OLSR, as well as other methods that enable backtracking such as subspace pursuit [8] or CoSaMP [19] is that IOLSR is able to run with a target residual norm (designated  $\epsilon_t$ ) similarly to OMP and OLS, rather than target sparsity. Target residual as a stopping condition is commonly used in many applications. Usually, IOLSR supplies a better approximation than OLSR, this is because of the fact that the IOLSR subspace is sequentially optimized (i.e. the back tracing is performed on subspaces of increasing sizes rather than on a subspace of constant size), in comparison to OLSR and other methods that begin with a given support. Both algorithms are described in table I, and MATLAB code is provided in the supplementaries.

#### IV. THEORETICAL PERFORMANCE GUARANTEES

The OLSR and IOLSR performance guarantees will be presented here using a series of lemmas and theorems. To this end the following notation will be used.  $L$  is the true support s.t.  $\mathbf{x}(j) \neq 0, \forall j \in L$ , and 0 otherwise.  $T$  is a set of  $k$  indices that represent the chosen support of an algorithm at a specific iteration.  $\tilde{L} \triangleq L \setminus T$ ,  $\tilde{T} \triangleq T \setminus L$ .  $\tilde{\mathbf{x}}(i) = \mathbf{x}(i) \forall i \in \tilde{L}$ , and 0 otherwise,  $\tilde{\mathbf{y}} = \mathbf{D} \tilde{\mathbf{x}}$ .  $\kappa$  the number of atoms in the true support that are not yet identified ( $\kappa = |\tilde{L}| = \|\tilde{\mathbf{x}}\|_0$ ).

In the proofs contained in this section we rely on the following lemma from [17].

**Lemma IV.1** (Lemma 2.1 in [17]). *Given  $\mathbf{x}_1, \mathbf{x}_2$ , such that  $\|\mathbf{x}_1\|_0 + \|\mathbf{x}_2\|_0 \leq 2k$ ,  $\mathbf{x}_1 \perp \mathbf{x}_2$ , and a dictionary  $\mathbf{D}$  with a RIP constant  $\delta < 1$  of order  $2k$ , then:*

$$|\cos \angle(\mathbf{D} \mathbf{x}_1, \mathbf{D} \mathbf{x}_2)| \leq \delta. \quad (\text{IV.1})$$

By using Lemma IV.1 and by noting that the support of  $\tilde{\mathbf{y}}$  and  $T$  are disjoint by definition, we have,

$$\|P_A \tilde{\mathbf{y}}\|^2 = \langle P_A \tilde{\mathbf{y}}, \tilde{\mathbf{y}} \rangle = \|P_A \tilde{\mathbf{y}}\| \|\tilde{\mathbf{y}}\| \cos \angle(P_A \tilde{\mathbf{y}}, \tilde{\mathbf{y}}) \leq \|P_A \tilde{\mathbf{y}}\| \|\tilde{\mathbf{y}}\| \delta.$$

Hence

$$\|P_A \tilde{\mathbf{y}}\|^2 \leq \delta^2 \|\tilde{\mathbf{y}}\|^2. \quad (\text{IV.2})$$

Noting that  $\|P_A \tilde{\mathbf{y}}\|^2 + \|R_A \tilde{\mathbf{y}}\|^2 = \|\tilde{\mathbf{y}}\|^2$  we get

$$\|R_A \tilde{\mathbf{y}}\|^2 \geq (1 - \delta^2) \|\tilde{\mathbf{y}}\|^2. \quad (\text{IV.3})$$

**Lemma IV.2.** *The square of the maximum cosine of the angle between  $R_A \mathbf{d}_{(i)}$ , where  $\mathbf{d}_{(i)}$  is an atom in  $\tilde{L}$ , and  $R_A \tilde{\mathbf{y}}$  obeys*

$$\max_{i \in \tilde{L}} (\cos \angle(R_A \mathbf{d}_{(i)}, R_A \tilde{\mathbf{y}}))^2 \geq \frac{1}{\kappa} c_\delta = \frac{1}{\kappa} (1 - \delta^2)(1 - \delta). \quad (\text{IV.4})$$

*Proof.* We prove by contradiction, inspired by the proof of Theorem 2.2 in [17]. Consider  $\|R_A \tilde{\mathbf{y}}\|_2$ :

$$\|R_A \tilde{\mathbf{y}}\|_2 = \frac{|\langle R_A \tilde{\mathbf{y}}, R_A \tilde{\mathbf{y}} \rangle|}{\|R_A \tilde{\mathbf{y}}\|_2} = \frac{\langle \sum_{i \in \tilde{L}} x(i) R_A \mathbf{d}_{(i)}, R_A \tilde{\mathbf{y}} \rangle}{\|R_A \tilde{\mathbf{y}}\|_2} \leq \frac{\sum_{i \in \tilde{L}} |x(i)| |\langle R_A \mathbf{d}_{(i)}, R_A \tilde{\mathbf{y}} \rangle|}{\|R_A \tilde{\mathbf{y}}\|_2}, \quad (\text{IV.5})$$

which leads to:

$$\|R_A \tilde{\mathbf{y}}\|_2 \stackrel{(a)}{\leq} \sum_{i \in \tilde{L}} |x(i)| \cos \alpha_i \stackrel{(b)}{<} \sqrt{\frac{1}{\kappa}} c_\delta \|\tilde{\mathbf{x}}\|_1 \stackrel{(c)}{\leq} \sqrt{c_\delta} \|\tilde{\mathbf{x}}\|_2, \quad (\text{IV.6})$$

where we define  $\alpha_i = \angle(R_A \mathbf{d}_{(i)}, R_A \tilde{\mathbf{y}})$ , and use the fact that  $\|R_A \mathbf{d}_{(i)}\|^2 \leq \|\mathbf{d}_{(i)}\|^2 = 1$  in transition (a). For step (b) we assume that for some positive constant  $c_\delta$ ,  $|\cos \alpha_i| < \sqrt{c_\delta/\kappa}$ , for all  $i \in \tilde{L}$ . For step (c) we use the relationship  $\|\tilde{\mathbf{x}}\|_1 \leq \sqrt{\kappa} \|\tilde{\mathbf{x}}\|_2$ .

On the other hand, By (IV.3) and the RIP we have:

$$\|R_A \tilde{\mathbf{y}}\| \geq \sqrt{1 - \delta^2} \|\tilde{\mathbf{y}}\|_2 \geq \sqrt{1 - \delta^2} \sqrt{1 - \delta} \|\tilde{\mathbf{x}}\|_2. \quad (\text{IV.7})$$

Combining (IV.7) and (IV.6) we get that in order to produce a contradiction we need to set  $\sqrt{1 - \delta^2} \sqrt{1 - \delta} = \sqrt{c_\delta}$ .  $\square$

The following Lemma gives an upper bound for the contribution to the residual of at least one of the erroneously selected atoms in the chosen support. We will later bound from below the contribution to the residual of a member of  $\tilde{L}$ .

**Lemma IV.3.** *For a given support  $T$  with  $\kappa$  atoms that are not in the true support of  $\mathbf{x}$ , there exists an atom  $\mathbf{d}_{(j)}$ ,  $j \in \tilde{T}$  such that:*

$$\frac{1}{\kappa} \frac{\delta^2}{1 - \delta} \|\tilde{\mathbf{y}}\|^2 \geq \langle \overrightarrow{R_T \mathbf{d}_{(j)}}, \mathbf{y} \rangle^2. \quad (\text{IV.8})$$

*Proof.* We start with a bound on the norm of  $P_T \tilde{\mathbf{y}}$  from below as follows:

$$\begin{aligned} \|P_T \tilde{\mathbf{y}}\|^2 &\stackrel{(a)}{\geq} (1 - \delta) \|\hat{\tilde{\mathbf{x}}}\|^2 \stackrel{(b)}{\geq} (1 - \delta) \sum_{j \in T} \frac{1}{\|R_{T \setminus j} \mathbf{d}_{(j)}\|^4} \langle R_{T \setminus j} \mathbf{d}_{(j)}, \tilde{\mathbf{y}} \rangle^2 \\ &\stackrel{(c)}{\geq} (1 - \delta) \sum_{j \in \tilde{T}} \frac{1}{\|R_{T \setminus j} \mathbf{d}_{(j)}\|^2} \langle R_{T \setminus j} \mathbf{d}_{(j)}, \mathbf{y} \rangle^2 = (1 - \delta) \sum_{j \in \tilde{T}} \langle \overrightarrow{R_{T \setminus j} \mathbf{d}_{(j)}}, \mathbf{y} \rangle^2, \end{aligned}$$

where (a) follows from the RIP, (b) follows from (II.6), and (c) is because we sum on fewer atoms  $|\tilde{T}| < |T|$ , reduce the power in the denominator from 4 to 2, and replace  $\tilde{\mathbf{y}}$  with  $\mathbf{y}$  since  $R_{T \setminus j} \perp \mathbf{d}_{(j)} \quad \forall j \in T \setminus \tilde{T}$ . The assertion in Lemma IV.3 now follows by noting that the smallest value of  $\langle \overrightarrow{R_{T \setminus j} \mathbf{d}_{(j)}}, \mathbf{y} \rangle^2$  for  $j \in \tilde{T}$  is no bigger than the average on  $\kappa$  atoms in  $\tilde{T}$ , and by noting that  $\delta^2 \|\tilde{\mathbf{y}}\|^2 \geq \|P_A \tilde{\mathbf{y}}\|^2$  from (IV.2).  $\square$

**Theorem IV.4** (RIP bound). *Given  $\mathbf{y} = \mathbf{D}\mathbf{x}$ , where  $\|\mathbf{x}\|_0 = k$  and  $\mathbf{D}$  satisfies the RIP of order  $2k$  with a constant  $\delta_{2k} \leq 0.445$  the OLSR and IOLSR algorithms yield perfect support reconstruction,  $T = L$ .*

*Proof of Theorem IV.4.* Assume that the algorithms converged to a support  $T$ , with  $\tilde{T} \neq \emptyset$ . Using Lemma IV.2 we get that there exists  $i \in \tilde{L} \neq \emptyset$  such that:

$$\langle \overrightarrow{R_T \mathbf{d}_{(i)}}, \mathbf{y} \rangle^2 \stackrel{(a)}{\geq} \langle \overrightarrow{R_T \mathbf{d}_{(i)}}, R_T \tilde{\mathbf{y}} \rangle^2 = \|R_T \tilde{\mathbf{y}}\|^2 |\cos \alpha_i|^2 \geq \|R_T \tilde{\mathbf{y}}\|^2 \frac{1}{\kappa} c_\delta, \quad (\text{IV.9})$$

where in transition (a) we used  $R_T \mathbf{y} = R_T \tilde{\mathbf{y}}$  since  $(\mathbf{y} - \tilde{\mathbf{y}}) \in \text{span}\{\mathbf{D}_T\}$ .

The OLSR algorithm stops when there are no atoms that can replace the one that has been extracted (the condition for replacement is that the new residual is lower than the previous one). From Lemma IV.3 we have that the last atom extracted satisfies (IV.8). We combine this with (IV.9) and (IV.3) that yield  $\|R_T \tilde{\mathbf{y}}\|^2 \frac{1}{\kappa} c_\delta \geq \|\tilde{\mathbf{y}}\|^2 \frac{1}{\kappa} (1 - \delta^2) c_\delta$ . Therefore we have that in order to contradict the assumption that the stopping criterion was reached with  $\mathbf{d}_{(i)}$  out of the selected support the following needs to hold:

$$\|\tilde{\mathbf{y}}\|^2 \frac{1}{\kappa} (1 - \delta^2) c_\delta \geq \frac{1}{\kappa} \frac{\delta^2}{1 - \delta} \|\tilde{\mathbf{y}}\|^2.$$

Substituting  $c_\delta$  from (IV.4) and reorganizing the terms we get the condition in Theorem IV.4.

IOLSR differs from OLSR in that the new candidate is first inserted into the support, and then an atom for elimination is selected (in OLSR it is the other way around, first an atom is removed from the support, and then a new one that improves the residual is inserted if it exists). The resulting bound is identical. The way to show it is similar except for a few changes. To this end denote by  $T$  the selected set at the beginning of the current iteration, set  $\tilde{\mathbf{y}}' = \tilde{\mathbf{y}} - \mathbf{d}_{(i)} \mathbf{x}(i)$  and note that:

$$\|R_T \tilde{\mathbf{y}}\|^2 \geq \|R_T \tilde{\mathbf{y}}'\|^2 - \langle \overrightarrow{R_T \mathbf{d}_{(i)}}, \tilde{\mathbf{y}}' \rangle^2 \stackrel{(a)}{\geq} \|R_{T \cup i} \tilde{\mathbf{y}}'\|^2 \stackrel{(b)}{\geq} (1 - \delta^2) \|\tilde{\mathbf{y}}'\|^2, \quad (\text{IV.10})$$

where (II.3) and (IV.3) were used in transitions (a) and (b) respectively. Now, combine (IV.9) with (IV.10) to get the required bound with  $\tilde{\mathbf{y}}'$  instead of  $\tilde{\mathbf{y}}$ .  $\square$

**Theorem IV.5.** *Let  $\gamma = c_\delta - \delta^2/c_\delta$ , with  $c_\delta$  defined by (IV.4). The number of iteration needed for the OLSR algorithm to converge to a solution with residual lower than  $\epsilon_t$  is bounded from above by*

$$b = k \left( 1 + \frac{1}{\gamma} \ln \left( \frac{\|\mathbf{y}\|^2}{\epsilon_t e^{c_\delta}} \right) \right). \quad (\text{IV.11})$$

Note that this is a worst case analysis. In practice, only several iterations are needed for convergence in addition to the OLS number of iterations that we use to get the initial  $k$ -sparse solution. To get a feeling of the bound, assume  $\delta = 1/3$  and  $\epsilon_t = 10^{-3}\|\mathbf{y}\|^2$ , this yields:

$$\gamma \sim 0.44, \quad b \sim 16k,$$

meaning that less than  $16k$  iterations will be needed for convergence, regardless of the value of  $k$ .

*Proof.* Consider the OLS algorithm, and recall (II.3) in Lemma II.3,

$$\left\langle \overrightarrow{R_{T_n} \mathbf{a}}, \mathbf{y} \right\rangle^2 = (\mathbf{y}^T \overrightarrow{R_{T_n} \mathbf{a}})^2 = \|R_{T_n} \mathbf{y}\|^2 - \|R_{T_{n+1}} \mathbf{y}\|^2 = \|P_{T_{n+1}} \mathbf{y}\|^2 - \|P_{T_n} \mathbf{y}\|^2,$$

where the subscripts  $n$  and  $n+1$  represent the set  $T$  in subsequent loop iterations. By using Lemma IV.2 we have  $\left\langle \overrightarrow{R_{T_{n+1}} \mathbf{a}}, \mathbf{y} \right\rangle^2 \geq \|R_{T_n} \mathbf{y}\|^2 \frac{c_\delta}{\kappa}$ , yielding the following relationship between the residuals in subsequent loop iterations of OLS:

$$\|R_{T_{n+1}} \mathbf{y}\|^2 \leq \left( 1 - \frac{c_\delta}{\kappa} \right) \|R_{T_n} \mathbf{y}\|^2.$$

We can replace  $1/\kappa$  with  $1/k$  to get a more restrictive bound  $(1 - c_\delta/k)^k \|\mathbf{y}\|^2 \geq \|R_{T_{n=k}} \mathbf{y}\|^2$ . For large enough  $k$  this converges to:  $e^{-c_\delta} \|\mathbf{y}\|^2 \geq \|R_{T_{n=k}} \mathbf{y}\|^2$ . Turning now to the replacement part of the OLSR algorithm, assume that  $tk$  ( $t > 0$ ) additional iterations were performed in this stage. From Lemma IV.3 we have that at each step we take out of the support an atom with at most  $\frac{\delta^2}{\kappa(1-\delta)} \|\tilde{\mathbf{y}}\|^2$  contribution to the residual, which can be further bounded by using (IV.3), yielding  $\frac{\delta^2}{\kappa c_\delta} \|R_{T_n} \tilde{\mathbf{y}}\|^2$ . Hence we get that at iteration  $tk$

$$\|R_{T_{n=tk+k}} \mathbf{y}\|^2 \leq \left( 1 - \frac{\gamma}{k} \right)^{tk} \|R_{T_{n=k}} \mathbf{y}\|^2 \leq \left( 1 - \frac{\gamma}{k} \right)^{tk} e^{-c_\delta} \|\mathbf{y}\|^2,$$

with  $\gamma = c_\delta - \delta^2/c_\delta$  (note that if  $\delta < 0.445$ ,  $\gamma > 0$ ). Thus for the estimation error to be less than  $\epsilon_t$  we need  $t \geq \gamma^{-1} \ln(\|\mathbf{y}\|^2/\epsilon_t e^{c_\delta})$ , with  $(1+t)k$  as an upper bound on the number of iterations.  $\square$

The remainder of this section deals with the case that the measurements are corrupted by an additive white Gaussian noise (AWGN). To this end denote the noise variance by  $\sigma^2$ , and recall that  $\mathbf{y}_0 = \mathbf{D}\mathbf{x}$ ,  $\mathbf{y} = \mathbf{y}_0 + \mathbf{w}$ ,  $\mathbf{w} \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$ .

**Theorem IV.6.** *Given measurements corrupted by AWGN with variance  $\sigma^2$ , and some parameter  $a \geq 0$ , a perfect support reconstruction is achieved by OLSR and IOLSR with probability exceeding  $1 - (\sqrt{\pi(1+a)} \log nn^a)^{-1}$  if*

$$\frac{\|\mathbf{y}_0\|^2}{\sigma^2 k} \geq \frac{(1 + \sqrt{1-\delta})^2 (2(1+a) \log n)}{((1-\delta)(1-\delta^2) - \delta)^2}.$$

The proof of Theorem IV.6 is relegated to the supplementary material.

The following is a consequence of Theorem IV.6:

**Corollary IV.6.1.** *The root mean squared error of the estimation of  $\mathbf{x}$  by the OLSR and IOLSR algorithms is bounded, with probability exceeding  $1 - (\sqrt{\pi(1+a)} \log nn^a)^{-1}$ , by*

$$\|\mathbf{x} - \hat{\mathbf{x}}\| \leq \sqrt{k} \sigma \left( \frac{1}{\sqrt{1-\delta^2}} + \frac{(1+\delta)(1+\sqrt{1-\delta})}{\sqrt{1-\delta}((1-\delta)(1-\delta^2) - \delta)} \sqrt{2(1+a) \log n} \right). \quad (\text{IV.12})$$

The proof of corollary IV.6.1 is relegated to the supplementary material. Its implication is that the IOLSR and OLSR error is proportional, up to a constant times  $\log n$ , to the error ( $\sqrt{k}\sigma$ ) of an oracle estimator that knows the true support. These bounds are similar to the ones achieved for other methods including SP and CoSaMP [6, 12].

## V. SIMULATIONS

In this section we demonstrate the performance of our algorithms numerically using MATLAB routines that are available in the supplementary material.

The results appear in Figure 1. There are three plots: (a) A phase transition diagram, following the methodology of [9]. We fix  $m$  and variate  $n$  and  $k$  according to two auxiliary variables. We use a threshold of  $\|\mathbf{x} - \hat{\mathbf{x}}\|^2/\|\mathbf{x}\|^2 \leq 10^{-4}$  and plot the resulting curve.  $\mathbf{D}$  is normalized Gaussian matrix, and  $\mathbf{x}$  is a Gaussian signal. Results are averaged over 50 realizations. For BP

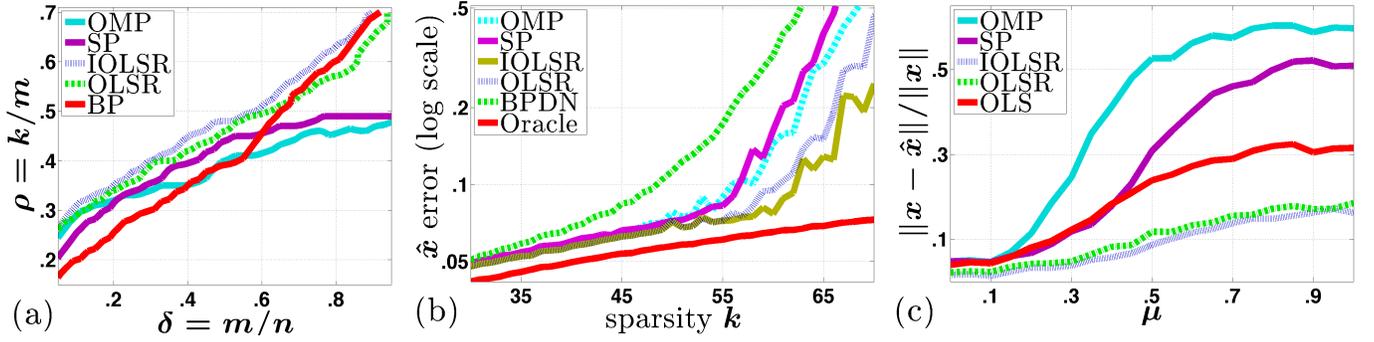


Fig. 1. (a) Phase transition diagram for  $m = 400$ , areas below the lines of each method represent an error threshold of  $\|\hat{x} - x\|^2 / \|x\|^2 \leq 10^{-4}$ . (b) Error ( $\|\hat{x} - x\|$ ) vs. sparsity in presence of noise  $\sigma = 0.1 \|\mathbf{y}_0\| / \sqrt{m}$ ,  $n = 600$ ,  $m = 200$  compared to an oracle that knows the true support. (c) Error as a function of a coherency damaging parameter  $\mu$ ;  $n = 300$ ,  $m = 100$ ,  $k = 30$ .

[7] we used the CVX package [13] to solve  $\min \|\mathbf{x}\|_1$  s.t.  $\mathbf{y} = \mathbf{D}\mathbf{x}$ , and then to improve its recovery we extract the maximal  $k$  values from the result and compute a new approximation for the sparse representation by least squares. (b) An experiment with noise, where we fix the size of the dictionary, corrupt the measurements by AWGN with  $\sigma = 0.1 \cdot \|\mathbf{y}_0\| / \sqrt{m}$ , and plot the reconstruction error vs. the sparsity. Results are averaged over 1000 instances of  $\mathbf{D}$  and  $\mathbf{x}$  drawn from a Gaussian distribution. For BPDN [7] we used the TFOCS package [1] to solve  $\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{y} - \mathbf{D}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1$  with  $\lambda = \sqrt{2 \log n} \sigma$  following [7]. (c) This plot demonstrates resilience to correlation in the dictionary. For a fixed  $n, m, k$  we generate  $\mathbf{D}$  from a Gaussian distribution. We then increase the correlation in the dictionary by performing  $\mathbf{d}_{(i)} = \mathbf{d}_{(i)} + \mu \mathbf{d}_{(i+1)}$  for each atom, repeating the process for 5 times to increase the effect. Results are averaged over 1000 realizations.

## VI. CONCLUDING REMARKS AND FUTURE WORK

We have presented here two new pursuits for sparse signal recovery, along with theoretical analysis of their properties. We maintain that IOLSR is a good alternative to OMP and OLS in applications, where the target of the sparse approximation is the magnitude of the residual (stemming perhaps from knowledge of noise variance), and where optimizing various method parameters may be hard. Both IOLSR and OLSR are simple to implement and demonstrate little overhead in regards to OMP. We posit that this kind of fast-converging approach can lead to interesting developments in applications such as dictionary learning where optimization is often carried one atom at a time, and in applications suffering from highly coherent dictionaries, as the correlation of an atom to the selected support is built into the algorithms (in the form of the vector  $\rho$ ) and does not require additional computations.

## REFERENCES

- [1] S. Becker, E. Candes, and M. Grant. TFOCS: Matlab software, version 1.4. <http://cvxr.com/tfocs>, October 2014.
- [2] T. Blumensath and M.E. Davies. On the difference between orthogonal matching pursuit and orthogonal least squares. 2007.
- [3] A.M. Bruckstein, D.L. Donoho, and M. Elad. From sparse solutions of systems of equations to sparse modeling of signals and images. *SIAM review*, 51(1):34–81, 2009.
- [4] T. Cai and A. Zhang. Sharp RIP bound for sparse signal and low-rank matrix recovery. *Applied and Computational Harmonic Analysis*, 35(1):74–93, 2013.
- [5] E. Candes and T. Tao. Decoding by linear programming. *IEEE Transactions on Information Theory*, 51(12):4203–4215, 2005.
- [6] E. Candes and T. Tao. The dantzig selector: Statistical estimation when  $p$  is much larger than  $n$ . *The Annals of Statistics*, 35(6):2313–2351, 2007.
- [7] S.S. Chen, D.L. Donoho, and M.A. Saunders. Atomic decomposition by basis pursuit. *SIAM review*, 43(1):129–159, 2001.
- [8] W. Dai and O. Milenkovic. Subspace pursuit for compressive sensing signal reconstruction. *Information Theory, IEEE Transactions on*, 55(5):2230–2249, 2009.
- [9] D. Donoho and J. Tanner. Observed universality of phase transitions in high-dimensional geometry, with implications for modern data analysis and signal processing. *Philosophical Transactions of the Royal Society of London A: Mathematical, Physical and Engineering Sciences*, 367(1906):4273–4293, 2009.
- [10] Y.C. Eldar and G. Kutyniok. *Compressed sensing: theory and applications*. Cambridge University Press, 2012.
- [11] S. Foucart. Sparse recovery algorithms: sufficient conditions in terms of restricted isometry constants. In *Approximation Theory XIII: San Antonio 2010*, pages 65–77. Springer, 2012.
- [12] R. Giryes and M. Elad. RIP-based near-oracle performance guarantees for SP, CoSaMP, and IHT. *IEEE Transactions on Signal Processing*, 60(3):1465–1468, March 2012.
- [13] M. Grant and S. Boyd. CVX: Matlab software for disciplined convex programming, version 2.1. <http://cvxr.com/cvx>, March 2014.
- [14] P. Jain, A. Tewari, and Inderjit S.D. Orthogonal matching pursuit with replacement. In J. Shawe-Taylor, R. S. Zemel, P. L. Bartlett, F. Pereira, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 24*, pages 1215–1223. Curran Associates, Inc., 2011.
- [15] P.A. Janakiraman and S. Renganathan. Recursive computation of pseudo-inverse of matrices. *Automatica*, 18(5):631–633, 1982.
- [16] M. Khan. Updating inverse of a matrix when a column is added. Technical report, UBC, 2008.
- [17] C. Ling-Hua and W. Jwo-Yuh. An improved rip-based performance guarantee for sparse signal recovery via orthogonal matching pursuit. *IEEE Transactions on Information Theory*, 60(9):5702–5715, Sept 2014.
- [18] A. Miller. *Subset selection in regression*. CRC Press, 2002.
- [19] D. Needell and J.A. Tropp. CoSaMP: Iterative signal recovery from incomplete and inaccurate samples. *Applied and Computational Harmonic Analysis*, 26(3):301–321, 2009.
- [20] L. Rebollo-Neira and D. Lowe. Optimized orthogonal matching pursuit approach. *IEEE signal processing Letters*, 9(4):137–140, 2002.
- [21] C.-B. Song, S.-T. Xia, and X.-J. Liu. Improved analysis for subspace pursuit algorithm in terms of restricted isometry constant. *IEEE Signal Processing Letters*, 21(11):1365–1369, 2014.
- [22] R. Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 267–288, 1996.

# Efficient Sparse Recovery Pursuits with Least Squares Supplemental Material

## S1. ALGORITHMS

### A. Orthogonal least squares

Cited in table S1 for completeness, The computationally demanding step in OLS is the calculation of  $\|R_T \mathbf{d}_{(i)}\|$  for each atom in  $\mathbf{D}$  at every loop iteration, which costs  $\mathcal{O}(mnk)$ . In section S1-B we propose a more computationally efficient way to calculate OLS using two auxiliary length- $n$  vectors.

### B. Fast orthogonal least squares

Designated Fast OLS, this method requires only two auxiliary length  $n$  vectors to be stored in memory. One for storing the residual norm of each atom in the current iteration and the other to store the magnitude of the correlation between the measurement residual and each atom. Utilizing the auxiliary vectors, the runtime is reduced to that of regular OMP. As was mentioned earlier, it resembles OOMP [20] in its mechanics, but it requires substantially less memory -  $\mathcal{O}(2n)$  in FOLS vs.  $\mathcal{O}(mn)$  in OOMP. The FOLS algorithm is depicted alongside OLS in table S1. The computational cost of the steps in the loop is  $\mathcal{O}(2n + mn + mk + k^2)$ . Note that there is a single step that requires  $\mathcal{O}(mn)$  flops, similar to OMP. The way FOLS operates is by keeping two vectors:

$$\mathbf{c}(j) = \langle R_T \mathbf{d}_{(j)}, \mathbf{y} \rangle, \quad \rho(j) = \|R_T \mathbf{d}_{(j)}\|^2$$

that are updated using (denote  $\mathbf{v} = R_T \mathbf{d}_{(i)}$ ,  $\tilde{\rho} = \mathbf{D}^T \mathbf{v}$ ,  $T' = T \cup i$ ):

$$\begin{aligned} \mathbf{c}(j) &= \langle R_{T'} \mathbf{d}_{(j)}, \mathbf{y} \rangle = \langle (R_T - \tilde{\mathbf{v}} \tilde{\mathbf{v}}^T) \mathbf{d}_{(j)}, \mathbf{y} \rangle = \langle R_T \mathbf{d}_{(j)}, \mathbf{y} \rangle - \langle \tilde{\mathbf{v}}, \mathbf{d}_{(j)} \rangle \langle \tilde{\mathbf{v}}, \mathbf{y} \rangle \\ &= \mathbf{c}(j) - \frac{\tilde{\rho}(j) \mathbf{c}(i)}{\|\mathbf{v}\|^2} \end{aligned} \tag{S1.1}$$

$$\begin{aligned} \rho(j) &= \|R_{T'} \mathbf{d}_{(j)}\|^2 = \langle (R_T - \tilde{\mathbf{v}} \tilde{\mathbf{v}}^T) \mathbf{d}_{(j)}, \mathbf{d}_{(j)} \rangle = \langle R_T \mathbf{d}_{(j)}, \mathbf{d}_{(j)} \rangle - \langle \tilde{\mathbf{v}}, \mathbf{d}_{(j)} \rangle^2 \\ &= \rho(j) - \frac{\tilde{\rho}(j)^2}{\|\mathbf{v}\|^2}. \end{aligned} \tag{S1.2}$$

## S2. PROOFS

*Proof of Theorem IV.6.* We describe the proof for  $\kappa = k$ , since this is the most restrictive case, it can be easily generalized and checked for the case  $\kappa < k$ . Hence,  $\tilde{\mathbf{y}} = \mathbf{y}$ . Denote for use in this proof:

$$\mathbf{w}_T = \frac{P_T \mathbf{w}}{\|\mathbf{y}_0\|}, \quad w_i = \frac{\langle \overrightarrow{R_T \mathbf{d}_{(i)}}, \mathbf{w} \rangle}{\|\mathbf{y}_0\|}. \tag{S2.1}$$

Consider the proof of Lemma IV.3. By using (S2.1) and noting that

$$\|P_T \mathbf{y}\| = \|P_T \mathbf{y}_0 + P_T \mathbf{w}\| \leq \|P_T \mathbf{y}_0\| + \|P_T \mathbf{w}\| \leq (\delta + \|\mathbf{w}_T\|) \|\mathbf{y}_0\|, \tag{S2.2}$$

Lemma IV.3 takes the form

$$\frac{\delta + \|\mathbf{w}_T\|}{\sqrt{k(1-\delta)}} \|\mathbf{y}_0\| \geq \min_{j \in T} \left| \langle \overrightarrow{R_T \mathbf{d}_{(j)}}, \mathbf{y} \rangle \right|. \tag{S2.3}$$

From Lemma IV.2 we have

$$\left| \langle \overrightarrow{R_T \mathbf{d}_{(i)}}, R_T \mathbf{y}_0 + \mathbf{w} \rangle \right| \geq \|R_T \mathbf{y}_0\| \cos \alpha_i - \|\mathbf{y}_0\| |w_i| \geq \|\mathbf{y}_0\| \left( \sqrt{\frac{c_\delta(1-\delta^2)}{k}} - |w_i| \right), \tag{S2.4}$$

for some  $i \in L$  (the true support). So, to contradict the assumption that the algorithms converged with  $\mathbf{d}_{(i)}$  out of the estimated support, combine (S2.3) and (S2.4) and replace the value for  $c_\delta$

$$(1-\delta)(1-\delta^2) - \delta \geq \|\mathbf{w}_T\| + \sqrt{k(1-\delta)} |w_i|. \tag{S2.5}$$

Note that  $w_i$  and  $\mathbf{w}_T$  reside in orthogonal spaces to each other and that  $\mathbf{w}_T$  is with dimension  $k$  whereas  $w_i$  is of dimension 1. Thus we can use the following confidence interval developed in [6]:

$$\Pr \left( \sup_{\mathbf{a} \in \mathbf{D}} |\langle \mathbf{a}, \mathbf{w} \rangle| > \sigma \sqrt{2(1+a) \log n} \right) \leq \left( \sqrt{\pi(1+a) \log nn^a} \right)^{-1}. \tag{S2.6}$$

We can provide a bound for the right hand side of (S2.5) with high probability:

$$\|\mathbf{w}_T\| + \sqrt{k(1-\delta)} |w_i| \leq \frac{\sqrt{k(1+\sqrt{1-\delta})} \sigma \sqrt{2(1+a) \log N}}{\|\mathbf{y}_0\|}. \tag{S2.7}$$

TABLE S1  
THE OLS AND FOLS ALGORITHMS

Orthogonal least squares (OLS)	Fast orthogonal least squares (FOLS)
Input: dictionary $\mathbf{D}$ , measurement $\mathbf{y}$ , either target cardinality $k$ or target residual norm $\epsilon_t$ Output: $\hat{\mathbf{x}}$ with $k$ elements or $\epsilon_t$ residual, and $T$ its support	
$T \leftarrow \{\}$ , $\epsilon_0 \leftarrow \ \mathbf{y}\ ^2$ , $\mathbf{r} \leftarrow \mathbf{y}$ <b>while</b> $ T  < k$ or $\epsilon_0 < \epsilon_t$ $i \leftarrow \arg \max_i \left\{ \langle \mathbf{r}, \mathbf{d}_{(i)} \rangle^2 / \ R_T \mathbf{d}_{(i)}\ ^2 \right\}$ $T \leftarrow T \cup \{i\}$ $\epsilon_0 \leftarrow \epsilon_0 - \langle \mathbf{r}, \mathbf{d}_{(i)} \rangle^2 / \ R_T \mathbf{d}_{(i)}\ ^2$ $\mathbf{r} \leftarrow R_T \mathbf{y}$ <b>end while</b> <b>return</b> $T$ , $\hat{\mathbf{x}} = \mathbf{D}_T^\dagger \mathbf{y}$	$T \leftarrow \{\}$ , $\epsilon_0 \leftarrow \ \mathbf{y}\ ^2$ , $\boldsymbol{\rho} \leftarrow \mathbf{1}_{n \times 1}$ , $\mathbf{c} \leftarrow \mathbf{D}^T \mathbf{y}$ <b>while</b> $ T  < k$ or $\epsilon_0 < \epsilon_t$ $i \leftarrow \arg \max_{i \notin T} \mathbf{c}(i)^2 / \boldsymbol{\rho}(i)$ $T \leftarrow T \cup \{i\}$ $\epsilon_0 \leftarrow \epsilon_0 - \mathbf{c}(i)^2 / \boldsymbol{\rho}(i)$ $\mathbf{v} \leftarrow R_T \mathbf{d}_{(i)}$ $\tilde{\boldsymbol{\rho}} \leftarrow \mathbf{D}^T \mathbf{v}$ $\mathbf{c} \leftarrow \mathbf{c} - \frac{\mathbf{c}(i)}{\ \mathbf{v}\ ^2} \tilde{\boldsymbol{\rho}}$ $\boldsymbol{\rho} \leftarrow \boldsymbol{\rho} - \frac{1}{\ \mathbf{v}\ ^2} \tilde{\boldsymbol{\rho}} \odot \tilde{\boldsymbol{\rho}}$ <b>end while</b> <b>return</b> $T$ , $\hat{\mathbf{x}} = \mathbf{D}_T^\dagger \mathbf{y}$

Combine (S2.5) and (S2.7) to conclude the proof.  $\square$

*Proof of Corollary IV.6.1.* The oracle square root of the MSE satisfies:

$$\sqrt{\frac{k}{1-\delta^2}} \sigma \geq \|\mathbf{x} - \hat{\mathbf{x}}_{\text{oracle}}\| \geq \sqrt{k} \sigma \quad (\text{S2.8})$$

This can be seen directly by noting that the oracle MSE satisfies  $\|\mathbf{x} - \hat{\mathbf{x}}_{\text{oracle}}\|^2 = \text{trace}\{(D_L^T D_L)^{-1}\} \sigma^2$ , and combining the results from (IV.3), and Lemma II.1. The worst case MSE on the other hand (assuming none of the vectors belonging to the true support were selected) satisfies (assume  $\mathbf{A} = \mathbf{D}_T$  is the selected support, and that  $\mathbf{A} \hat{\mathbf{x}}_{\text{worst case}} = \hat{\mathbf{y}}_{\text{worst case}}$ ):

$$\begin{aligned} \|\mathbf{x} - \hat{\mathbf{x}}_{\text{worst case}}\| &= \|\mathbf{x} - (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T (\mathbf{y}_0 + \mathbf{w})\| \\ &\leq \|(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}_0\| + \|\mathbf{x}\| + \|(\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{w}\| \\ &\leq \|\hat{\mathbf{x}}_{\text{worst case}}\| + \|\mathbf{x}\| + \sqrt{\frac{k}{1-\delta^2}} \sigma \\ &\leq \frac{1}{\sqrt{1-\delta}} (\|\mathbf{y}_0\| + \|\hat{\mathbf{y}}_0\|) + \sqrt{\frac{k}{1-\delta^2}} \sigma \\ &\leq \frac{1+\delta}{\sqrt{1-\delta}} \|\mathbf{y}_0\| + \sqrt{\frac{k}{1-\delta^2}} \sigma \end{aligned} \quad (\text{S2.9})$$

Combining the condition from Theorem IV.6 and the expression from (S2.8) we get that (S2.9) can be written as:

$$\|\mathbf{x} - \hat{\mathbf{x}}\| \leq \sqrt{\frac{k}{1-\delta^2}} \sigma + \frac{1+\delta}{\sqrt{1-\delta}} \|\mathbf{y}_0\| \cdot \mathbf{1}_{\left\{ \sigma \geq \frac{((1-\delta)(1-\delta^2)-\delta)\|\mathbf{y}_0\|}{\sqrt{k}(1+\sqrt{1-\delta})\sqrt{2(1+a)} \log N} \right\}} \quad (\text{S2.10})$$

where  $\mathbf{1}_{\{\cdot\}}$  is the indicator function. The result is linear in  $\sigma$ , except for a step in a certain threshold. To get (IV.12), approximate (S2.10) as a sum of two slopes, the first for the oracle MSE, and the second to approximate the step at  $\sigma = \frac{((1-\delta)(1-\delta^2)-\delta)\|\mathbf{y}_0\|}{\sqrt{k}(1+\sqrt{1-\delta})\sqrt{2(1+a)} \log N}$   $\square$