



**Yale University**  
**Department of Computer Science**

Greedy Basis Pursuit

Patrick S. Huggins      Steven W. Zucker

YALEU/DCS/TR-1359  
June 2006

# Greedy Basis Pursuit\*

Patrick S. Huggins<sup>†</sup>      Steven W. Zucker<sup>‡</sup>

## Abstract

We introduce Greedy Basis Pursuit (GBP), a new algorithm for computing signal representations using overcomplete dictionaries. GBP is rooted in computational geometry and exploits an equivalence between minimizing the  $\ell^1$ -norm of the representation coefficients and determining the intersection of the signal with the convex hull of the dictionary. GBP unifies the different advantages of previous algorithms: like standard approaches to Basis Pursuit, GBP computes representations that have minimum  $\ell^1$ -norm; like greedy algorithms such as Matching Pursuit, GBP builds up representations, sequentially selecting atoms. We describe the algorithm, demonstrate its performance, and provide code. Experiments show that GBP can provide a fast alternative to standard linear programming approaches to Basis Pursuit.

## 1 Introduction

The problem of computing sparse signal representations using an overcomplete dictionary arises in a wide range of signal processing applications [82, 31, 50], including image [7], audio [63], and video [3] compression and source localization [66]. The goal is to represent a given signal as a linear superposition of a small number of stored signals, called *atoms*, drawn from a larger set, called the *dictionary*. In traditional signal representation methods, such as the DCT or various wavelet transforms, the dictionary is simply a basis: the number of atoms in the dictionary is equal to the dimensionality of the signal space and representation is unique. By contrast, in an overcomplete dictionary the number of atoms is greater than the dimensionality of the signal space and representation is no longer unique; this enables flexibility in representation [67], ‘shiftability’ [88], and the use of multiple bases [57, 92],

---

\*This work is partially supported by DARPA and AFOSR.

<sup>†</sup>Department of Computer Science, Yale University, CT

<sup>‡</sup>Department of Computer Science, Yale University, CT

but it requires a criterion to select from among the (many) possible representations. A natural one is sparsity, by which the representation selected is the one that uses as few atoms as possible.

Computing sparse representations is NP-hard [73, 28], and so several (heuristic) methods have been developed [67, 78, 16, 51]. These methods optimize various measures of sparsity, typically functions of the representation coefficients [61, 60], using, for example, greedy algorithms [67], gradient descent [64], linear programming [18], and global optimization [81]. Currently, the two most popular algorithms are Matching Pursuit [67] and Basis Pursuit [17, 18].

Matching Pursuit (MP) is a greedy algorithm: a signal representation is iteratively built up by selecting the atom that maximally improves the representation at each iteration. While there is no guarantee that MP computes sparse representations, MP is easily implemented, converges quickly, and has good approximation properties [67, 95, 53]. Moreover, one variant of MP, Orthogonal Matching Pursuit (OMP) [78], can be shown to compute nearly sparse representations under some conditions [97].

Basis Pursuit (BP), instead of seeking sparse representations directly, seeks representations that minimize the  $\ell^1$ -norm of the coefficients. By equating signal representation with  $\ell^1$ -norm minimization, BP reduces signal representation to linear programming [17, 18], which can be solved by standard methods [98]. Furthermore, BP can compute sparse solutions in situations where greedy algorithms fail [18]. Recent theoretical work shows that representations computed by BP are guaranteed to be sparse under certain conditions [34, 33, 46].

While applying standard linear programming methods to compute minimum  $\ell^1$ -norm signal representations is natural, such methods were developed with very different problems in mind and may not be ideally suited to the representation problem. For example, if the dictionary is not sparse then the (normally fast) interior point methods advocated for BP [18] can be slow. Furthermore, the design required to produce examples on which greedy algorithms fail yet BP succeeds suggests that a greedy strategy could be successfully applied to minimum  $\ell^1$ -norm representation.

In this article we develop a new method for computing sparse signal representations, which we call Greedy Basis Pursuit (GBP). Like BP, GBP minimizes the  $\ell^1$ -norm of the representation coefficients. However, unlike standard linear programming approaches to BP, GBP proceeds much like MP, building up the representation by iteratively selecting atoms.

While algorithmically similar to MP, GBP differs from MP in two key ways: (1) GBP uses a novel criterion for selecting the next atom in the rep-

resentation. The criterion is based on computational geometry, and effects a search for the intersection between the signal vector and the convex hull of the dictionary. (2) GBP may discard atoms that it has already selected; this is crucial, as it allows GBP to overcome the ‘mistakes’ that MP makes in atom selection when compared to BP [18].

While GBP returns the signal representation with the minimum  $\ell^1$ -norm, and thus GBP enjoys the theoretical benefits of BP, the greedy strategy of GBP leads to computational gains when compared to standard linear programming methods. Experiments show our implementation of GBP to be faster than off-the-shelf linear programming packages on some signal representation problems, particularly high-dimensional problems with very overcomplete dictionaries.

The remainder of this paper is organized as follows. In Section 1.1 we formally state the sparse signal representation problem. In Section 2 we review current approaches to the problem. Section 3 provides the geometric interpretation of Basis Pursuit that underlies GBP. In Section 4 we describe the Greedy Basis Pursuit algorithm. Section 5 present the results of experiments with GBP. We discuss GBP in Section 6 and conclude in Section 7.

## 1.1 Problem Statement

Given a signal  $\mathbf{x}$  and a dictionary  $\mathcal{D}$  we seek a sparse representation of  $\mathbf{x}$ . We assume that  $\mathbf{x}$  consists of  $d$  real valued measurements, that is,  $\mathbf{x} \in \mathbb{R}^d$ , for example, a sound wave sampled at  $d$  points. We assume that  $\mathcal{D}$  consists of  $n$  atoms and is overcomplete, that is,  $\mathcal{D} = \{\psi_i\}_{i=1}^n$  and  $n > d$ , and that the atoms are also  $d$ -dimensional and have unit norm, that is,  $\forall \psi_i \in \mathcal{D}, \psi_i \in \mathbb{R}^d$  and  $\|\psi_i\|_2 = 1$ . A *representation* of  $\mathbf{x}$  is a set of indices  $\mathcal{I}$ , where  $\mathcal{I} \subseteq \{1, \dots, n\}$  is a set of indices into  $\mathcal{D}$ , and a corresponding set of coefficients  $\mathcal{A} = \{\alpha_i\}_{i \in \mathcal{I}}$  such that

$$\mathbf{x} = \sum_{i \in \mathcal{I}} \alpha_i \psi_i \quad (1)$$

A representation is *sparse* if the number of atoms used,  $|\mathcal{I}|$  (here  $|\cdot|$  denotes cardinality), is minimized over all possible representations.

Equivalently, in matrix notation, given a (column) vector  $\mathbf{x} \in \mathbb{R}^d$  corresponding to the signal, and a  $d \times n$  matrix  $\mathbf{D}$  corresponding to the dictionary, where the  $i$ th column of  $\mathbf{D}$  is the atom  $\psi_i$ , the sparse signal representation problem is then to compute a (column) vector  $\alpha \in \mathbb{R}^n$  solving

$$\text{Minimize } \|\alpha\|_0 \quad \text{subject to } \mathbf{D}\alpha = \mathbf{x} \quad (2)$$

where  $\|\alpha\|_0$  is the  $\ell^0$ -norm of  $\alpha$ , defined to be the number of nonzero entries of  $\alpha$ .

BP replaces the  $\ell^0$ -norm with the  $\ell^1$ -norm, seeking representations that minimize  $\sum_{i \in I} |\alpha_i|$ . In matrix form this corresponds to

$$\text{Minimize } \|\alpha\|_1 \quad \text{subject to } \mathbf{D}\alpha = \mathbf{x} \quad (3)$$

## 2 Related work

The design of GBP draws on previous work in sparse signal representation, particularly the contrast between MP and BP, and on ideas from subset selection, which we summarize here. We also highlight some unexplored connections between sparse signal representation and linear programming.

### 2.1 Matching Pursuit

Matching Pursuit (MP) [67] is the prototypical greedy algorithm [20] applied to sparse signal representation. MP is currently the most popular algorithm for computing sparse signal representations using an overcomplete dictionary, and is used in a variety of applications [7, 79, 3]. MP has also spawned several variants [41, 58, 42], including Orthogonal Matching Pursuit (OMP) [78, 29], which itself has several variants [49, 21, 83].

MP computes a signal representation by greedily constructing a sequence of approximations to the signal,  $\tilde{\mathbf{x}}^{(0)}, \tilde{\mathbf{x}}^{(1)}, \tilde{\mathbf{x}}^{(2)}, \dots$ , where each consecutive approximation is closer to the signal. MP begins with an ‘empty’ representation,  $\tilde{\mathbf{x}}^{(0)} = 0$ , and at each iteration augments the current representation by selecting the atom from the dictionary which is closest to the residual,  $\tilde{\mathbf{x}}^{(t+1)} = \tilde{\mathbf{x}}^{(t)} + \alpha^{(t)}\psi^{(t)}$ , where  $\psi^{(t)}$  maximizes  $\langle \psi_i, \mathbf{x} - \tilde{\mathbf{x}} \rangle$  over all  $\psi_i \in \mathcal{D}$ .

MP is easy to implement, has a guaranteed exponential rate of convergence [67, 95, 53], and recovers relatively sparse solutions [97], particularly compared to earlier approaches such as the Method-of-Frames [26, 18].

A fundamental drawback of MP (and its variants) is its inability to compute truly sparse representations. It is possible to construct signal representation problems where, because of its greediness, MP initially selects an atom that is not part of the optimal sparse representation; as a result, many of the subsequent atoms selected by MP simply compensate for the poor initial selection [30, 18]. This shortcoming motivated the development of BP, which succeeds on these problems [18]; recent theoretical work explains this phenomenon [34, 33, 46].

These problems are also motivation for the development of GBP. Here MP fails because of its poor initial selection of atoms; however, the atoms initially selected by MP are not necessarily bad in general, after all, these problems are specially designed for MP to fail on. For MP to succeed on these problems, it would need to either make ‘better’ atom selections or be able to discard ‘bad’ atoms to recover from poor selections (or both). GBP adapts the greedy strategy to incorporate both of these ideas and compute the same representations as BP.

## 2.2 Basis Pursuit

Basis Pursuit (BP) [16, 17, 18] approaches sparse signal representation by changing the problem to one of minimizing the  $\ell^1$ -norm of the representation coefficients. As noted above, this has theoretical advantages over greedy approaches. Algorithmically, BP equates sparse signal representation with linear programming.

A linear program is defined as follows: Given a matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ , a (column) vector  $\mathbf{b} \in \mathbb{R}^m$ , and a (column) vector  $\mathbf{c} \in \mathbb{R}^n$ , compute a (column) vector  $\mathbf{x} \in \mathbb{R}^n$  satisfying

$$\text{Minimize } \mathbf{c}^T \mathbf{x} \quad \text{subject to } \mathbf{A} \mathbf{x} = \mathbf{b}, x_i \geq 0 \quad (4)$$

The signal representation problem is posed in BP as a linear program with the following assignments (the variables on the right hand side are as defined in Section 1.1 and the variables on the left hand side plug into the linear program above):

$$\begin{aligned} \mathbf{A} &\leftarrow [\psi_1 \ \psi_2 \ \cdots \ \psi_n \ -\psi_1 \ -\psi_2 \ \cdots \ -\psi_n] \\ \mathbf{b} &\leftarrow \mathbf{x} \\ \mathbf{c} &\leftarrow [1 \ 1 \ \cdots \ 1] \\ \mathbf{x} &\leftarrow [\alpha_1 \ \alpha_2 \ \cdots \ \alpha_n] \end{aligned}$$

Minimizing  $\mathbf{c}^T \mathbf{x}$  is equivalent to minimizing the  $\ell^1$ -norm of the coefficients. Note that  $\mathbf{A}$ , corresponding to the dictionary, is doubled to include the negative of each atom; this is due to the linear programming constraint that coefficients be positive.

Chen *et al.* [17, 18] describe two algorithms for BP, BP-Simplex and BP-Interior, which are the well-known simplex and interior point methods of linear programming [98] applied to signal representation. The choice of which BP algorithm to use depends on the structure of the dictionary: for

dictionaries that have fast transforms, BP-Interior exploits these transforms in the solution of the corresponding linear program. However, the running time of linear programming is still typically an order of magnitude slower than that of MP [18].

While standard linear programming methods have been highly tuned over time, they are not necessarily ideally suited to the specific problem of computing signal representations. For example, many linear programming methods assume that the matrix  $\mathbf{A}$  is sparse, as is the case for constraints that arise in typical operations research problems, while this may not be the case in signal representation problems. This raises the possibility that alternative approaches could prove more efficient for the particular problem of signal representation. Some inspiration for an alternative approach is provided by Chen *et al.* [18], who contrast MP and BP-Simplex, characterizing MP as a ‘build-up’ approach and BP-Simplex as a ‘swap-down’ approach. If  $\mathbf{A}$  is not sparse, then the swaps (or pivots) executed by BP-Simplex can be costly, in the computation of an individual swap, in the number of swaps, and in the computation of an initial basis. GBP instead takes the ‘build-up’ approach to solving linear programming.

### 2.3 Subset selection

Sparse signal representation is closely related to the problem of subset selection for regression, i.e., determining the optimal subset of variables on which to regress a data set [69]. In sparse signal representation, the signal corresponds to the data set, while the atoms correspond to the variables. In fact, MP was inspired by Projection Pursuit [45, 56], in particular its use as a regression algorithm [44]. Given this connection, it should not be surprising that some algorithmic ideas in sparse signal representation correspond to earlier work in regression. For example, in Forward Selection the optimal subset is constructed by starting with the empty subset and iteratively adding variables to it, selecting at each iteration the variable that accounts for most of the residual variance; this is essentially what OMP does. Backward Elimination, which starts with the full set of variables and iteratively pares it down, has similarly been adapted for signal representation [54, 22].

One standard algorithm for subset selection in regression which appears to have no analogue in sparse signal representation is Efroymson’s algorithm [38], also called step-wise regression, proceeds like Forward Selection, but, like Backward Elimination, drops variables from the subset as they become irrelevant. GBP follows a similar strategy, iteratively selecting atoms and occasionally discarding them.

## 2.4 Linear programming

While Basis Pursuit represents the first formal casting of signal representation as linear programming, linear programming has long been used in sparse signal representation, particularly for deconvolution in various applications [37, 5, 74]. It is therefore not surprising that developments in sparse signal representation closely parallel earlier developments in linear programming.

Examining the literature in linear programming reveals that MP and OMP have linear programming analogues: MP is technically equivalent to one of the earliest (1948) methods developed for linear programming, called von Neumann’s algorithm [23]. Similarly, OMP is equivalent to a phase I algorithm [62] for the simplex method.

GBP builds up a solution to a linear programming problem; several linear programming methods adopt a similar strategy, solving increasingly complex problems as constraints or variables are iteratively introduced [93, 87, 76]; see also [55]. We remark that one method, an interior point method called the gravitational method [72, 15], can be shown to be equivalent to GBP when applied to the problem dual to (4). Empirically, the gravitational method is faster than standard methods on some problems [15], which is consistent with our results.

## 3 The Geometry of Basis Pursuit

GBP is based on computational geometry, specifically on the following geometric interpretation of BP. Given a signal  $\mathbf{x}$  and a dictionary  $\mathcal{D}$ , let  $\mathbf{conv}(\mathcal{D})$  denote the convex hull of  $\mathcal{D}$ ; *the vertices of the facet of  $\mathbf{conv}(\mathcal{D})$  intersected by the vector  $\mathbf{x}$  are the atoms in the minimum  $\ell^1$ -norm representation of  $\mathbf{x}$ .*

To see this, treat the signal as a vector and the atoms as points in  $\mathbb{R}^d$ . First consider the set of signals that have representations  $\alpha$  such that  $\|\alpha\|_1 = 1$ . By definition, this is the convex hull of the dictionary

$$\mathbf{conv}(\mathcal{D}) = \left\{ \mathbf{x} \mid \mathbf{x} = \sum_{i \in \mathcal{I}} \alpha_i \psi_i \text{ and } \sum_{i \in \mathcal{I}} \alpha_i = 1, \alpha_i > 0 \right\}$$

Note that because  $\|\psi_i\|_2 = 1$ ,  $\mathbf{conv}(\mathcal{D})$  is a polytope inscribed in the unit sphere. Let  $\mathbf{x}_{\mathcal{D}}$  be the point of intersection between the vector  $\mathbf{x}$  and the boundary of  $\mathbf{conv}(\mathcal{D})$ .  $\mathbf{x}_{\mathcal{D}}$  lies on the boundary of  $\mathbf{conv}(\mathcal{D})$  and can be represented as a linear combination of the vertices of the facet of  $\mathbf{conv}(\mathcal{D})$



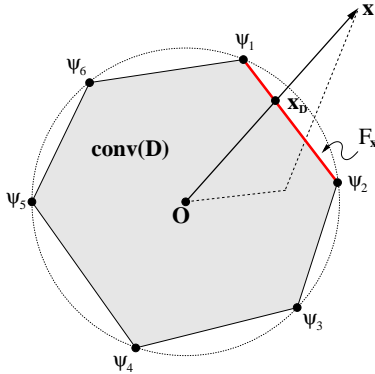


Figure 1: *A geometric interpretation of Basis Pursuit.* The signal vector  $\mathbf{x}$  intersects the facet  $F_{\mathbf{x}}$  of the convex hull of the dictionary, shown in gray. The vertices of  $F_{\mathbf{x}}$ ,  $\psi_1$  and  $\psi_2$ , are the atoms in the Basis Pursuit representation of  $\mathbf{x}$ .

containing  $\mathbf{x}_{\mathcal{D}}$ ; call this facet  $F_{\mathbf{x}}$ . This representation is the minimum  $\ell^1$ -norm representation: its  $\ell^1$ -norm is 1, and it is impossible to construct a representation with  $\ell^1$ -norm less than 1. The minimum  $\ell^1$ -norm representation of  $\mathbf{x}$  is simply a scaling of the minimum  $\ell^1$ -norm representation of  $\mathbf{x}_{\mathcal{D}}$ , and the atoms in the representation are the same. See Figure 1. (Note that if we know the atoms in a representation of  $\mathbf{x}$  it is straightforward to calculate the corresponding coefficients.)

Thus BP is equivalent to finding the facet of  $\mathbf{conv}(\mathcal{D})$  which intersects  $\mathbf{x}$ . Computing this intersection is known to reduce to linear programming [85]; to our knowledge, the converse is known [12] but never utilized to solve linear programming. We use this equivalence to drive GBP.

A previous geometric interpretation of sparse representation [11] recognizes that in two dimensions BP computes representations with atoms that ‘enclose’  $\mathbf{x}$ . The interpretation provided here can be viewed as the generalization of this notion to higher dimensions.

## 4 The Greedy Basis Pursuit Algorithm

Given the equivalence between BP and finding the facet of the convex hull of the dictionary that intersects the signal vector, we propose Greedy Basis Pursuit (GBP). GBP computes the minimum  $\ell^1$ -norm representation by

searching for this facet directly.

The main idea behind GBP is to find the facet of interest by iteratively ‘pushing’ a hyperplane onto the surface of the convex hull of the dictionary until it coincides with the supporting hyperplane containing the facet. This approach is inspired by gift-wrapping methods [14, 59, 94] for the convex hull problem in computational geometry [86]. To adapt gift-wrapping to the problem of finding a particular facet, we need to specify how the initial hyperplane is chosen and the direction in which the ‘wrapping’ proceeds at each iteration. Below we describe the GBP algorithm, prove its convergence, and discuss implementation issues.

## 4.1 The main algorithm

GBP takes as input a signal  $\mathbf{x} \in \mathbb{R}^d$  and an overcomplete dictionary  $\mathcal{D} = \{\psi_i\}_{i=1}^n$ , where  $n > d$  and  $\forall i, \psi_i \in \mathbb{R}^d$  and  $\|\psi_i\|_2 = 1$ , and outputs a representation of  $\mathbf{x}$  as a set of indices  $\mathcal{I} \subseteq \{1, \dots, n\}$  and a corresponding set of coefficients  $\mathcal{A} = \{\alpha_i\}_{i \in \mathcal{I}}$  such that  $\mathbf{x} = \sum_{i \in \mathcal{I}} \alpha_i \psi_i$ .

GBP greedily searches for the facet of  $\mathbf{conv}(\mathcal{D})$  that intersects  $\mathbf{x}$ , call it  $F_{\mathbf{x}}$ . GBP proceeds by iteratively constructing a sequence of hyperplanes,  $H^{(0)}, H^{(1)}, H^{(2)}, \dots$ , supporting  $\mathbf{conv}(\mathcal{D})$ . (We use the superscript  $(t)$  to denote iteration  $t$ .) At each iteration, GBP maintains a set of indices  $\mathcal{I}^{(t)}$  and a set of coefficients  $\mathcal{A}^{(t)}$ , defining an approximation to  $\mathbf{x}$ :  $\tilde{\mathbf{x}}^{(t)} = \sum_{i \in \mathcal{I}^{(t)}} \alpha_i \psi_i$ , and a normal vector  $\mathbf{n}^{(t)}$ . The current hyperplane  $H^{(t)}$  is defined to have normal  $\mathbf{n}^{(t)}$  and contain the set  $\{\psi_i\}_{i \in \mathcal{I}^{(t)}}$ . Each consecutive hyperplane  $H^{(t+1)}$  is a rotation of the current hyperplane  $H^{(t)}$  determined by  $\tilde{\mathbf{x}}^{(t)}$ . GBP stops when  $H^{(t)}$  contains  $F_{\mathbf{x}}$  (and therefore  $\tilde{\mathbf{x}}^{(t)} = \mathbf{x}$ ).

### 4.1.1 Initialization

As we do not *a priori* know the orientation of  $F_{\mathbf{x}}$ , we optimistically choose the initial supporting hyperplane  $H^{(0)}$  to have normal  $\mathbf{n}^{(0)} = \mathbf{x}/\|\mathbf{x}\|_2$ . In general  $H^{(0)}$  will intersect only one vertex of  $\mathbf{conv}(\mathcal{D})$ , in particular the atom  $\psi_{i_0}$ , where  $i_0 = \arg \max_i \langle \psi_i, \mathbf{n}^{(0)} \rangle$ . To see this, consider a hyperplane with normal  $\mathbf{n}^{(0)}$  at some distance greater than 1 away from the origin; if we move this hyperplane in the negative normal direction (towards the origin), the first point of  $\mathbf{conv}(\mathcal{D})$  it will intersect is  $\psi_{i_0}$ . (Note that this is also the first atom selected by MP and OMP.) Thus  $\mathcal{I}^{(0)} = \{i_0\}$ ; this gives us  $\alpha_{i_0} = \langle \psi_{i_0}, \mathbf{x} \rangle$ ,  $\mathcal{A}^{(0)} = \{\alpha_{i_0}\}$ , and  $\tilde{\mathbf{x}}^{(0)} = \alpha_{i_0} \psi_{i_0}$ . For convenience, we denote the set of currently selected atoms by  $\Psi^{(t)} = \{\psi_i\}_{i \in \mathcal{I}^{(t)}}$ .

### 4.1.2 Iteration

Each consecutive hyperplane  $H^{(t+1)}$  is constructed by rotating  $H^{(t)}$  in a 2-dimensional plane around a pivot point until another vertex of  $\mathbf{conv}(\mathcal{D})$  is intersected. The plane of rotation and the pivot point are defined in terms of  $\tilde{\mathbf{x}}^{(t)}$ . We define  $\tilde{\mathbf{x}}^{(t)}$  to be the best current approximation to  $\mathbf{x}$  using  $\Psi^{(t)}$  and positive coefficients, that is,  $\tilde{\mathbf{x}}^{(t)} = \sum_{i \in \mathcal{I}^{(t)}} \alpha_i \psi_i$ , where  $\alpha_i > 0$  and  $\|\tilde{\mathbf{x}}^{(t)} - \mathbf{x}\|_2$  is minimized. Note that  $\tilde{\mathbf{x}}^{(t)}$  is the projection of  $\mathbf{x}$  onto the convex cone spanned by  $\Psi^{(t)}$  with the origin at the apex; we provide details on computing  $\tilde{\mathbf{x}}^{(t)}$  in Section 4.1.3. Let  $\tilde{\mathbf{x}}_H^{(t)}$  be the intersection of the vector  $\tilde{\mathbf{x}}^{(t)}$  with  $H^{(t)}$ . If  $d_H^{(t)}$  is the distance from the hyperplane to the origin (in the normal direction), e.g.,  $d_H^{(t)} = \langle \psi_i, \mathbf{n} \rangle$ ,  $i \in \mathcal{I}^{(t)}$ , then

$$\tilde{\mathbf{x}}_H^{(t)} = \left( d_H^{(t)} / \langle \tilde{\mathbf{x}}^{(t)}, \mathbf{n}^{(t)} \rangle \right) \tilde{\mathbf{x}} \quad (5)$$

Let  $\mathbf{r}^{(t)}$  denote the residual vector,  $\mathbf{r}^{(t)} = \mathbf{x} - \tilde{\mathbf{x}}^{(t)}$ . Define  $\mathbf{v}^{(t)}$  to be the unit vector in the direction of  $\mathbf{r}^{(t)}$  projected onto  $H^{(t)}$ .

$$\mathbf{v}^{(t)} = \frac{\mathbf{r}^{(t)} - \langle \mathbf{r}^{(t)}, \mathbf{n}^{(t)} \rangle \mathbf{n}^{(t)}}{\|\mathbf{r}^{(t)} - \langle \mathbf{r}^{(t)}, \mathbf{n}^{(t)} \rangle \mathbf{n}^{(t)}\|} \quad (6)$$

The plane of rotation is the 2-dimensional plane defined by the point  $\tilde{\mathbf{x}}_H^{(t)}$  and the vectors  $\mathbf{n}^{(t)}$  and  $\mathbf{v}^{(t)}$ . The pivot point around which  $H$  is rotated is  $\tilde{\mathbf{x}}_H^{(t)}$ .

To compute the first vertex which the hyperplane intersects under this rotation, we order the atoms by the angle  $\theta$  they form with  $\mathbf{v}$ , where  $\theta$  is given by

$$\theta_i = \arctan(\langle \psi_i - \tilde{\mathbf{x}}_H^{(t)}, \mathbf{n}^{(t)} \rangle / \langle \psi_i - \tilde{\mathbf{x}}_H^{(t)}, \mathbf{v}^{(t)} \rangle)$$

The atom selected is then  $\psi_k$  where

$$k = \arg \min_i \theta_i \quad (7)$$

Once selected, the atom  $\psi_k$  is added to the set  $\Psi^{(t)}$  and a new approximation to  $\mathbf{x}$  is computed,  $\tilde{\mathbf{x}}^{(t+1)}$ . In this new approximation, some atoms in  $\Psi^{(t)} \cup \{\psi_k\}$  may be extraneous; they are discarded to form  $\Psi^{(t+1)}$ .

The new hyperplane  $H^{(t+1)}$  can now be computed; it has normal

$$\mathbf{n}^{(t+1)} = -\langle \psi_k - \tilde{\mathbf{x}}_H^{(t)}, \mathbf{n}^{(t)} \rangle \mathbf{v}^{(t)} + \langle \psi_k - \tilde{\mathbf{x}}_H^{(t)}, \mathbf{v}^{(t)} \rangle \mathbf{n}^{(t)} \quad (8)$$

and contains  $\tilde{\mathbf{x}}_H^{(t+1)}$ .



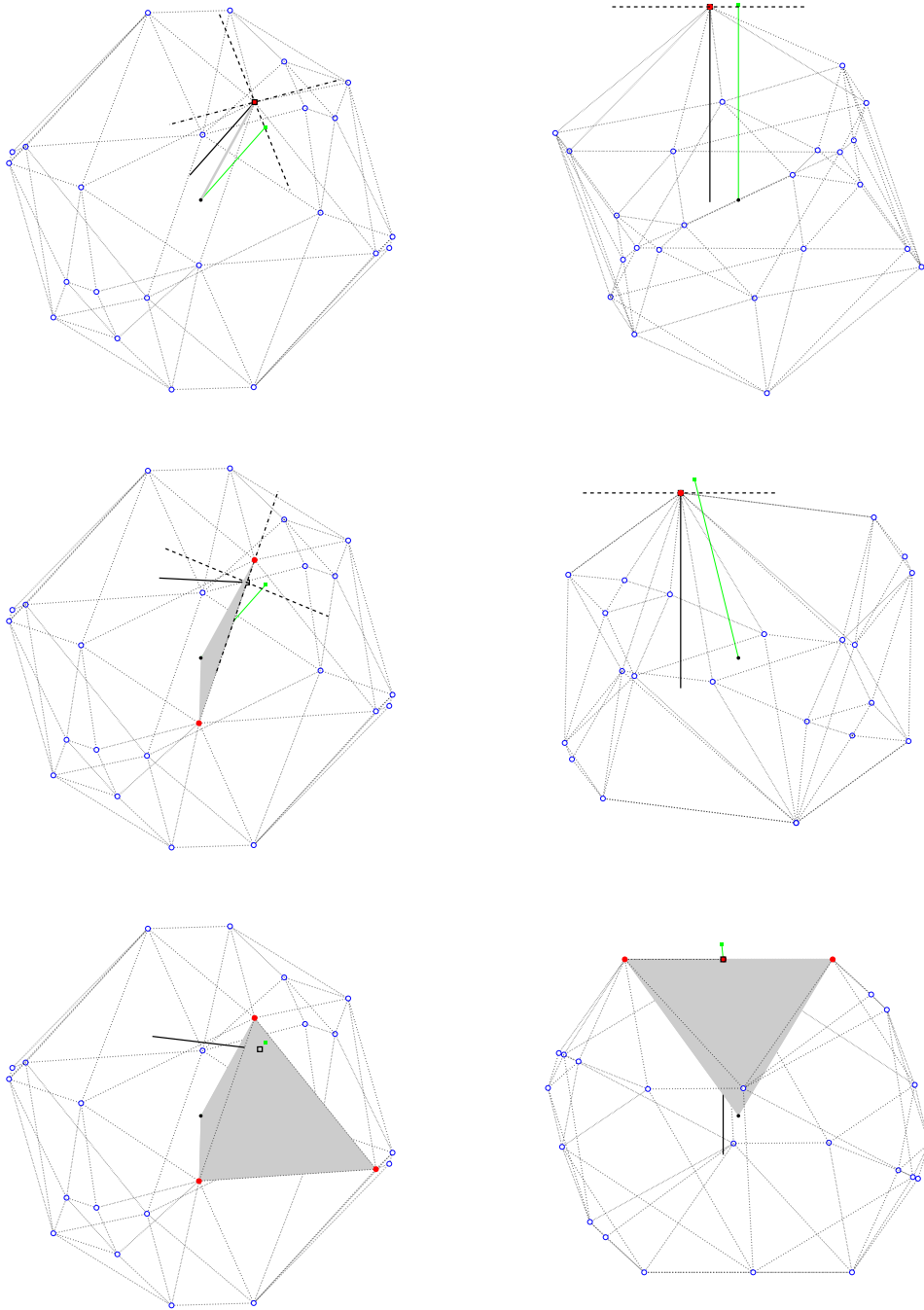


Figure 3: GBP in action on a 3-dimensional problem. See text for details.

the following two equations:

$$\langle \psi_i, \tilde{\psi}_i^{\perp(t)} \rangle = 1 \quad (9)$$

$$\langle \psi_i, \tilde{\psi}_j^{\perp(t)} \rangle = 0, \text{ if } i \neq j \quad (10)$$

The biorthogonal vector  $\tilde{\psi}_i^{\perp(t)}$  can be understood as the component of  $\psi_i^{(t)}$  that is orthogonal to all of the other vectors in  $\Psi^{(t)}$ , appropriately scaled. That is, if we express an atom  $\psi_i \in \Psi^{(t)}$  as

$$\psi_i = \psi_i^{\parallel(t)} + \psi_i^{\perp(t)} \quad (11)$$

where  $\psi_i^{\parallel(t)}$  is the component of  $\psi_i$  lying in the span of  $\Psi^{(t)} - \{\psi_i\}$

$$\psi_i^{\parallel(t)} = \sum_{j \in \mathcal{I}^{(t)}, j \neq i} \beta_{ij}^{(t)} \psi_j \quad (12)$$

and  $\psi_i^{\perp(t)}$  is orthogonal to the span of  $\Psi^{(t)} - \{\psi_i\}$ , then the biorthogonal vector to  $\psi_i^{(t)}$  is given by

$$\tilde{\psi}_i^{\perp(t)} = \psi_i^{\perp(t)} / \|\psi_i^{\perp(t)}\|^2 \quad (13)$$

Given the biorthogonal system, we can compute the current approximation to  $\mathbf{x}$  as

$$\tilde{\mathbf{x}}^{(t)} = \sum_{i \in \mathcal{I}^{(t)}} \alpha_i^{(t)} \psi_i \quad \text{where } \alpha_i^{(t)} = \langle \mathbf{x}, \tilde{\psi}_i^{\perp(t)} \rangle \quad (14)$$

The biorthogonal system and  $\tilde{\mathbf{x}}^{(t)}$  can be updated as atoms are added to and subtracted from  $\Psi^{(t)}$ . Such adaptive biorthogonalization methods have recently been applied to MP [83, 4] and are standard in linear programming ([98], Chapter 8). We present pseudocode for adding an atom in Algorithm 2 and for subtracting an atom in Algorithm 3.

## 4.2 Analysis

By construction, GBP computes the minimum  $\ell^1$ -norm representation of a given signal. To prove this we show that GBP converges to an exact representation in a finite number of steps and that the representation corresponds to a facet of the convex hull of the dictionary.

First, we prove that GBP converges to an exact representation. At each iteration of GBP there is a decrease in approximation error, as stated in the following theorem.

**Theorem 1.** *Given a signal  $\mathbf{x} \in \mathbb{R}^d$  and a dictionary  $\mathcal{D} = \{\psi_i\}_{i=1}^n$ , where  $n \geq 2d$ ,  $\forall \psi_i \in \mathcal{D}$ ,  $\psi_i \in \mathbb{R}^d$  and  $\|\psi_i\|_2 = 1$ , if  $\psi_i \in \mathcal{D}$  then  $-\psi_i \in \mathcal{D}$ , and the atoms are in general position, if GBP is run with  $\mathcal{D}$  and  $\mathbf{x}$  as input and if  $\tilde{\mathbf{x}}^{(t)} \neq 0$ , then at iteration  $t + 1$  of GBP,  $0 \leq \|\mathbf{x} - \tilde{\mathbf{x}}^{(t+1)}\|_2 < \|\mathbf{x} - \tilde{\mathbf{x}}^{(t)}\|_2$ .*

*Proof.* At iteration  $t$ , let  $S$  be the hypersphere centered at  $\mathbf{x}$  with radius  $\|\mathbf{x} - \tilde{\mathbf{x}}^{(t)}\|_2$ , let  $\psi_k$  be the next atom selected by GBP, and let  $T$  be the tangent plane to  $S$  at  $\tilde{\mathbf{x}}^{(t)}$ .  $T$  contains the origin (if it did not, then some scaling of  $\tilde{\mathbf{x}}^{(t)}$  would be a better approximation to  $\mathbf{x}$ ), and thus bisects the unit sphere. Because the atoms are in general position and  $\psi_i \in \mathcal{D}$  if  $-\psi_i \in \mathcal{D}$ , if  $|\Psi^{(t)}| < d$ , then there will be at least one atom in the same half-space of  $T$  as  $\mathbf{x}$ . (Note that if  $|\Psi^{(t)}| = d$ , we are done, as we would also have  $\tilde{\mathbf{x}} = \mathbf{x}$ .)

$\psi_k$  lies in the same half-space of  $T$  as  $\mathbf{x}$ : by construction, there is no atom  $\psi_0$  such that  $\langle \psi_0 - \tilde{\mathbf{x}}_H^{(t)}, \mathbf{n}^{(t)} \rangle > 0$ , by general position, there is no atom  $\psi_0$  such that  $\langle \psi_0 - \tilde{\mathbf{x}}_H^{(t)}, \mathbf{n}^{(t)} \rangle = 0$  and  $\langle \psi_0 - \tilde{\mathbf{x}}_H^{(t)}, \mathbf{v}^{(t)} \rangle > 0$ , and, by the ordering of atoms by step 2(b) of GBP, GBP selects an atom in the same half-space of  $T$  as  $\mathbf{x}$ , if one exists.

$\psi_k$  lies in the same half-space as  $\mathbf{x}$ , we can find a point  $\tilde{\mathbf{x}} + \epsilon(\psi_k - \tilde{\mathbf{x}}^{(t)})$  that is interior to  $S$  and therefore closer to  $\mathbf{x}$  than  $\tilde{\mathbf{x}}^{(t)}$ . Therefore  $\|\mathbf{x} - \tilde{\mathbf{x}}^{(t+1)}\| < \|\mathbf{x} - \tilde{\mathbf{x}}^{(t)}\|$ .  $\square$

Theorem 1 also implies that GBP does not cycle. GBP may select the same atom more than once, that is, GBP may select an atom, discard it, and select it again (this behaviour depends on the shape of the facets of  $\mathbf{conv}(\mathcal{D})$ ), but GBP will never revisit the same state. Because there are a finite number of states and GBP improves at each iteration, GBP converges. By the same arguments as Theorem 1, at convergence the final supporting hyperplane contains a facet of  $\mathbf{conv}(\mathcal{D})$  and thus GBP computes the minimum  $\ell^1$ -norm representation.

The duality of GBP to the gravitational method [72] of linear programming, implies that the computational complexity of GBP is exponential in the worst-case [71]. Current results on the simplex algorithm suggest that GBP is likely to be polynomial in the average [12] and smoothed [91] cases.

### 4.3 Implementation Issues

We briefly describe two obstacles that any implementation of GBP may encounter, degeneracy and numerical instability, and our approach to handling them.

### 4.3.1 Degeneracy

Degeneracy occurs when the atoms of the dictionary are not in general position, that is, a  $k$ -face of  $\mathbf{conv}(\mathcal{D})$  contains more than  $k + 1$  atoms. Degeneracy can occur if the dictionary is specially designed, for example, if the atoms are defined to be the vertices of a hypercube inscribed in the unit hypersphere. If GBP encounters degeneracy, the updates described in Section 4.1.3 will fail, resulting in an error. Although GBP does not currently include a mechanism to detect and handle degeneracy, incorporating such a feature is possible. A simple solution is to perturb the atoms of the dictionary sufficiently.

### 4.3.2 Numerical instability

Numerical instability can occur in the biorthogonalization stage of GBP. Let  $\mathbf{\Psi}$  be a matrix corresponding to  $\Psi^{(t)}$  for some  $t$ , where each row of  $\mathbf{\Psi}$  is an atom in  $\Psi^{(t)}$ , and let  $\tilde{\mathbf{\Psi}}^\perp$  denote the corresponding matrix of biorthogonal vectors. If at any iteration the matrix  $\mathbf{\Psi}$  is ill-conditioned, the computation of the biorthogonal vectors we have described may be unstable (similar difficulties arise in Gram-Schmidt orthogonalization [84, 10]). One work around is to compute a full biorthogonalization at each iteration, or at least whenever instability is detected. However, a full biorthogonalization can be costly, as it is typically computed via the pseudoinverse [52]: since  $\mathbf{\Psi} \left( \tilde{\mathbf{\Psi}}^\perp \right)^T = \mathbf{I}$ , where  $\mathbf{I}$  denotes the identity matrix, we can compute  $\tilde{\mathbf{\Psi}}^\perp$  as  $(\mathbf{\Psi}^+)^T$ , where ‘+’ denotes the pseudoinverse.

We instead opt to compute the biorthogonalization using an iterative pseudoinverse technique [6]. This technique takes an initial estimate of the pseudoinverse and iteratively updates it, converging to the true pseudoinverse. If the initial estimate is sufficiently close to the true pseudoinverse, then the iterative pseudoinverse computation is substantially faster than the standard pseudoinverse. This approach is well suited to GBP, as the adaptive biorthogonalization already provides such an estimate.

The iterative pseudoinverse algorithm proceeds as follows. Given a matrix  $\mathbf{\Psi}$  and an initial estimate of the pseudoinverse  $\mathbf{\Psi}^{+(0)}$ , the updates to  $\mathbf{\Psi}^+$  are computed by

$$\mathbf{\Psi}^{+(t+1)} \leftarrow \mathbf{\Psi}^{+(t)} \left( 2\mathbf{I} - \mathbf{\Psi} \mathbf{\Psi}^{+(t)} \right)$$

(Note that here  $t$  denotes the iteration of the pseudoinverse algorithm, not the iteration of GBP.) For a detailed analysis of this algorithm, see



[90]. While a classic technique, this algorithm is the subject of ongoing research [77, 80].

Our implementation of GBP tests if  $\Psi \left( \tilde{\Psi}^\perp \right)^T = \mathbf{I}$  within a specified level of tolerance after each adaptive biorthogonalization. If the test fails, the iterative pseudoinverse algorithm is applied.

## 5 Results

We examine the performance of GBP. We compared the running time of GBP to that of standard linear programming algorithms on three data sets, random data, speech data, and seismic data, described below. We also provide an example of GBP's performance on a single signal and contrast it with that of Matching Pursuit.

In each experiment, we measured the running times of GBP and standard linear programming algorithms on the signal representation problems described below. The algorithms we compared were GBP, two variants of the simplex method, and an interior point method.

The implementation of GBP used was our own, written entirely in Matlab. The linear programming solvers used were those included in the Matlab Optimization Toolbox 3.0 [1], and a freely available Matlab implementation [70] of the revised simplex method [24]. The Optimization Toolbox version of the simplex method is the classical simplex method [25], with the initial basis determined as in [8]. The Optimization Toolbox version of the interior point method is essentially LIPSOL [99], a freely available interior point solver that implements Mehrotra's predictor-corrector method [68, 65].

For each problem, all algorithms were run and timed. All algorithms were run under Matlab 7 on a 1.5GHz Pentium M processor running Windows XP, with 1.25GB memory. On all problems all algorithms returned identical representations.

### 5.1 Running times: Random data

The random data set consisted of 3000 randomly generated signal representation problems, varying both the dimension of the signal space and the overcompleteness of the dictionary. Each problem consisted of a randomly generated signal and a randomly generated dictionary. The dimension  $d$  of the problems varied over the set  $\{8, 16, 32, 64, 128, 256\}$ . The overcompleteness  $k$  of the dictionaries varied over the set  $\{2, 4, 8, 16, 32\}$ . In each problem, the signal  $\mathbf{x}$  was randomly generated to be uniformly distributed

on the unit hypersphere in  $\mathbb{R}^d$ . The dictionary for each problem had  $2kd$  atoms; the first  $kd$  of these atoms were generated in the same fashion as the signal, the second  $kd$  atoms were the negatives of the first  $kd$  atoms. Additionally, the dictionary of each problem was perturbed: To each atom was added Gaussian noise with variance 0.000001, after which the atom was normalized to lie on the unit hypersphere; this was necessary to ensure that the linear programming algorithms could compute the requisite matrix inverses. For each  $d$ - $k$  pair, 100 problems were generated.

Figure 4 shows running times of the three algorithms as a function of overcompleteness for each dimension; the curve plotted shows the mean running time of each algorithm over the 100 problems of the specified dimension and overcompleteness, with error bars showing the corresponding minimum and maximum running times. (We do not show the results of the revised simplex method here, as it was outperformed by the Matlab’s simplex algorithm.)

## 5.2 Running times: Speech data

The speech data set consisted of 100 signal representation problems. Each problem consisted of a signal randomly drawn from the TIMIT database [48] and an overcomplete multiscale Gabor dictionary.

Each signal comprised 256 samples ( $d = 256$ ) and was randomly selected from the ‘train’ subset of the TIMIT database. The signals were mean centered and normalized. Some samples are shown in Figure 5.

The dictionary used was a  $9 \times$  overcomplete multiscale Gabor dictionary (4608 atoms). The dictionary consisted of several fixed scale critically sampled cosine Gabor bases. Each atom was defined by the parameters  $t$  and  $f$  as  $G(t, f) = \exp^{-\sigma^2 t^2} \cos(2\pi ft)$ , where  $t \in \{0 : \Delta t : 1\}$  and  $f \in \{0 : \Delta f : d/2\}$ , with  $\Delta t = 2^j/d$ ,  $\sigma = \sqrt{\pi/2}/\Delta t$ , and  $\Delta f = \sigma/\sqrt{2\pi}$ ; the scale parameter  $j$  varied over  $\{0, 1, \dots, 8\}$ . See [47, 40] for details and other sampling schemes. Once the atoms were defined, they were perturbed as in the random data case. Some samples from the final dictionary are shown in Figure 6.

We show the running times of GBP, LIPSOL, and the revised simplex method on the sound data set in Table 1. (The revised simplex method outperformed Matlab’s simplex method.) We show the mean, minimum, and maximum running times for each algorithm on the 100 signals.

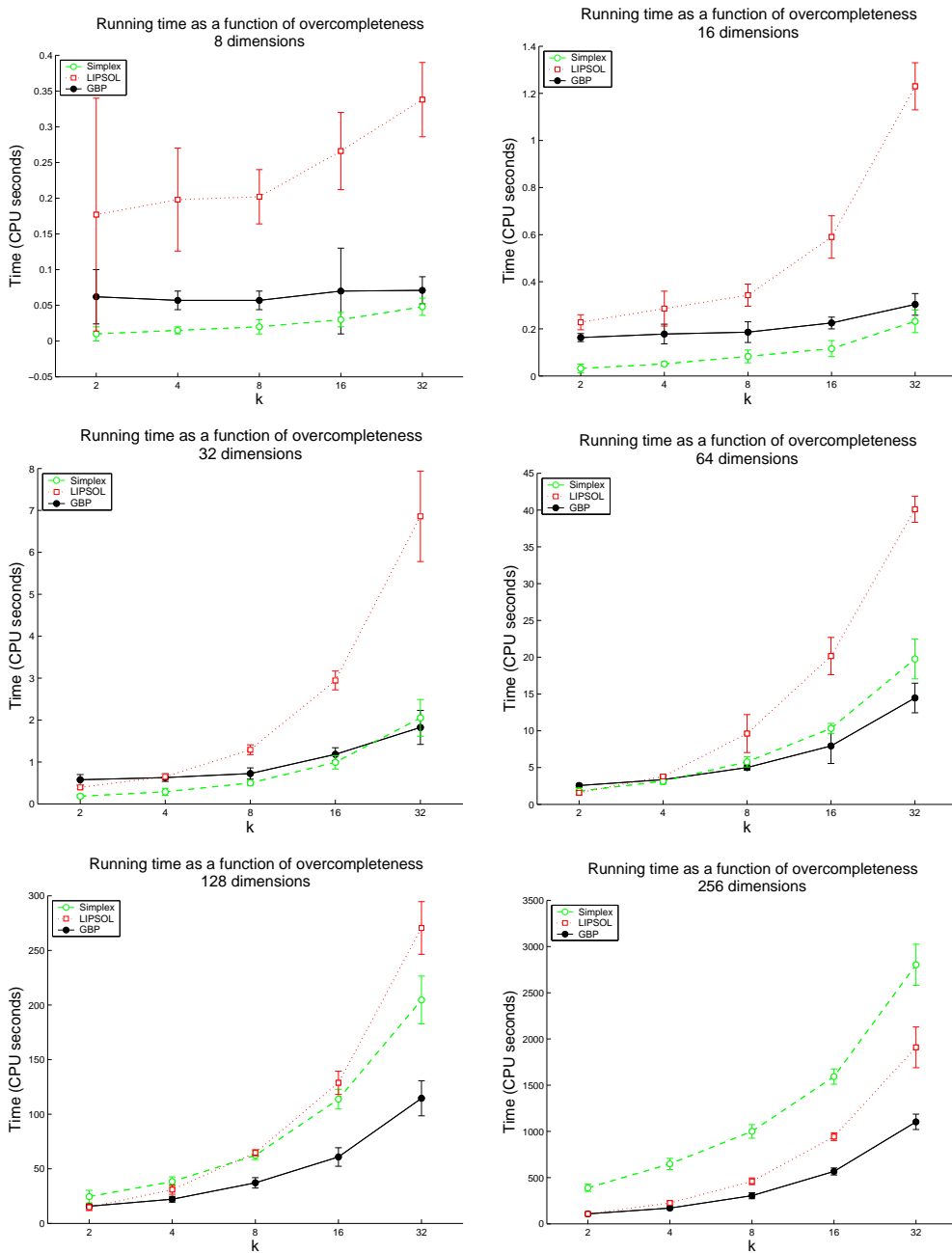


Figure 4: Running times for GBP, LIPSOL, and the simplex method (Matlab) on the random data set, plotted as a function of overcompleteness for each dimension. Note that GBP's performance improves relative to the other methods as the dimensionality of the problems increase.

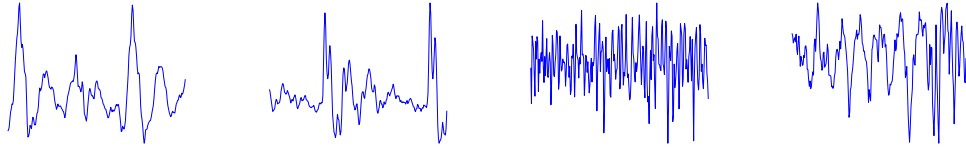


Figure 5: Four signals drawn from the speech data set. Each signal consists of 256 samples.

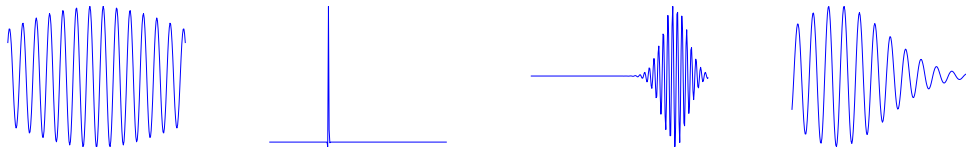


Figure 6: Four atoms drawn from the multiscale Gabor cosine dictionary.

Algorithm	Min	Mean	Max
GBP	40.41	48.72	56.35
LIPSOL	58.62	75.97	155.56
Revised Simplex	441.66	1297.65	2700.51

Table 1: Running times of GBP, LIPSOL, and the revised simplex method on the sound data set, in CPU seconds.



Figure 7: Four signals drawn from the seismic data set. Each signal consists of 256 samples.

Algorithm	Min	Mean	Max
GBP	42.29	48.83	55.05
LIPSOL	60.52	70.36	112.90
Revised Simplex	2233.45	2489.05	2831.59

Table 2: Running times of GBP, LIPSOL, and the revised simplex method on the seismic data set, in CPU seconds.

### 5.3 Running times: Seismic data

The seismic data consists of 100 signal representation problem. Each problem consists of a 256 sample signal of seismic recordings from the North Sea, 4 times downsampled from the original data [89]; some samples are shown in Figure 7. The dictionary used was the same as used in the speech experiment above. We show the running times of GBP, LIPSOL, and the revised simplex method on the seismic data set in Table 2.

### 5.4 Example: Speech signal

Figure 8 provides an example comparing GBP to MP on a 1024-dimensional signal (Figure 8, top left), selected from the TIMIT speech database [48], using a multiscale Gabor dictionary ( $n = 22528$ ), similar to the one used for the sound data. (Note that the other BP methods were unable to compute representations on problems of this size in our environment.) Examining the approximation error of each algorithm as a function of iteration (Figure 8, top right), we observe that while the approximation error of GBP decreases somewhat more slowly than that of MP (note also that each iteration of GBP is more costly), the error of GBP does appear to decrease approximately exponentially. Furthermore, the representation computed by GBP is considerably sparser than that of MP, as indicated by the sorted-amplitudes-curves and the  $\ell^1$ -norm of the representations. The sorted-amplitudes-curves [61, 57] (Figure 8, bottom left) are plots of the

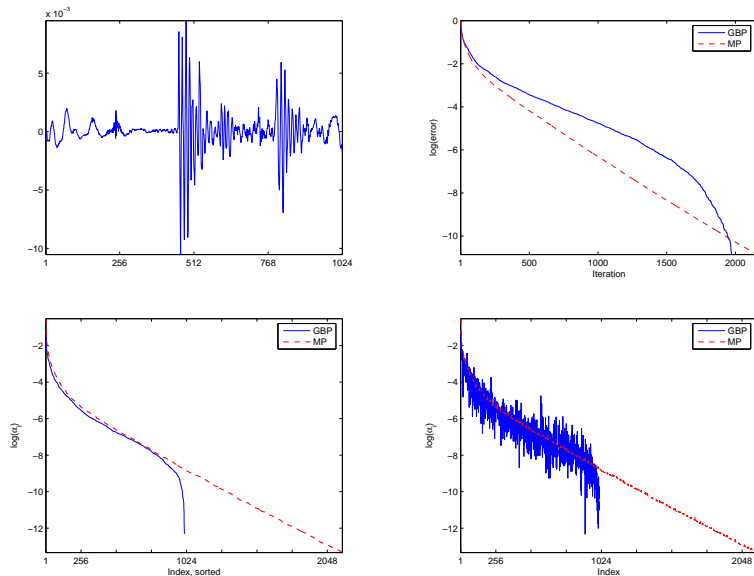


Figure 8: An example comparing GBP to MP on a speech signal (a). (b) The log of the error ( $\ell^2$ -norm of the residual) as a function of iteration. (c) The sorted-amplitudes-curves; observe that GBP produces a sparser representation than MP. (d) The (final) coefficient values, in order of atom selection. (Note that the coefficient values change in GBP at each iteration.) See text for discussion.

logarithm of the final coefficients, sorted in descending order; the rates of decrease indicate the relative sparsity of the representations. The  $\ell^1$ -norm of the representation coefficients are 6.33 and 8.83, for GBP and MP, respectively. (Note that the results for GBP would be the same as those for standard linear programming methods.) A notable feature of GBP is its ‘greediness’: the coefficients in the order of atom selection track the sorted-amplitudes-curve, that is, GBP tends to select significant atoms early on (Figure 8, bottom right). This demonstrates that it is possible to compute Basis Pursuit signal representations and to be greedy at the same time.

## 6 Discussion

Our results show that GBP provides a fast alternative to standard linear programming methods for sparse signal representation problems, particularly when the dimension of the signal space is high and the dictionary is

very overcomplete. While there are a variety of factors which may contribute to the results, there are several algorithmic reasons why we expect GBP to perform well relatively.

The efficient solution of linear programming problems depends in a complicated way on the problem, the method of solution and its implementation, and the available resources; see Bixby [9]. Thus the relative success of GBP compared to the linear programming methods implemented in the Matlab Optimization Toolbox is partially a function of the specific methods used and their implementation. There are many possible linear programming methods, and many available implementations of those methods [43]; making an exhaustive comparison of GBP against these methods impossible. However, even if the linear programming methods against which GBP was compared do not represent the current state-of-the-art, GBP itself has the potential for significant speed increases through more efficient implementation.

Algorithmically, GBP has several advantages over standard linear program solvers. First, most linear program solvers assume, for historical reasons, that the constraint matrix is sparse, and they therefore rely on techniques that exploit this sparsity, whether or not sparsity is actually present [19]. The signal representation problems considered here are not particularly sparse: the random dictionaries used in the random data set are certainly not sparse, while the Gabor dictionary used with the sound data set is somewhat sparse, however the sparsity is not as structured as necessary for certain fast algorithms to be applicable [17, 18]. GBP does not exploit sparsity, and therefore does not suffer when it is not present. Second, GBP is efficient in the search for the next atom to select, because this search is based on a geometric criterion that involves 2 projections per possible atom. Simplex methods can be inefficient at this task as the search can involve evaluating more than 2, even  $d$ , projections per possible atom; see [98, 96]. Third, the updates in GBP are seldom of a full basis, further reducing computation. Finally, the complexity of the simplex method depends on the closeness of the initial solution to the optimal solution, which in turn depends on the phase I algorithm by which the initial basis is selected. GBP does not depend on an initial solution; in fact, GBP can be interpreted as a combined phase I/ phase II linear programming algorithm.

One area which we have not explored that merits further investigation is the dependence of the performance of GBP (and other sparse representation algorithms) on the structure of the dictionary. For example, a dictionary optimized for use with MP [27] or OMP [39] may well have very different properties from one optimized for BP. The design of dictionaries has only recently received attention in the signal processing community [27, 39, 2]

(for work in neural computation, see [75, 64]); our work suggests that the geometric properties of dictionaries play a crucial role in both the efficiency of representation algorithms and the quality of the resulting representations. Indeed, geometric considerations have already led to a better theoretical understanding of sparse signal representation [36, 35].

As noted, part of the motivation for the development of BP is the observation that MP and OMP can fail to find sparse, in the  $\ell^0$ -norm sense, signal representations [18], with much theoretical work showing under exactly what conditions BP finds sparse representations, i.e., when the minimal  $\ell^1$ -norm solution is equivalent to the minimal  $\ell^0$ -norm solution [34, 33, 46]. These findings have made BP useful for areas beyond signal representation, including compressed sensing [32] and error correcting codes [13], thus GBP may prove useful in these domains.

## 7 Conclusions

We have described GBP, a new algorithm for Basis Pursuit, and demonstrated that it is faster than standard linear programming methods on some problems, particularly in high-dimensional signal spaces using very overcomplete dictionaries. A Matlab implementation of GBP is currently available online at: <http://www.cs.yale.edu/~huggins/gbp.html>

Computational geometry has traditionally been the preserve of computer science, particularly computer graphics and theoretical computer science; its use here in the development of GBP highlights the relevance of computational geometry to signal processing. GBP also illustrates the interplay between signal processing and linear programming. That an efficient linear programming algorithm falls naturally out of sparse signal representation is surprising, and suggests that researchers in signal processing should not view linear programming, or optimization in general, as a black box: on one hand signal processing naturally defines a set of problems that can serve to drive research in linear programming, on the other hand, given the historical parallels, optimization research deserves deeper examination by the signal processing community.

## Acknowledgements

We thank Mauro Maggioni for many helpful discussions, Ohad Ben-Shahar, Pavel Dimitrov, and Gang Li for their advice, Karl Skretting for providing



the seismic data, and Dan Spielman for his insight into linear programming. Research supported by DARPA and AFOSR.

## References

- [1] *Optimization Toolbox User's Guide*. The MathWorks, 2003.
- [2] M. Aharon, M. Elad, A.M. Bruckstein, and Y. Katz. The K-SVD: An algorithm for designing of overcomplete dictionaries for sparse representation. *IEEE Trans. On Signal Processing*, In press.
- [3] O.K. Al-Shaykh, E. Miloslavsky, T. Nomura, R. Neff, and A. Zakhor. Video compression using matching pursuits. *IEEE Transactions on Circuits and Systems for Video Technology*, 9(1), 1999.
- [4] M. Andrle, L. Rebollo-Neira, and E. Sagianos. Backward-optimized orthogonal matching pursuit approach. *IEEE Signal Processing Letters*, 2004.
- [5] J.B. Bednar, R. Yarlagadda, and T. Watt.  $L_1$  deconvolution and its application to seismic signal processing. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 34(6), 1986.
- [6] A. Ben-Israel. An iterative method for computing the generalized inverse of an arbitrary matrix. *Mathematics of Computation*, 19:452–455, 1965.
- [7] F. Bergeaud and S. Mallat. Matching pursuit of images. In *Proceedings of the International Conference on Image Processing*, 1995.
- [8] R.E. Bixby. Implementing the simplex method: The initial basis. *ORSA Journal on Computing*, 4(3), 1992.
- [9] R.E. Bixby. Solving real-world linear programs: A decade and more of progress. *Operations Research*, 50(1):3–15, 2002.
- [10] A Björck. Numerics of Gram-Schmidt orthogonalization. *Linear Algebra and its Applications*, 197-198:297–316, 1994.
- [11] P. Bofill and M. Zibulevsky. Underdetermined blind source separation using sparse representations. *Signal Processing*, 81:2353–2362, 2001.
- [12] K.H. Borgwardt. *The Simplex Method - A Probabilistic Analysis*. Springer-Verlag, 1987.
- [13] E. Candes, M. Rudelson, R. Vershynin, and T. Tao. Error correction via linear programming. In *FOCS*, 2005.
- [14] D.R. Chand and S.S. Kapur. An algorithm for convex polytopes. *Journal of the Association for Computing Machinery*, 17(1):78–86, 1970.
- [15] S.Y. Chang and K.G. Murty. The steepest descent gravitational method for linear programming. *Discrete Applied Mathematics*, 25:211–239, 1989.
- [16] S. Chen and D. Donoho. Basis pursuit. In *Twenty-Eighth Asilomar Conference on Signals, Systems & Computers*, 1994.
- [17] S.S. Chen. *Basis Pursuit*. PhD thesis, Stanford University, Department of Statistics, 1995.
- [18] S.S. Chen, D.L. Donoho, and M.A. Saunders. Atomic decomposition by basis pursuit. *SIAM Review*, 43(1):129–159, 2001.

- [19] V. Chvátal. *Linear Programming*. Freeman, 1983.
- [20] T.H. Cormen, C.E. Leiserson, and R.L. Rivest. *Introduction to Algorithms*. McGraw-Hill, 1990.
- [21] S.F. Cotter, J. Adler, B.D. Rao, and K. Kreutz-Delgado. Forward sequential algorithms for best basis selection. *IEE Proc.-Vis. Signal Processing*, 146(5), 1999.
- [22] S.F. Cotter, K. Kreutz-Delgado, and B.D. Rao. Backward sequential elimination for sparse vector subset selection. *Signal Processing*, 81:1849–1864, 2001.
- [23] G.B. Dantzig. Converting a converging algorithm into a polynomial bounded algorithm. Technical Report SOL 91-5, Systems Optimization Laboratory, Stanford University, 1991.
- [24] G.B. Dantzig and W. Orchard-Hays. The product form for the inverse in the simplex method. *Mathematical Tables and Other Aids to Computation*, 8:64–67, 1954.
- [25] G.B. Dantzig, A. Orden, and P. Wolfe. The generalized simplex method for minimizing a linear form under linear inequality restraints. *Pacific Journal of Mathematics*, 5:183–195, 1955.
- [26] I. Daubechies. Time-frequency localization operators: A geometric phase space approach. *IEEE Transactions on Information Theory*, 34(4), 1988.
- [27] G. Davis. *Adaptive Nonlinear Approximations*. PhD thesis, New York University, Department of Mathematics, 1994.
- [28] G. Davis, S. Mallat, and M. Avellaneda. Adaptive greedy approximation. *Constructive Approximation*, 13:57–98, 1997.
- [29] G. Davis, S. Mallat, and Z. Zhang. Adaptive time-frequency decompositions with matching pursuit. *Optical Engineering*, 33(7):2183–2191, 1994.
- [30] R.A. DeVore and V.N. Temlyakov. Some remarks on greedy algorithms. *Advances in Computational Mathematics*, 5:173–187, 1996.
- [31] D.L. Donoho. Sparse components of images and optimal atomic decompositions. *Constructive Approximation*, 17(3):353–382, 2001.
- [32] D.L. Donoho. Compressed sensing. Technical report, Stanford University, Department of Statistics, September 2004.
- [33] D.L. Donoho and M. Elad. Optimally sparse representation in general (nonorthogonal) dictionaries via  $\ell^1$  minimization. *PNAS*, 100(5):2197–2202, 2003.
- [34] D.L. Donoho and X. Huo. Uncertainty principles and ideal atomic decomposition. *IEEE Transactions on Information Theory*, 47(7), 2001.
- [35] D.L. Donoho and J. Tanner. Neighborliness of randomly projected simplices in high dimensions. *Proceedings of the National Academy of Sciences*, 102(27):9452–9457, 2005.
- [36] D.L. Donoho and J. Tanner. Sparse nonnegative solution of underdetermined linear equations by linear programming. *Proceedings of the National Academy of Sciences*, 102(27):9446–9451, 2005.
- [37] B. Drachman. Two methods to deconvolve:  $L_1$ -method using simplex algorithm and  $L_2$ -method using least squares and a parameter. *IEEE Transactions on Antennas and Propagation*, 32(3), 1984.

- [38] M.A. Efronson. Multiple regression analysis. In A. Ralston and H.S. Wilf, editors, *Mathematical Methods for Digital Computers*, pages 191–203. Wiley, 1960.
- [39] K. Engan, S.O. Aase, and J.H. Husøy. Multi-frame compression: Theory and design. *Signal Processing*, 80(10):2121–2140, 2001.
- [40] H.G. Feichtinger and T. Strohmer, editors. *Gabor Analysis and Algorithms, Theory and Applications*. Birkhäuser, 1998.
- [41] H.G. Feichtinger, A. Türk, and T. Strohmer. Hierarchical parallel matching pursuit. In *Proc. SPIE: Image Reconstruction and Restoration*, pages 222–232, 1994.
- [42] S.E. Ferrando, E.J. Doolittle, A.J. Bernal, and L.J. Bernal. Probabilistic matching pursuit with Gabor dictionaries. *Signal Processing*, 80:2099–2120, 2000.
- [43] R. Fourer. Software survey: Linear programming. *OR/MS Today*, June 2005.
- [44] J.H. Friedman and W. Stuetzle. Projection pursuit regression. *Journal of the American Statistical Association*, 76(376), 1981.
- [45] J.H. Friedman and J.W. Tukey. A projection pursuit algorithm for exploratory data analysis. *IEEE Transactions on Computers*, C-23:881–890, 1974.
- [46] J.-J. Fuchs. On sparse representations in arbitrary redundant bases. *IEEE Trans. on Information Theory*, 50(6), 2004.
- [47] D. Gabor. Theory of communication. *Journal of the IEE*, 93:429–457, 1946.
- [48] J.S. Garofolo, L.F. Lamel, W.M. Fisher, J.G. Fiscus, D.S. Pallett, and N.L. Dahlgren. DARPA TIMIT acoustic-phonetic continuous speech corpus CD-ROM. Technical Report NISTIR-4930, US Dept. of Commerce, National Institute of Standards and Technology, 1993.
- [49] M. Gharavi-Alkhansari and T.S. Huang. A fast orthogonal matching pursuit algorithm. In *Proceedings of the 1998 IEEE International Conference on Acoustics, Speech, and Signal Processing*, pages 1389–1392, 1998.
- [50] A.C. Gilbert and J.A. Tropp. Applications of sparse approximations in communications. In *Proceedings of IEEE International Symposium on Information Theory*, 2005.
- [51] I.F. Gorodnitsky and B.D. Rao. Sparse signal reconstruction from limited data using FOCUSS: A re-weighted minimum norm algorithm. *IEEE Trans. on Signal Processing*, 45(3), 1997.
- [52] T.N.E. Greville. The pseudoinverse of a rectangular or singular matrix and its application to the solution of systems of linear equations. *SIAM Review*, 1(1), 1959.
- [53] R. Gribonval and P. Vandergheynst. On the exponential convergence of matching pursuits in quasi-incoherent dictionaries. Technical Report 1619, IRISA, 2004.
- [54] G. Harikumar, C. Couvreur, and Y. Bresler. Fast optimal and suboptimal algorithms for sparse solutions to linear inverse problems. In *Proceedings of the 1998 IEEE International Conference on Acoustics, Speech, and Signal Processing*, pages 1877–1881, 1998.
- [55] A.J. Hoffman. On greedy algorithms that succeed. In I. Anderson, editor, *Surveys in Combinatorics 1985*, pages 97–112, 1985.
- [56] P.J. Huber. Projection pursuit. *The Annals of Statistics*, 13(2):435–475, 1985.

- [57] X. Huo. *Sparse Image Representation via Combined Transforms*. PhD thesis, Stanford University, Department of Statistics, 1999.
- [58] S. Jaggi, W.C. Karl, S. Mallat, and A.S. Willsky. High resolution pursuit for feature extraction. *Applied and Computational Harmonic Analysis*, 5:428–449, 1998.
- [59] R.A. Jarvis. On the identification of the convex hull of a finite set of points in the plane. *Information Processing Letters*, 2:18–21, 1973.
- [60] J. Karvanen and A. Cichocki. Measuring sparseness of noisy signals. In *4th International Symposium on Independent Component Analysis and Blind Signal Separation*, pages 125–130, 2003.
- [61] K. Kreutz-Delgado and B.D. Rao. Measures and algorithms for best basis selection. In *Proceedings of the 1998 IEEE International Conference on Acoustics, Speech, and Signal Processing*, 1998.
- [62] S.A. Leichter, G.B. Dantzig, and J.W. Davis. A strictly improving phase I algorithm using least-squares subproblems. Technical Report SOL 92-1, Systems Optimization Laboratory, Stanford University, 1992.
- [63] M.S. Lewicki. Efficient coding of natural sounds. *Nature Neuroscience*, 5(4):356–363, 2002.
- [64] M.S. Lewicki and T.J. Sejnowski. Learning overcomplete representations. *Neural Computation*, 12:337–365, 2000.
- [65] I.J. Lustig, R.E. Marsten, and D.F. Shanno. On implementing Mehrotra’s predictor-corrector interior point method for linear programming. *SIAM J. Optimization*, 2(3):435–449, 1992.
- [66] D. Malioutov, M. Cetin, and A.S. Willsky. A sparse signal reconstruction perspective on source localization with sensor arrays. *IEEE Transactions on Signal Processing*, 53(8), 2005.
- [67] S.G. Mallat and Z. Zhang. Matching pursuits with time-frequency dictionaries. *IEEE Transactions on Signal Processing*, 41(12), 1993.
- [68] S. Mehrotra. On the implementation of a primal-dual interior point method. *SIAM Journal on Optimization*, 2(4):575–601, 1992.
- [69] A.J. Miller. *Subset Selection in Regression, Second Edition*. CRC Press, 2002.
- [70] S.S. Morgan. A comparison of simplex method algorithms. Master’s thesis, University of Florida, 1997.
- [71] T.L. Morin, N. Prabhu, and Z. Zhang. Complexity of the gravitational method for linear programming. *Journal of Optimization Theory and Applications*, 108(3):633–658, 2001.
- [72] K.G. Murty. The gravitational method for linear programming. *Opsearch*, 23:206–214, 1986.
- [73] B.K. Natarajan. Sparse approximate solutions to linear systems. *SIAM Journal on Computing*, 24(2):227–234, 1995.
- [74] M.S. O’Brien, A.N. Sinclair, and S.M. Kramer. Recovery of a sparse spike time series by  $L_1$  norm deconvolution. *IEEE Transactions on Signal Processing*, 42(12), 1994.
- [75] B.A. Olshausen and D.J. Field. Sparse coding with an overcomplete basis set: A strategy employed by V1? *Vision Research*, 37(23):3311–3325, 1997.

- [76] P.-Q. Pan. A basis-deficiency-allowing variation of the simplex method for linear programming. *Computers and Mathematics with Applications*, 36(3), 1998.
- [77] V. Pan and R. Schreiber. An improved Newton iteration for the generalized inverse of a matrix, with applications. *SIAM J. Sci. Stat. Comput.*, 12(5):1109–1130, 1991.
- [78] Y.C. Pati, R. Rezaifar, and P.S. Krishnaprasad. Orthogonal matching pursuit: Recursive function approximation with applications to wavelet decomposition. In *Twenty-Seventh Asilomar Conference on Signals, Systems & Computers*, 1993.
- [79] P.J. Phillips. Matching pursuit filters applied to face identification. *IEEE Transactions on Image Processing*, 7(8), 1998.
- [80] W.H. Pierce. A self-correcting matrix iteration for the Moore-Penrose generalized inverse. *Linear Algebra and Its Applications*, 244:357–363, 1996.
- [81] S. Qian and D. Chen. Signal representation using adaptive normalized Gaussian functions. *Signal Processing*, 36:1–11, 1994.
- [82] B.D. Rao. Signal processing with the sparseness constraint. In *Proceedings of the 1998 IEEE International Conference on Acoustics, Speech, and Signal Processing*, 1998.
- [83] L. Rebollo-Neira and D. Lowe. Optimized orthogonal matching pursuit approach. *IEEE Signal Processing Letters*, 9(4), 2002.
- [84] J.R. Rice. Experiments on Gram-Schmidt orthogonalization. *Math. Comp.*, 20:325–328, 1966.
- [85] R. Seidel. Constructing higher-dimensional convex hulls at logarithmic cost per face. In *Proc. 18th ACM Symposium on the Theory of Computation*, pages 404–413, 1986.
- [86] R. Seidel. Convex hull computations. In J.E. . E. Goodman and J. O’Rourke, editors, *Handbook of Discrete and Computational Geometry, Second Edition*. CRC Press, 2004.
- [87] A.P. Sethi and G.L. Thompson. The pivot and probe algorithm for solving a linear program. *Mathematical Programming*, 29:219–233, 1984.
- [88] E.P. Simoncelli, W.T. Freeman, E.H. Adelson, and D.J. Heeger. Shiftable multi-scale transforms. *IEEE Transactions on Information Theory*, 38(2), 1992.
- [89] K. Skretting. Pre-stack/post-stack seismic data from the North Sea. <http://www.ux.his.no/karlsk/sdata/>.
- [90] T. Söderström and G.W. Stewart. On the numerical properties of an iterative method for computing the Moore-Penrose generalized inverse. *SIAM J. Numer. Anal.*, 11(1), 1974.
- [91] D.A. Spielman and S. Teng. Smoothed analysis of algorithms: Why the simplex algorithm usually takes polynomial time. In *STOC’01*, pages 296–305, 2001.
- [92] J.L. Starck, M. Elad, and D. Donoho. Redundant multiscale transforms and their application for morphological component separation. *Advances in Imaging and Electron Physics*, 2004.
- [93] J.J. Stone. The cross-section method. Technical Report P-1490, The RAND Corporation, 1958.
- [94] G. Swart. Finding the convex hull facet by facet. *Journal of Algorithms*, 6:17–48, 1985.

- [95] V.N. Temlyakov. Greedy algorithms and  $m$ -term approximation with regard to redundant dictionaries. *Journal of Approximation Theory*, 98(1):117–145, 1999.
- [96] T. Terlaky and S. Zhang. A survey on pivot rules for linear programming. Technical Report 91-99, Delft University of Technology, Faculty of Technical Mathematics and Informatics, 1991.
- [97] J.A. Tropp. Greed is good: Algorithmic results for sparse approximation. *IEEE Transactions on Information Theory*, 50(10), 2004.
- [98] R.J. Vanderbei. *Linear Programming: Foundations and Extensions, Second Edition*. Kluwer Academic Publishers, 2001.
- [99] Y. Zhang. User's guide to LIPSOL Linear-programming Interior Point SOLvers v0.4. *Optimization Methods and Software*, 11-12, 1999.

---

**Algorithm 1** Greedy Basis Pursuit

---

**Input**

- A signal  $\mathbf{x} \in \mathbb{R}^d$
- A dictionary  $\mathcal{D} = \{\psi_i\}_{i=1}^n$ ,
- A threshold  $\epsilon \geq 0$

**Output**

A representation of  $\mathbf{x}$ , consisting of

- A set of indices  $\mathcal{I} \subseteq \{1, \dots, n\}$
- A set of coefficients  $\mathcal{A} = \{\alpha_i\}_{i \in \mathcal{I}}$

such that  $\mathbf{x} - \sum_{i \in \mathcal{I}} \alpha_i \psi_i < \epsilon$

**Procedure**

1. Initialize
  - (a) Select the first atom
$$k \leftarrow \arg \max_{i \in \{1, \dots, n\}} \langle \mathbf{x}, \psi_i \rangle$$
  - (b) Compute the initial approximation
$$\alpha_k \leftarrow \langle \mathbf{x}, \psi_k \rangle, \quad \mathcal{I}^{(0)} \leftarrow \{k\}, \quad \mathcal{A}^{(0)} \leftarrow \{\alpha_k\}$$
  - (c) Initialize the biorthogonal system
$$\tilde{\Psi}^\perp \leftarrow \{\psi_k\}$$
  - (d) Initialize the hyperplane
$$\tilde{\mathbf{x}}^{(0)} \leftarrow \alpha_k \psi_k, \quad \mathbf{n} \leftarrow \mathbf{x} / \|\mathbf{x}\|, \quad \mathbf{r} \leftarrow \mathbf{x} - \tilde{\mathbf{x}}$$
2. Repeat until  $\|\mathbf{r}\| < \epsilon$ 
  - (a) Compute the center and plane of rotation
$$\tilde{\mathbf{x}}_H \leftarrow (\langle \psi_i, \mathbf{n} \rangle / \langle \tilde{\mathbf{x}}, \mathbf{n} \rangle) \tilde{\mathbf{x}}, \quad i \in \mathcal{I}$$
$$\mathbf{v} \leftarrow (\mathbf{r} - \langle \mathbf{r}, \mathbf{n} \rangle \mathbf{n}) / \|\mathbf{r} - \langle \mathbf{r}, \mathbf{n} \rangle \mathbf{n}\|$$
  - (b) Project atoms into the  $\mathbf{n}$ - $\mathbf{v}$ -plane and select the next atom
$$k \leftarrow \arg \min_{i \in \{1, \dots, n\}} \tan^{-1} \frac{\langle \psi_i - \tilde{\mathbf{x}}_H, \mathbf{n} \rangle}{\langle \psi_i - \tilde{\mathbf{x}}_H, \mathbf{v} \rangle}$$
  - (c) Compute the new representation and update the biorthogonal system
$$\{\mathcal{I}, \mathcal{A}, \tilde{\Psi}^\perp\} \leftarrow \text{AddAtom}(\mathbf{x}, \mathcal{I}, \mathcal{A}, \psi_k, \tilde{\Psi}^\perp)$$
  - (d) Discard any extraneous atoms
$$\mathbf{while} \quad \exists \alpha_i \leq 0, i \in \mathcal{I} \quad \mathbf{do}$$
$$\{\mathcal{I}, \mathcal{A}, \tilde{\Psi}^\perp\} \leftarrow \text{SubtractAtom}(\mathbf{x}, \mathcal{I}, \mathcal{A}, \psi_j, \tilde{\Psi}^\perp)$$
  - (e) Update the hyperplane parameters
$$\tilde{\mathbf{x}} \leftarrow \sum_{i \in \mathcal{I}} \alpha_i \psi_i$$
$$\mathbf{n} \leftarrow -\frac{\langle \psi_k - \tilde{\mathbf{x}}_H, \mathbf{n} \rangle \mathbf{v} + \langle \psi_k - \tilde{\mathbf{x}}_H, \mathbf{v} \rangle \mathbf{n}}{\langle \psi_k - \tilde{\mathbf{x}}_H, \mathbf{n} \rangle \mathbf{v} + \langle \psi_k - \tilde{\mathbf{x}}_H, \mathbf{v} \rangle \mathbf{n}}$$
$$\mathbf{r} \leftarrow \mathbf{x} - \tilde{\mathbf{x}}$$

---

**Algorithm 2** AddAtom

---

**Input**

- The signal  $\mathbf{x}$
- The current representation  $\mathcal{I}, \mathcal{A}$
- The atom to add  $k, \psi_k$
- The current biorthogonal vectors  $\tilde{\Psi}^\perp$

**Output**

- The updated representation  $\mathcal{I}, \mathcal{A}$
- The updated biorthogonal vectors  $\tilde{\Psi}^\perp$

**Procedure**

1. Compute the new biorthogonal vector  
 $\forall i \in \mathcal{I}, \beta_i \leftarrow \langle \tilde{\psi}_i^\perp, \psi_k \rangle$   
 $\psi_k^\perp \leftarrow \psi_k - \sum_{i \in \mathcal{I}} \beta_i \psi_i$   
 $\tilde{\psi}_k^\perp \leftarrow \psi_k^\perp / \|\psi_k^\perp\|_2$
  2. Update the biorthogonal system  
 $\forall i \in \mathcal{I}, \tilde{\psi}_i^\perp \leftarrow \tilde{\psi}_i^\perp - \beta_i \tilde{\psi}_k^\perp$   
 $\tilde{\Psi}^\perp \leftarrow \tilde{\Psi}^\perp \cup \{\tilde{\psi}_k^\perp\}$
  3. Update the representation  
 $\alpha_k \leftarrow \langle \mathbf{x}, \tilde{\psi}_k^\perp \rangle$   
 $\forall i \in \mathcal{I}, \alpha_i \leftarrow \alpha_i - \beta_i \alpha_k$   
 $\mathcal{I} \leftarrow \mathcal{I} \cup \{k\}$   
 $\mathcal{A} \leftarrow \mathcal{A} \cup \{\alpha_k\}$
-



---

**Algorithm 3** SubtractAtom

---

**Input**

- The signal  $\mathbf{x}$
- The current representation  $\mathcal{I}, \mathcal{A}$
- The index of the atom to subtract  $k$
- The current biorthogonal vectors  $\tilde{\Psi}^\perp$

**Output**

- The updated representation  $\mathcal{I}, \mathcal{A}$ ,
- The updated biorthogonal vectors  $\tilde{\Psi}^\perp$

**Procedure**

1. Delete the atom from the representation  
 $\mathcal{I} \leftarrow \mathcal{I} - \{k\}$   
 $\mathcal{A} \leftarrow \mathcal{A} - \{\alpha_k\}$
  2. Update the biorthogonal system  
 $\tilde{\Psi}^\perp \leftarrow \tilde{\Psi}^\perp - \{\tilde{\psi}_k^\perp\}$   
 $\forall i \in \mathcal{I}, \gamma_i \leftarrow \langle \tilde{\psi}_k^\perp, \tilde{\psi}_i^\perp \rangle / \|\tilde{\psi}_k^\perp\|_2^2$   
 $\forall i \in \mathcal{I}, \tilde{\psi}_i^\perp \leftarrow \tilde{\psi}_i^\perp - \gamma_i \tilde{\psi}_k^\perp$
  3. Update the representation  
 $\forall i \in \mathcal{I}, \alpha_i \leftarrow \alpha_i - \alpha_k \gamma_i$
-