

# Fast Bayesian Matching Pursuit

Philip Schniter, Lee C. Potter, and Justin Ziniel

Dept. of ECE, The Ohio State University, Columbus, OH 43210.

e-mail: schniter@ece.osu.edu, potter@ece.osu.edu, ziniel.1@osu.edu.

**Abstract**—A low-complexity recursive procedure is presented for minimum mean squared error (MMSE) estimation in linear regression models. A Gaussian mixture is chosen as the prior on the unknown parameter vector. The algorithm returns both an approximate MMSE estimate of the parameter vector and a set of high posterior probability mixing parameters. Emphasis is given to the case of a sparse parameter vector. Numerical simulations demonstrate estimation performance and illustrate the distinctions between MMSE estimation and MAP model selection. The set of high probability mixing parameters not only provides MAP basis selection, but also yields relative probabilities that reveal potential ambiguity in the sparse model.<sup>1</sup>

## I. INTRODUCTION

Sparse linear regression is a topic of long standing interest in statistics and signal processing. The linear regression model is

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \boldsymbol{\nu}, \quad (1)$$

with unknown parameter vector  $\mathbf{x}$ , unit norm columns in the regressor matrix  $\mathbf{A}$ , and additive noise  $\boldsymbol{\nu}$ . We provide a brief, and necessarily incomplete, survey of existing approaches, with an emphasis on themes relevant to the proposed estimator.

Algorithmic approaches have been proposed over several decades, providing greedy heuristic solutions. Examples include CLEAN [1], iteratively re-weighted least-squares [2], and orthogonal matched pursuit (OMP) [3]. Tropp and Gilbert [4], for example, provide sufficient conditions on the sparsity of  $\mathbf{x}$  and correlation among columns of  $\mathbf{A}$  such that the greedy OMP provides correct model selection with high probability in the noiseless measurement case.

In addition to greedy approaches, penalized least-squares solutions for  $\mathbf{x}$  have likewise been presented in the past four decades. In this class of approaches, parameters are found via the optimization

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{x}\|_p^p, \quad (2)$$

or, equivalently for some  $\epsilon > 0$

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \|\mathbf{x}\|_p \quad \text{s.t.} \quad \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 < \epsilon. \quad (3)$$

Ridge regression [5] (Tikhonov regularization) adopts  $p = 2$ , while basis pursuit [6] and Lasso [7] use  $p = 1$ . Equation (2) has been widely adopted, for example in radar imaging [8], image reconstruction [9], [10], and elsewhere [11], [12]. Elegant

recent results by several authors [13]–[15] have demonstrated sufficient conditions on  $\mathbf{A}$ ,  $\boldsymbol{\nu}$ , and sparsity of  $\mathbf{x}$  such that the convex problem in (3) for  $p = 1$  provides the unique solution to the non-convex task

$$\min \|\mathbf{x}\|_0 \quad \text{s.t.} \quad \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 < \epsilon. \quad (4)$$

These proofs have validated the widespread use of (2)–(3), providing a deeper understanding, spurring a resurgent interest, and promoting the interpretation as “compressive sensing.” The large class of methods adopting (2) may be interpreted as implicitly seeking the Bayesian MAP estimate of  $\mathbf{x}$  under a sparsity inducing prior

$$p(\mathbf{x}) \propto \exp \left\{ -\frac{\lambda}{2} \|\mathbf{x}\|_p^p \right\}. \quad (5)$$

The method of sparse Bayesian learning [16], [17] explicitly adopts a Bayesian framework with  $x_i$  independent, zero-mean, Gaussian with unknown variance  $\sigma_i^2$ . The unknown variances are given the Gamma conjugate prior, and an expectation-maximization (EM) iteration computes a MAP estimate of  $\mathbf{x}$ .

In the literature, primary focus is placed on the detection of the few significant entries of the sparse  $\mathbf{x}$ —a task alternatively known as model selection or basis selection. In contrast, we adopt a minimum mean-squared error (MMSE) estimation formulation and focus on accurately inferring  $\mathbf{x}$  from the noisy observations,  $\mathbf{y}$ . The MMSE estimation approach was likewise adopted in a crisp exposition by Larsson and Selén [18]; as they noted, the Bayesian formulation requires *a priori* assumptions that are explicitly stated and admit unambiguous interpretation. (We specifically identify similarities to [18] in Section V.)

As a byproduct of approximating the proposed MMSE estimation algorithm, we also provide exact ratios of posterior probabilities for a set of high probability solutions to the detection problem. These relative probabilities serve to reveal potential ambiguity among multiple models, due to low signal-to-noise ratio and/or significant correlation among columns in the regressor matrix,  $\mathbf{A}$ .

The remainder of the paper is organized as follows. In Section II, we state the signal model and explicitly identify the assumed priors. In Section III, we describe our proposed technique, and in Section IV we investigate its performance numerically. In Section V, we discuss our findings, and in Section VI we conclude.

<sup>1</sup>This work was supported by the National Science Foundation under Grant 0237037, the Office of Naval Research grant N00014-07-1-0209, the Wright Brothers Institute grant GRT00009715, and AFOSR under award FA9550-06-1-0324.

## II. SIGNAL MODEL

Consider observing  $\mathbf{y} \in \mathbb{R}^M$ , a noisy linear combination of the parameters in  $\mathbf{x} \in \mathbb{R}^N$ :

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \boldsymbol{\nu}, \quad (6)$$

where the noise  $\boldsymbol{\nu}$  is assumed to be white Gaussian with variance  $\sigma^2$ , i.e.,  $\boldsymbol{\nu} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_M)$ , and the columns of  $\mathbf{A}$  are taken to be unit-norm. Our focus is on the under-determined case (i.e.,  $N \gg M$ ) with a suitably *sparse* parameter vector  $\mathbf{x}$  (i.e.,  $\|\mathbf{x}\|_0 \ll N$ ).

To model sparsity, we assume that the parameters are generated from a Gaussian mixture density:

$$\mathbf{x}|\mathbf{s} \sim \mathcal{N}(\mathbf{0}, \mathbf{R}(\mathbf{s})), \quad (7)$$

where the covariance  $\mathbf{R}(\mathbf{s})$  is determined by a discrete random vector  $\mathbf{s} = [s_0, \dots, s_{N-1}]^T$  of mixture parameters. For simplicity, we take  $\mathbf{R}(\mathbf{s})$  to be diagonal with  $[\mathbf{R}(\mathbf{s})]_{n,n} = \sigma_{s_n}^2$ , implying that  $\{x_n|s_n\}_{n=0}^{N-1}$  are independent with  $x_n|s_n \sim \mathcal{N}(0, \sigma_{s_n}^2)$ . Also for simplicity, we assume that the mixture parameters  $\{s_n\}_{n=0}^{N-1}$  are  $^2$  Bernoulli( $p_1$ ). To model sparse  $\mathbf{x}$ , we choose  $\sigma_0^2 = 0$  and  $p_1 \ll 1$ .

From the model assumptions it can be seen that

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{x} \end{bmatrix} | \mathbf{s} \sim \mathcal{N} \left( \mathbf{0}, \begin{bmatrix} \Phi(\mathbf{s}) & \mathbf{A}\mathbf{R}(\mathbf{s}) \\ \mathbf{R}(\mathbf{s})\mathbf{A}^T & \mathbf{R}(\mathbf{s}) \end{bmatrix} \right), \quad (8)$$

where

$$\Phi(\mathbf{s}) := \mathbf{A}\mathbf{R}(\mathbf{s})\mathbf{A}^T + \sigma^2 \mathbf{I}_M. \quad (9)$$

## III. ESTIMATION OF BASIS AND PARAMETERS

In this section, we propose an efficient search procedure to find the most probable basis configurations along with their respective posterior probabilities. These posteriors can then be used to compute an MMSE estimate of the sparse parameters  $\mathbf{x}$ .

### A. Basis Selection Metric

As a consequence of the model described in Section II, the nonzero locations in  $\mathbf{s}$  specify which of the basis elements (i.e., columns of  $\mathbf{A}$ ) are “active.” Thus, basis selection reduces to estimation of  $\mathbf{s}$ . Because we have adopted a probabilistic model for  $\{\mathbf{s}, \mathbf{y}\}$ , we can not only compute *which* of the basis configurations are most likely, but also *how likely* these bases are. The latter is accomplished through the estimation of dominant posteriors  $p(\mathbf{s}|\mathbf{y})$ .

The posterior can be written, via Bayes rule, as

$$p(\mathbf{s}|\mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{s})p(\mathbf{s})}{\sum_{\mathbf{s}' \in \mathcal{S}} p(\mathbf{y}|\mathbf{s}')p(\mathbf{s}')}, \quad (10)$$

where  $\mathcal{S} = \{0, 1\}^N$ , which shows that estimating  $\{p(\mathbf{s}|\mathbf{y})\}_{\mathbf{s} \in \mathcal{S}}$  reduces to estimating  $\{p(\mathbf{y}|\mathbf{s})p(\mathbf{s})\}_{\mathbf{s} \in \mathcal{S}}$ . The size of  $\mathcal{S}$  makes it impractical to compute  $\{p(\mathbf{s}|\mathbf{y})\}$  or  $\{p(\mathbf{y}|\mathbf{s})p(\mathbf{s})\}$  for all  $\mathbf{s} \in \mathcal{S}$ . However, the set  $\mathcal{S}_*$  responsible for the *dominant* posteriors

<sup>2</sup>In other words,  $s_n$  is binary with  $\Pr\{s_n = 1\} = p_1$  and  $\Pr\{s_n = 0\} = 1 - p_1$ .

can be quite small and therefore practical to compute. Working in the log domain, we find

$$\mu(\mathbf{s}) := \ln p(\mathbf{y}|\mathbf{s})p(\mathbf{s}) \quad (11)$$

$$= \ln p(\mathbf{y}|\mathbf{s}) + \sum_{n=0}^{N-1} \ln p(s_n) \quad (12)$$

$$= \ln p(\mathbf{y}|\mathbf{s}) + \|\mathbf{s}\|_0 \ln p_1 + (N - \|\mathbf{s}\|_0) \ln(1 - p_1) \quad (13)$$

$$= \ln p(\mathbf{y}|\mathbf{s}) + \|\mathbf{s}\|_0 \ln \frac{p_1}{1-p_1} + N \ln(1 - p_1) \quad (14)$$

$$= -\frac{M}{2} \ln 2\pi - \frac{1}{2} \ln \det(\Phi(\mathbf{s})) - \frac{1}{2} \mathbf{y}^T \Phi(\mathbf{s})^{-1} \mathbf{y} + \|\mathbf{s}\|_0 \ln \frac{p_1}{1-p_1} + N \ln(1 - p_1). \quad (15)$$

We will refer to  $\mu(\mathbf{s})$  as the *basis selection metric*.

### B. MMSE Parameter Estimation

For applications in which the identification of the most probable basis is the primary objective, the sparse coefficients  $\mathbf{x}$  can be regarded as nuisance parameters. In other applications, however, estimation of  $\mathbf{x}$  is the primary goal.

The MMSE estimate of  $\mathbf{x}$  from  $\mathbf{y}$  is

$$\hat{\mathbf{x}}_{\text{mmse}} := \mathbb{E}\{\mathbf{x}|\mathbf{y}\} = \sum_{\mathbf{s} \in \mathcal{S}} p(\mathbf{s}|\mathbf{y}) \mathbb{E}\{\mathbf{x}|\mathbf{y}, \mathbf{s}\} \quad (16)$$

where from (8) it is straightforward (e.g., [19, p. 155]) to obtain

$$\mathbb{E}\{\mathbf{x}|\mathbf{y}, \mathbf{s}\} = \mathbf{R}(\mathbf{s})\mathbf{A}^T \Phi(\mathbf{s})^{-1} \mathbf{y}. \quad (17)$$

Although exact evaluation of (16) involves a summation over  $2^N$  terms, which may be computationally infeasible, the MMSE estimate can be closely approximated using only the dominant posteriors:

$$\hat{\mathbf{x}}_{\text{ammse}} := \sum_{\mathbf{s} \in \mathcal{S}_*} p(\mathbf{s}|\mathbf{y}) \mathbb{E}\{\mathbf{x}|\mathbf{y}, \mathbf{s}\}. \quad (18)$$

Likewise, the covariance of the corresponding estimation error can be closely approximated as

$$\text{Cov}\{\mathbf{x}|\mathbf{y}\} \approx \sum_{\mathbf{s} \in \mathcal{S}_*} p(\mathbf{s}|\mathbf{y}) [\text{Cov}\{\mathbf{x}|\mathbf{y}, \mathbf{s}\} + (\hat{\mathbf{x}}_{\text{ammse}} - \mathbb{E}\{\mathbf{x}|\mathbf{y}, \mathbf{s}\})(\hat{\mathbf{x}}_{\text{ammse}} - \mathbb{E}\{\mathbf{x}|\mathbf{y}, \mathbf{s}\})^T] \quad (19)$$

$$\text{Cov}\{\mathbf{x}|\mathbf{y}, \mathbf{s}\} = \mathbf{R}(\mathbf{s}) - \mathbf{R}(\mathbf{s})\mathbf{A}^T \Phi(\mathbf{s})^{-1} \mathbf{A}\mathbf{R}(\mathbf{s}). \quad (20)$$

Note that, in evaluating (18)-(20), the primary challenge becomes that of obtaining  $p(\mathbf{s}|\mathbf{y})$  and  $\Phi(\mathbf{s})^{-1}$  for each  $\mathbf{s} \in \mathcal{S}_*$ . In the sequel, we propose a fast algorithm to search for the dominant basis configurations  $\mathcal{S}_*$  that, as a byproduct, also generates  $p(\mathbf{s}|\mathbf{y})$  and  $\Phi(\mathbf{s})^{-1}$  for each of the  $\mathbf{s}$  returned by the search.

### C. Bayesian Matching Pursuit

We now describe an efficient means of determining  $\mathcal{S}_*$ , the set of mixture parameters  $\mathbf{s}$  yielding the dominant values of  $p(\mathbf{s}|\mathbf{y})$ , or, equivalently, the dominant values of  $\mu(\mathbf{s})$ . First we present a prosaic description of the search heuristic; the detailed algorithm will be specified in Section III-E.

The search starts with  $\mathbf{s} = \mathbf{0}$  and first “turns on” one mixture parameter at a time, yielding a set of  $N$  binary vectors  $\mathbf{s}$  which

we refer to as  $\mathcal{S}^{(1)}$ . The metrics  $\mu(\mathbf{s})$  for the  $N$  vectors in  $\mathcal{S}^{(1)}$  are then computed, and the elements of  $\mathcal{S}^{(1)}$  with the  $D$  largest metrics are collected in  $\mathcal{S}_*^{(1)}$ . For each candidate in  $\mathcal{S}_*^{(1)}$ , all locations of a second active mixture parameter are then considered, yielding  $(N-1) + (N-2) + \dots + (N-D) = ND - \frac{D(D+1)}{2}$  unique binary vectors to store in  $\mathcal{S}^{(2)}$ . The metrics  $\mu(\mathbf{s})$  for all vectors in  $\mathcal{S}^{(2)}$  are then computed, and the elements of  $\mathcal{S}^{(2)}$  with the  $D$  largest metrics are collected in  $\mathcal{S}_*^{(2)}$ . Then, for each candidate vector in  $\mathcal{S}_*^{(2)}$ , all possibilities of a third active mixture parameter are considered, and those with the  $D$  largest metrics are stored in  $\mathcal{S}_*^{(3)}$ . The process continues until  $\mathcal{S}_*^{(P)}$  is computed, where  $P$  can be chosen<sup>3</sup> to make  $\Pr(\|\mathbf{s}\|_0 > P)$  sufficiently small.<sup>4</sup> Note that  $\mathcal{S}_*^{(P)}$  constitutes the algorithm's final estimate of  $\mathcal{S}_*$ . Henceforth we denote this final estimate by  $\hat{\mathcal{S}}_*$ .

#### D. Fast Metric Update

For use with the aforementioned Bayesian matching pursuit (BMP) algorithm, we propose a fast metric update which computes the change in  $\mu(\cdot)$  that results from the activation of a single mixture parameter. More precisely, if we denote by  $\mathbf{s}_n$  the vector identical to  $\mathbf{s}$  except for the  $n^{\text{th}}$  coefficient, which is active in  $\mathbf{s}_n$  but inactive in  $\mathbf{s}$  (i.e.,  $[\mathbf{s}_n]_n = 1$  and  $[\mathbf{s}]_n = 0$ ), then we seek an efficient method of computing  $\Delta_n(\mathbf{s}) := \mu(\mathbf{s}_n) - \mu(\mathbf{s})$ . Note that the metric at the root node (i.e.,  $\mathbf{s} = \mathbf{0}$ ) is

$$\mu(\mathbf{0}) = -\frac{M}{2} \ln 2\pi - \frac{M}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} \|\mathbf{y}\|_2^2 + N \ln(1 - p_1) \quad (21)$$

via (15) and the fact that  $\Phi(\mathbf{0}) = \sigma^2 \mathbf{I}_M$ .

To derive the fast metric update, we start with the property that, for any  $n$  and  $\mathbf{s}$ ,

$$\Phi(\mathbf{s}_n) = \Phi(\mathbf{s}) + \sigma_1^2 \mathbf{a}_n \mathbf{a}_n^T, \quad (22)$$

from which the matrix inversion lemma implies

$$\Phi(\mathbf{s}_n)^{-1} = \Phi(\mathbf{s})^{-1} - \sigma_1^2 \beta_n \mathbf{b}_n \mathbf{b}_n^T \quad (23)$$

$$\mathbf{b}_n := \Phi(\mathbf{s})^{-1} \mathbf{a}_n \quad (24)$$

$$\beta_n := (1 + \sigma_1^2 \mathbf{a}_n^T \mathbf{b}_n)^{-1}. \quad (25)$$

Equations (22)-(25) then imply

$$\mathbf{y}^T \Phi(\mathbf{s}_n)^{-1} \mathbf{y} = \mathbf{y}^T (\Phi(\mathbf{s})^{-1} - \sigma_1^2 \beta_n \mathbf{b}_n \mathbf{b}_n^T) \mathbf{y} \quad (26)$$

$$= \mathbf{y}^T \Phi(\mathbf{s})^{-1} \mathbf{y} - \sigma_1^2 \beta_n (\mathbf{y}^T \mathbf{b}_n)^2 \quad (27)$$

$$\ln \det(\Phi(\mathbf{s}_n)) = \ln \det(\Phi(\mathbf{s}) + \sigma_1^2 \mathbf{a}_n \mathbf{a}_n^T) \quad (28)$$

$$= \ln \left[ (1 + \sigma_1^2 \mathbf{a}_n^T \Phi(\mathbf{s})^{-1} \mathbf{a}_n) \det(\Phi(\mathbf{s})) \right] \quad (29)$$

$$\|\mathbf{s}_n\|_0 \ln \frac{p_1}{1-p_1} = (\|\mathbf{s}\|_0 + 1) \ln \frac{p_1}{1-p_1} \quad (30)$$

$$= \|\mathbf{s}\|_0 \ln \frac{p_1}{1-p_1} + \ln \frac{p_1}{1-p_1}, \quad (31)$$

<sup>3</sup>One could also determine the stopping parameter  $P$  adaptively.

<sup>4</sup>Notice that  $\|\mathbf{s}\|_0$  is Binomial( $N, p_1$ ) distribution. When  $Np_1 > 5$ , it is common to use the approximation  $\|\mathbf{s}\|_0 \sim \mathcal{N}(Np_1, Np_1(1-p_1))$ , in which case  $\Pr(\|\mathbf{s}\|_0 > P) = \frac{1}{2} \operatorname{erfc} \left( \frac{P - Np_1}{\sqrt{2Np_1(1-p_1)}} \right)$ .

which, combined with (15), yield

$$\mu(\mathbf{s}_n) = \mu(\mathbf{s}) + \underbrace{\frac{1}{2} \ln \beta_n + \frac{\sigma_1^2}{2} \beta_n (\mathbf{y}^T \mathbf{b}_n)^2 + \ln \frac{p_1}{1-p_1}}_{\Delta_n(\mathbf{s})}. \quad (32)$$

In summary,  $\Delta_n(\mathbf{s})$  in (32) quantifies the change in our basis selection metric  $\mu(\cdot)$  due to the activation of the  $n^{\text{th}}$  tap of  $\mathbf{s}$ .

#### E. Fast Bayesian Matching Pursuit

Notice that the cost of computing  $\{\beta_n\}_{n=0}^{N-1}$  via (24)-(25) is  $\mathcal{O}(NM^2)$  if standard matrix multiplication is used. As we now describe, the complexity of this operation can be made linear in  $M$  by exploiting the structure of  $\Phi(\mathbf{s})^{-1}$ .

Say that  $\mathbf{t} = [t_1, t_2, \dots, t_p]^T$  contains the indices of active elements in  $\mathbf{s}$ . Then, from (23),

$$\Phi(\mathbf{s})^{-1} = \frac{1}{\sigma^2} \mathbf{I}_M - \sigma_1^2 \sum_{i=1}^p \beta^{(i)} \mathbf{b}^{(i)} \mathbf{b}^{(i)T}, \quad (33)$$

where  $\mathbf{b}^{(i)}$  and  $\beta^{(i)}$  denote the values of  $\mathbf{b}$  and  $\beta$  generated while activating index  $t_i$  in the mixture parameter vector defined by the active indices  $[t_1, \dots, t_{i-1}]$ . From (24), we are required to compute

$$\mathbf{b}_n = \frac{1}{\sigma^2} \mathbf{a}_n - \sigma_1^2 \sum_{i=1}^p \beta^{(i)} \mathbf{b}^{(i)} \underbrace{\mathbf{b}^{(i)T} \mathbf{a}_n}_{:= c_n^{(i)}} \quad (34)$$

when activating the  $n^{\text{th}}$  tap in  $\mathbf{s}$ . The key observation is that the coefficients  $\{c_n^{(i)}\}_{n=0}^{N-1}$  need only be computed once, i.e., when index  $t_i$  is activated. Furthermore,  $\{c_n^{(i)}\}_{n=0}^{N-1}$  only need to be computed for *surviving* indices  $t_i$ . These tricks form the foundation of the Fast Bayesian Matching Pursuit (FBMP) algorithm outlined in Table I. From the table, it is straightforward to verify that the number of multiplications required by the algorithm is  $\mathcal{O}(NMPD)$ .

## IV. NUMERICAL EXPERIMENTS

### A. FBMP Behavior

Numerical experiments were conducted to investigate the performance of FBMP from Table I for various values of model and algorithmic parameters, and the results are reported in Figs. 1-5. Unless otherwise noted, the experiments used  $N = 256$ ,  $M = 64$ ,  $\text{SNR} = 15$  dB,  $p_1 = 0.04$ , and  $P = \lceil \operatorname{erfc}^{-1}(2P_0) \sqrt{2Np_1(1-p_1)} + Np_1 \rceil$  where  $P_0 = 0.00005$  is the target value of  $\Pr\{\|\mathbf{s}\|_0 > P\}$  as suggested in Section III-C. Here we use  $\text{SNR} := \frac{\sigma_1^2 p_1 N}{\sigma^2 M}$ , as motivated by the unit-norm assumption on the columns of  $\mathbf{A}$ . The plots represent an average of 200 independent model realizations. For each realization of  $\mathbf{A}$ , an i.i.d. normal matrix was drawn and then scaled to make each of its columns unit-norm. Note that the average number of active coefficients  $\mathbb{E}\{\|\mathbf{x}\|_0\} = p_1 N$  is approximately equal to 10 when  $p_1 = 0.04$  and  $N = 256$ . When referring to the ‘‘normalized mean-squared error’’ (NMSE) of an estimate  $\hat{\mathbf{x}}$ , we mean  $\mathbb{E}\{\|\hat{\mathbf{x}} - \mathbf{x}\|_2^2 / \|\mathbf{x}\|_2^2\}$ .

In Fig. 1, we plot NMSE versus observation length  $M$  for FBMP under several values of the search parameter  $D$ . Recall that  $D$  effects a tradeoff between search accuracy and search complexity (the latter of which is expected to grow linearly

```

 $\mu_{0,1} = -\frac{M}{2} \ln 2\pi - \frac{M}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} \|\mathbf{y}\|_2^2 + N \ln(1 - p_1);$ 
for  $n = 1 : N,$ 
   $\tilde{\mathbf{b}}_{1,n} = \sigma^{-2} \mathbf{a}_n;$ 
   $\tilde{\beta}_{1,n} = (1 + \sigma_1^2 \mathbf{a}_n^T \tilde{\mathbf{b}}_{1,n})^{-1};$ 
   $\tilde{\mu}_{1,n} = \mu_{0,1} + \frac{1}{2} \log \tilde{\beta}_{1,n} + \frac{\sigma_1^2}{2} \tilde{\beta}_{1,n} (\mathbf{y}^T \tilde{\mathbf{b}}_{1,n})^2 + \log \frac{p_1}{1-p_1};$ 
end
for  $q = 1 : D,$ 
   $n_* = n$  corresponding to  $q^{\text{th}}$  largest  $\tilde{\mu}_{1,n};$ 
   $\mu_{1,q} = \tilde{\mu}_{1,n_*};$ 
   $\mathbf{b}_{1,q}^{(1)} = \tilde{\mathbf{b}}_{1,n_*}; \mathbf{c}_{1,q}^{(1)} = \mathbf{A}^T \mathbf{b}_{1,q}^{(1)}; \beta_{1,q}^{(1)} = \tilde{\beta}_{1,n_*}; t_{1,q}^{(1)} = n_*;$ 
end
for  $p = 2 : P,$ 
  for  $d = 1 : D,$ 
    for  $n = 1 : N,$ 
       $\tilde{\mathbf{b}}_{d,n} = \sigma^{-2} \mathbf{a}_n - \sum_{i=1}^{p-1} \mathbf{b}_{p-1,d}^{(i)} \beta_{p-1,d}^{(i)} [\mathbf{c}_{p-1,d}^{(i)}]_n;$ 
       $\tilde{\beta}_{d,n} = (1 + \sigma_1^2 \mathbf{a}_n^T \tilde{\mathbf{b}}_{d,n})^{-1};$ 
       $\tilde{\mu}_{d,n} = \mu_{p-1,d} + \frac{1}{2} \log \tilde{\beta}_{d,n} + \frac{\sigma_1^2}{2} \tilde{\beta}_{d,n} (\mathbf{y}^T \tilde{\mathbf{b}}_{d,n})^2 + \log \frac{p_1}{1-p_1};$ 
      if  $n \in \{t_{p-1,d}^{(i)}\}_{i=1}^{p-1}$  then  $\tilde{\mu}_{d,n} = -\infty;$ 
    end
  end
  for  $q = 1 : D,$ 
     $\{d_*, n_*\} = \{d, n\}$  corresponding to  $q^{\text{th}}$  largest  $\tilde{\mu}_{d,n};$ 
     $\mu_{p,q} = \tilde{\mu}_{d_*, n_*};$ 
     $\mathbf{b}_{p,q}^{(p)} = \tilde{\mathbf{b}}_{d_*, n_*}; \mathbf{c}_{p,q}^{(p)} = \mathbf{A}^T \mathbf{b}_{p,q}^{(p)}; \beta_{p,q}^{(p)} = \tilde{\beta}_{d_*, n_*}; t_{p,q}^{(p)} = n_*;$ 
    for  $i = 1 : p-1,$ 
       $\mathbf{b}_{p,q}^{(i)} = \mathbf{b}_{p-1,q}^{(i)}; \mathbf{c}_{p,q}^{(i)} = \mathbf{c}_{p-1,q}^{(i)}; \beta_{p,q}^{(i)} = \beta_{p-1,q}^{(i)}; t_{p,q}^{(i)} = t_{p-1,q}^{(i)};$ 
    end
  end
end
end

```

TABLE I  
FAST BAYESIAN MATCHING PURSUIT

in  $D$ ). There we see that NMSE performance improves as  $M$  gets larger, i.e., as the average number of unknown parameters per measurement  $\frac{p_1 N}{M}$  decreases. For  $D = 1$  (i.e., the simplest possible search), Fig. 1 shows a “knee” in the curve at  $M = 64$  (i.e.,  $\frac{p_1 N}{M} = 0.16$ ) below which NMSE degrades quickly. By increasing search complexity  $D$ , the knee shifts so that the FBMP is robust to a wider range of  $M$  (e.g.,  $M = 48$  or  $\frac{p_1 N}{M} = 0.21$  when  $D = 5$ ). The benefits of increased  $D$  diminish quickly for  $D > 5$ , however.

In Fig. 2, we plot the number of active basis elements missed from FBMP’s estimate of the MAP basis configuration:

$$\hat{\mathbf{s}}_{\text{map}} := \operatorname{argmax}_{\mathbf{s} \in \hat{\mathcal{S}}_*} p(\mathbf{s} | \mathbf{y}). \quad (35)$$

In particular, the traces in Fig. 2 show number-of-misses versus observation length  $M$  for FBMP under several values of search parameter  $D$ . Because the number-of-misses in Fig. 2 closely parallel the NMSEs in Fig. 1, we conjecture that the sub-optimality of FBMP’s greedy search is to blame for the relatively large NMSE values that occur when  $M < 64$  (i.e., when  $\frac{p_1 N}{M} > 0.16$ ).

In Fig. 3, we plot NMSE versus  $p_1 N$ , the expected number of active coefficients, for FBMP under several values of search parameter  $D$ . There we see that NMSE performance quickly degrades as  $p_1 N$  increases above  $p_1 N = 10$  (i.e., above  $\frac{p_1 N}{M} = 0.16$ ), mirroring the results in Figs. 1-2. As in Fig. 1,

when  $\frac{p_1 N}{M} > 0.16$ , increasing  $D$  from 1 to 10 can yield an NMSE improvement of 3 dB. When  $\frac{p_1 N}{M} \leq 0.16$ , however,  $D = 1$  appears to suffice.

In Fig. 4, we plot NMSE versus SNR for FBMP under several values of search parameter  $D$  (where  $\{M, p_1\}$  correspond to the aforementioned breakpoints in the NMSE-vs- $M$  and NMSE-vs- $p_1 N$  curves). Figure 4 shows a satisfying linear relationship between NMSE and SNR (in dB). As expected, the effect of increasing  $D$  from 1 to 10 is negligible because  $\frac{p_1 N}{M} = 0.16$ ; a more significant effect would be expected if  $\frac{p_1 N}{M}$  had been larger.

In Fig. 5, we plot NMSE versus SNR for two FBMP-derived estimates: the (approximate) MMSE estimate  $\hat{\mathbf{x}}_{\text{ammse}}$  from (18) and the quasi-MAP estimate  $\hat{\mathbf{x}}_{\text{amap}}$  from (36):

$$\hat{\mathbf{x}}_{\text{amap}} := \mathbb{E}\{\mathbf{x} | \mathbf{y}, \hat{\mathbf{s}}_{\text{map}}\}. \quad (36)$$

Whereas  $\hat{\mathbf{x}}_{\text{ammse}}$  is the *average* of the conditional MMSE estimates  $\mathbb{E}\{\mathbf{x} | \mathbf{y}, \mathbf{s}\}$  over  $\mathbf{s} \in \hat{\mathcal{S}}_*$ , the estimate  $\hat{\mathbf{x}}_{\text{amap}}$  is MMSE conditioned on (FBMP’s estimate of) the MAP basis-configuration  $\hat{\mathbf{s}}_{\text{map}}$ . In terms of average NMSE, Fig. 5 demonstrates that  $\hat{\mathbf{x}}_{\text{ammse}}$  are about 1 dB better than  $\hat{\mathbf{x}}_{\text{map}}$  at  $\text{SNR} \leq 10$  dB and about 0.5 dB better at  $\text{SNR} > 10$  dB. The improvement reflects the advantage of allowing for *uncertainty* in the estimated basis.

Finally, in Fig. 6, we plot average FBMP runtime versus search parameter  $D$ . As expected from the algorithmic description in Table I, the runtime scales linearly in  $D$ .

## B. Comparison To Other Algorithms

In Figs. 7-8 we compare FBMP to several other popular sparse estimation algorithms, including SparseBayes [16], OMP [4], StOMP [20], GPSR-Basic [21], and BCS [22]. Unless otherwise noted, the model parameters were set at  $N = 512$ ,  $M = 128$ ,  $p_1 = 0.04$ , and  $\sigma^2 = 0.001$  (which corresponds to  $\text{SNR} = 19$  dB at the nominal values of  $N$ ,  $M$ , and  $p_1$ ). Our plots represent an average of 100 independent model realizations.

For FBMP, we used *non-optimized* MATLAB code (which we plan to optimize in the near future), and unless otherwise noted used  $D = 5$  and the same  $P$  specified in Section IV-A. For the other algorithms, we used the publicly available implementations that were found at the web-sites listed in our bibliography. The algorithmic parameters were chosen largely in accordance with suggested values provided by the authors of the software, or in accordance with values used in examples that accompanied the algorithms. The SparseBayes algorithm was tested with the initial hyper-parameter set to  $\alpha = 1$ . StOMP was tested using the “False Alarm Control” thresholding strategy, with the thresholding parameter set to  $\frac{M}{NQ} (1 - \frac{1}{M} \|\mathbf{x}\|_0)$ , where the default number of iterations,  $Q = 10$ , was used. The  $\ell_1$ -penalty in the GPSR algorithm was chosen as  $\tau = 0.1 \|\mathbf{A}^H \mathbf{y}\|_\infty$ , and the MSE kept for comparison purposes was the smaller of the MSEs of the un-debiased and debiased estimates. The BCS algorithm was tested with the “Adaptive CS” option turned off.

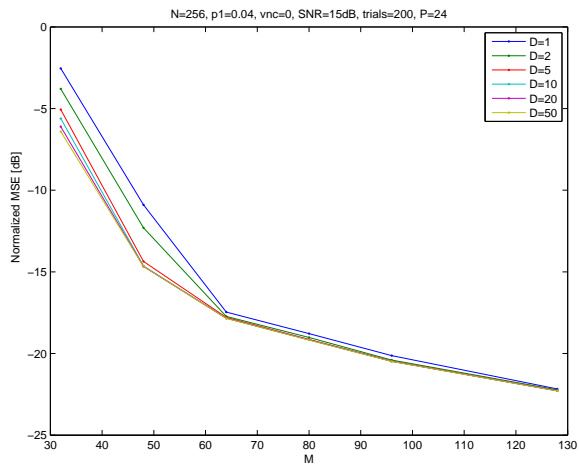


Fig. 1. Normalized MSE versus observation length  $M$  for FBMP under several values of search parameter  $D$ . From top to bottom, the traces reflect increasing  $D$ . (See the graph title for configuration.)

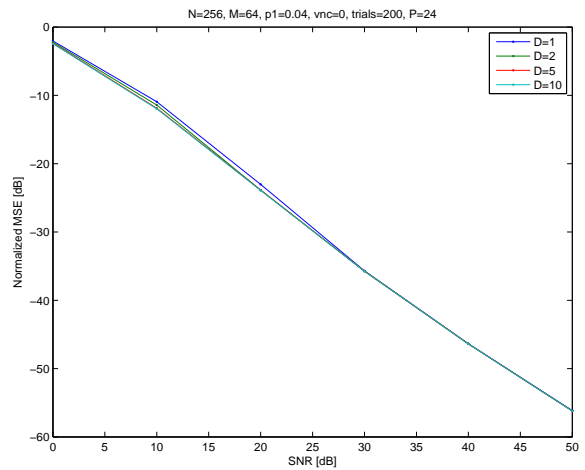


Fig. 4. Normalized MSE versus SNR for FBMP under several values of search parameter  $D$ . From top to bottom, the traces reflect increasing  $D$ . (See the graph title for configuration.)

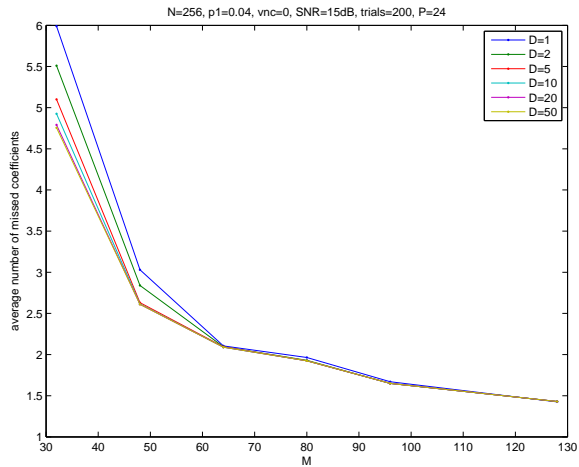


Fig. 2. Number of active coefficients in  $\mathbf{x}$  missing from FBMP's estimate of the most probably basis configuration versus observation length  $M$  and under several values of search parameter  $D$ . From top to bottom, the traces reflect increasing  $D$ . (See the graph title for configuration.)

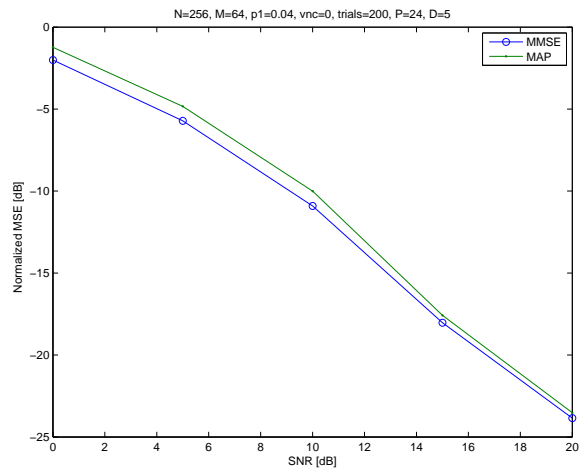


Fig. 5. Normalized MSE versus SNR of FBMP-returned MMSE estimate  $\hat{\mathbf{x}}_{\text{ammse}}$  and quasi-MAP estimate  $\hat{\mathbf{x}}_{\text{amap}}$ . (See the graph title for configuration.)

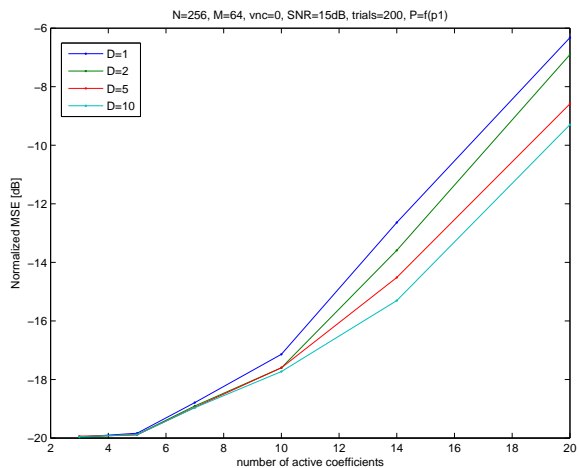


Fig. 3. Normalized MSE versus number of active coefficients (i.e.,  $\|\mathbf{x}\|_0$ ) for FBMP under several values of search parameter  $D$ . From top to bottom, the traces reflect increasing  $D$ . (See the graph title for configuration.)

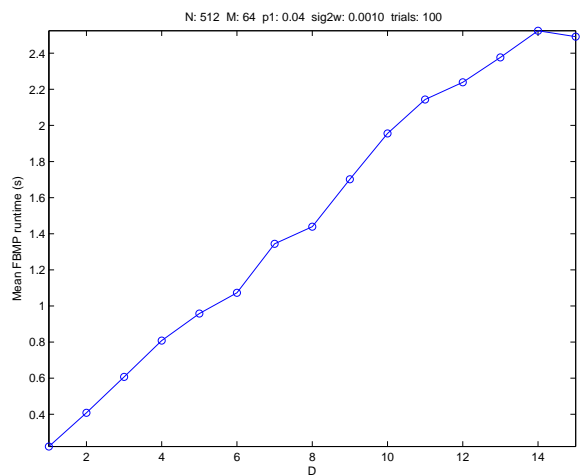


Fig. 6. Average FBMP runtime versus search parameter  $D$ . (See the graph title for configuration.)

In Fig. 7 we plot NMSE versus observation length  $M$  (at  $\sigma^2 = 0.001$ ) for the various sparse estimation algorithms. There we see that FBMP achieved significantly lower NMSE than the other algorithms over the examined range of  $M$ . In particular, it outperformed BCS by approximately 3 dB, it outperformed OMP by 3 dB at small  $M$  and 10 dB at large  $M$ , and it outperformed the other algorithms by even more. In Fig. 8 we plot NMSE versus SNR (at  $M = 128$ ) for the various sparse estimation algorithms. Again, the NMSEs achieved by FBMP were significantly lower than those achieved by the other algorithms. At SNR = 22 dB, FBMP outperformed BCS by approximately 3 dB and the other algorithms by > 9 dB; at SNR = 15 dB, FBMP outperformed all other algorithms by > 6 dB; and, at SNR = 3 dB, FBMP outperformed GPSR by approximately 1 dB and the other algorithms by > 5 dB.

Finally, in Fig. 9, we plot average runtime versus observation length  $M$  for the various sparse estimation algorithms. For FBMP, we used  $D = 1$ . Figure 9 shows that FBMP is about an order of magnitude faster than SparseBayes, on the same order of complexity as BCS, and about an order of magnitude slower than OMP, StOMP, and GPSR. We anticipate that optimized FBMP code will yield improved runtimes.

## V. DISCUSSION

The Bayesian framework provides a report on the confidence of estimates of both the coefficients  $\mathbf{x}$  and the basis configuration  $\mathbf{s}$ . In particular, the basis selection metric  $\mu(\mathbf{s})$  provides a posterior confidence label for a candidate basis configuration  $\mathbf{s}$ , in addition to providing the MMSE estimate  $\hat{\mathbf{x}}_{\text{mmse}}$  through (16). Specifically, from (10), we can write the posterior probability of basis configuration  $\mathbf{s}$  as

$$p(\mathbf{s}|\mathbf{y}) = \frac{\exp\{\mu(\mathbf{s})\}}{\sum_{\mathbf{s}' \in \mathcal{S}} \exp\{\mu(\mathbf{s}')\}} \approx \frac{\exp\{\mu(\mathbf{s})\}}{\sum_{\mathbf{s}' \in \mathcal{S}_*} \exp\{\mu(\mathbf{s}')\}}, \quad (37)$$

where the approximation in (37) includes only the basis configurations  $\mathcal{S}_* \subset \mathcal{S}$  that account for the dominant values of  $\exp\{\mu(\mathbf{s})\}$ . Likewise, (19) provides an approximate error covariance for the MMSE estimate  $\hat{\mathbf{x}}_{\text{mmse}}$ . These posterior confidence values reflect the ambiguity inherently present in the sparse inference problem—an ambiguity especially evident when the SNR is low and/or the correlation among the columns of  $\mathbf{A}$  is high.

Standard errors for estimated  $\hat{\mathbf{x}}$  are largely absent in the compressive sensing literature. Exceptions are found in [7], [22] which give the error covariance for the simple linear problem conditioned on *perfect knowledge of the active basis elements*. As noted by Tibshirani [7], such a measure of posterior uncertainty has dubious value, because “a difficulty with this formula is that it gives an estimated variance of 0 for predictors with”  $s_i = 0$ . In this light, we expect certain advantages for algorithms that consider the active basis as implicitly uncertain.

A Gaussian mixture model similar to that in Section II was likewise adopted by Larsson and Selén [18], who also constructed the MMSE estimate in the manner of (18) but with an  $\mathcal{S}_*$  that contains exactly one sequence  $\mathbf{s}$  for each Hamming

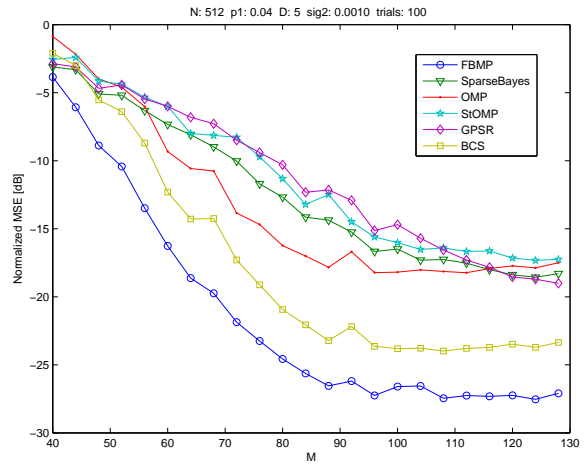


Fig. 7. Normalized MSE versus observation length  $M$  for several algorithms. (See the graph title for configuration.)

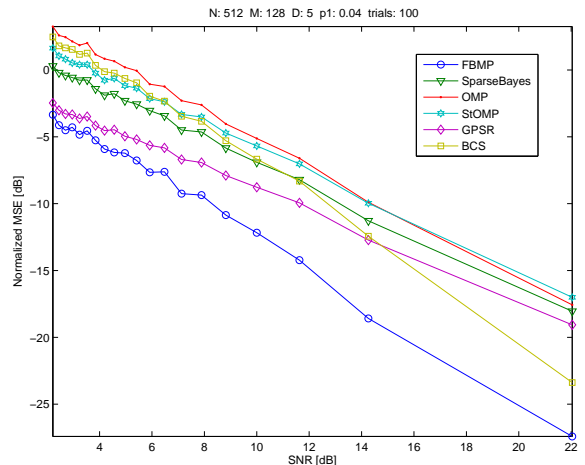


Fig. 8. Normalized MSE versus SNR for several algorithms. (See the graph title for configuration.)

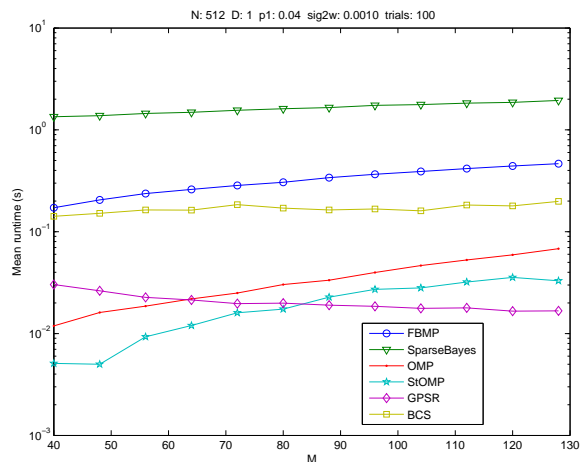


Fig. 9. Average runtime versus observation length  $M$  for several algorithms. (See the graph title for configuration.)

weight 0 to  $N$ . They proposed to find these  $s$  via greedy deflation, i.e., starting with an all-active basis configuration and recursively deactivating one element at a time. Thus, the  $D = 1$  version of the BMP heuristic from Section III-C recalls the heuristic of [18], but in reverse. Note, however, that the *fast*  $D = 1$  BMP presented in Section III-E has a complexity of only  $\mathcal{O}(NMP)$ , in comparison to  $\mathcal{O}(N^3M^2)$  for the technique in [18]. Given the typically large values of  $N$  encountered in practice, the complexity of FBMP can be several orders of magnitude lower than that of [18].

As a caveat, we should emphasize that our greedy FBMP search returns only  $\hat{S}_*$ , an *estimate* of the dominant subset  $S_*$ , along with the values of  $\mu(s)$  for  $s \in \hat{S}_*$ . Thus, while the values  $\mu(s)$  returned by FBMP can be used to compute exact *ratios* between the posterior probabilities of the configurations in  $\hat{S}_*$ , the *absolute* posteriors of these configurations (as approximated by (37) with  $\hat{S}_*$  in place of  $S_*$ ) will only be accurate when  $\hat{S}_*$  indeed contains  $S_*$ . For example, if FBMP somehow missed the MAP configuration  $\hat{s}_{\text{map}}$  (i.e.,  $\hat{s}_{\text{map}} \notin \hat{S}_*$ ), then we would expect a large discrepancy between  $\sum_{s' \in S_*} \exp\{\mu(s')\}$  and  $\sum_{s' \in \hat{S}_*} \exp\{\mu(s')\}$  which would in turn corrupt the FBMP estimates of  $p(s|\mathbf{y})$  and  $\text{Cov}\{\mathbf{x}|\mathbf{y}\}$ . Fortunately, the proposed greedy FBMP basis search seems to perform quite well, as least for  $\frac{p_1 N}{M} \leq 0.16$  (as suggested by Fig. 2).

Although the model in Section II assumed that each  $x_i$  is generated according to a binary mixture of zero-mean Gaussians, one can imagine extending the model to, e.g., a mixture of finitely many Gaussians with non-zero means. In this case, one would need to generalize the BMP search heuristic of Section III-C to handle several types of active coefficient (e.g., one for each allowed mean).

## VI. CONCLUSION

In this paper, we proposed an algorithm for joint basis selection and sparse parameter estimation which we call fast Bayesian matching pursuit (FBMP). In brief, FBMP models each unknown coefficient  $x_i$  as either inactive or active (with prior probability  $p_1$ ), where an i.i.d. Gaussian distribution (with zero mean and variance  $\sigma_1^2$ ) is assigned to the values of active coefficients. The observation  $\mathbf{y}$  is then modeled as an AWGN-corrupted version of the unknown coefficients that has been mixed by a known matrix  $\mathbf{A}$ . FBMP navigates through the tree of active/inactive configurations  $\mathcal{S}$  with the goal of finding the configurations with dominant posterior probability,  $S_*$ . The search is controlled by a parameter  $D$  which effects a tradeoff between complexity and accuracy. Numerical experiments suggest that the estimates returned by FBMP outperform (in normalized MSE) those of other popular algorithms (e.g., SparseBayes, OMP, StOMP, GPSR-Basic, BCS) by several dB in typical situations.

We plan, in the near future, to extend FBMP to the case where the active coefficients are complex Gaussian with non-zero means chosen from a finite set according to some prior probabilities. An example of the non-zero mean sparse signal model can be found in electron paramagnetic resonance

imaging [23], where micro-liter particulate probes are inserted into a tumor and fill less than 0.25% volume in the field of view. The fabrication of the paramagnetic signal probes results in variable shape, size and electron spin density, giving rise to a non-zero-mean and nearly Gaussian distribution of signal strength in a very few active voxels.

## REFERENCES

- [1] J. Hogbom, "Aperture synthesis with a non-regular distribution of interferometer baselines," *Astrophys. J. Suppl. Ser.*, vol. 15, pp. 417–426, 1974.
- [2] H. Lee, D. Sullivan, and T. Huang, "Improvement of discrete band-limited signal extrapolation by iterative subspace modification," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, pp. 1569 – 1572, 1987.
- [3] Y. C. Pati, R. Rezaifar, and P. S. Krishnaprasad., "Orthogonal matching pursuit: Recursive function approximation with applications to wavelet decomposition," in *Proc. 27th Ann. Asilomar Conf. Signals, Systems, and Computers*, 1993.
- [4] J. Tropp and A. Gilbert, "Signal recovery from random measurements via orthogonal matching pursuit," *IEEE Trans. Information Theory*, vol. 53, pp. 4655–4666, 2007. (software available at <http://sparselab.stanford.edu/>).
- [5] A. E. Hoerl and R. W. Kennard, "Ridge regression: Biased estimation for nonorthogonal problems," *Technometrics*, vol. 12, pp. 55–67, 1970.
- [6] S. S. Chen, D. L. Donoho, and M. A. Saunders, "Atomic decomposition by basis pursuit," *SIAM Journal on Scientific Computing*, vol. 20, no. 1, pp. 33–61, 1998.
- [7] R. Tibshirani, "Regression shrinkage and selection via the lasso," *J. R. Statist. Soc. B*, vol. 58, no. 1, pp. 267 – 288, 1996.
- [8] M. Çetin and W. C. Karl, "Feature-enhanced synthetic aperture radar image formation based on nonquadratic regularization," *IEEE Trans. Image Process.*, vol. 10, pp. 623–631, 2001.
- [9] C. Bouman and K. Sauer, "A generalized Gaussian image model for edge-preserving MAP estimation," *IEEE Trans. Image Process.*, vol. 2, pp. 296–310, 1993.
- [10] A. H. Delaney and Y. Bresler, "A fast and accurate Fourier algorithm for iterative parallel-beam tomography," *IEEE Trans. Image Process.*, vol. 5, pp. 740–753, May 1996.
- [11] S. Levy and P. K. Fullagar, "Reconstruction of a sparse spike train from a portion of its spectrum and application to high-resolution deconvolution," *Geophysics*, vol. 46, no. 9, pp. 1235–1243, 1981.
- [12] B. D. Rao and K. Kreutz-Delgado, "An affine scaling methodology for best basis selection," *IEEE Trans. Signal Processing*, vol. 47, pp. 187–200, 1999.
- [13] D. L. Donoho, M. Elad, and V. N. Temlyakov, "Stable recovery of sparse overcomplete representations in the presence of noise," *IEEE Trans. Information Theory*, vol. 52, no. 1, pp. 6–18, 2006.
- [14] J. A. Tropp, "Just relax: Convex programming methods for identifying sparse signal," *IEEE Trans. Info. Theory*, vol. 51, no. 3, pp. 1030–1051, 2006.
- [15] E. Candès, J. Romberg, and T. Tao, "Stable signal recovery from incomplete and inaccurate measurements," *Communications on Pure and Applied Mathematics*, vol. 59, no. 8, pp. 1207–1223, 2006.
- [16] M. E. Tipping, "Sparse Bayesian learning and the relevance vector machine," *J. Machine Learning Res.*, vol. 1, pp. 211–244, 2001. (software available at <http://www.miketipping.com/index.php?page=rvm>).
- [17] D. Wipf and B. Rao, "Sparse Bayesian learning for basis selection," *IEEE Trans. Signal Process.*, vol. 52, pp. 2153 – 2164, 2004.
- [18] E. Larsson and Y. Selén, "Linear regression with a sparse parameter vector," *IEEE Trans. Signal Process.*, vol. 55, pp. 451 – 460, 2007.
- [19] H. V. Poor, *An Introduction to Signal Detection and Estimation*. New York: Springer, 2nd ed., 1994.
- [20] D. L. Donoho, Y. Tsaig, I. Drori, and J.-C. Starck, "Sparse solution of underdetermined linear equations by stagewise orthogonal matching pursuit," Tech. Rep. 2006-02, Dept. of Statistics, Stanford University, Stanford, CA, 2006. (software available at <http://sparselab.stanford.edu/>).

- [21] M. A. T. Figueiredo, R. D. Nowak, and S. J. Wright, "Gradient projection for sparse reconstruction: Application to compressed sensing and other inverse problems," *IEEE Journal of Selected Topics in Signal Processing*, vol. 1, no. 4, pp. 586–597, 2007. (software available at <http://www.lx.it.pt/~mtf/GPSR/>).
- [22] S. Ji and L. Carin, "Bayesian compressive sensing and projection optimization," in *Proc. 24th Intern. Conf. Machine Learning (ICML)*, 2007. (software available at <http://www.ece.duke.edu/~shji/BCS.html>).
- [23] P. Kuppasamy, "EPR spectroscopy in biology and medicine," *Antioxid. Redox Signa*, vol. 6, pp. 583–585, 2004.