

Near-Oracle Performance of Greedy Block-Sparse Estimation Techniques from Noisy Measurements

Zvika Ben-Haim, *Student Member, IEEE*, and Yonina C. Eldar, *Senior Member, IEEE*

Abstract—This paper examines the ability of greedy algorithms to estimate a block sparse parameter vector from noisy measurements. In particular, block sparse versions of the orthogonal matching pursuit and thresholding algorithms are analyzed under both adversarial and Gaussian noise models. In the adversarial setting, it is shown that estimation accuracy comes within a constant factor of the noise power. Under Gaussian noise, the Cramér–Rao bound is derived, and it is shown that the greedy techniques come close to this bound at high SNR. The guarantees are numerically compared with the actual performance of block and non-block algorithms, highlighting the advantages of block sparse techniques.

I. INTRODUCTION

The success of signal processing techniques depends to a large extent on the availability of an appropriate model which captures our knowledge of the system under consideration and translates it to a productive mathematical framework. There is consequently an ongoing search for mathematical models which can accurately describe real-world signals. In recent years, much research has been devoted to the sparse representation model, which stems from the observation that many signals can be approximated using a small number of elements, or “atoms,” chosen from a large dictionary [1]–[3]. Thus, we may write $\mathbf{y} = \mathbf{D}\mathbf{x} + \mathbf{w}$, where the signal \mathbf{y} is a linear combination of a small number of columns of the dictionary matrix \mathbf{D} , corrupted by noise \mathbf{w} . Since only a small number of elements of \mathbf{D} are required for this representation, the vector \mathbf{x} is sparse, i.e., most of its entries equal 0. It turns out that the sparsity assumption can be used to accurately estimate \mathbf{x} from \mathbf{y} , even when the number of possible atoms (and thus, the length of \mathbf{x}) is greater than the number of measurements in \mathbf{y} [2], [4], [5]. This model has been used to great advantage in many fundamental fields of signal processing, including compressed sensing [1], [2], denoising [6], deblurring [7], and interpolation [8].

The assumption of sparsity is an example of a much more general class of signal models which can be described as a union of subspaces [9]–[11]. Indeed, each support pattern defines a subspace of the space of possible parameter vectors. Saying that the parameter contains no more than k nonzero entries is equivalent to stating that \mathbf{x} belongs to the union of all such subspaces. Unions of subspaces are proving to be a powerful generalization of the sparsity model. Apart from

ordinary sparsity, unions of subspaces have been applied to estimate signals as diverse as pulse streams [12], [13], multi-band communications [14]–[16], and block sparse vectors [11], [17]–[19], the latter being the focus of this paper. The common thread running through these applications is the ability to exploit the union of subspaces structure in order to achieve accurate reconstruction of signals from a very low number of measurements.

The block sparsity model is based on the realization that in many practical sparse representation settings, not all support patterns are equally likely. Specifically, if a particular element of \mathbf{x} is nonzero, then in many cases “similar” elements in \mathbf{x} are also nonzero. The precise definition of similarity is context-dependent. For example, in Fourier-based dictionaries, neighboring frequency bins are often jointly nonzero, while in wavelet-based dictionaries, nonzero entries in a certain detail level are likely to be correlated with nonzeros in higher detail levels. Consequently, the sparsity model does not incorporate all of the structure present in the signal. The block sparsity approach aims to partially overcome this drawback by partitioning the vector \mathbf{x} into blocks, each of which contains a small number of elements. The structure imposed by the block sparsity model is that no more than a small number k of blocks are nonzero. The model thus favors the use of related atoms, rather than sporadic dictionary columns. Consequently, block sparsity is well-suited for those situations described above, in which specific atoms tend to be used together.

The usefulness of a model depends on the existence of efficient and effective methods for estimating a signal \mathbf{x} from its measurements. Fortunately, estimators designed for the ordinary sparsity model can be readily adapted to the block sparse setting. Thus, previous work has described techniques such as block orthogonal matching pursuit (BOMP) [19] and the mixed ℓ_2/ℓ_1 -optimization (L-OPT) [11], [18], the latter being a block version of the Lasso. In this paper, we also describe a block-sparse version of the thresholding algorithm, which we refer to as block-thresholding (BTH). The BOMP and BTH approaches are representatives of a class of so-called greedy algorithms, which attempt to identify the support of \mathbf{x} by choosing at each step the most likely candidate. In this paper we restrict attention to these greedy techniques, which are simpler (and more naive) than convex relaxation techniques such as L-OPT, and are therefore more suitable for implementation in large-scale or computationally parsimonious settings.

Having described various estimation algorithms, it is natural to ask what can be guaranteed analytically about the performance of these methods in practice. For example, in the ordinary (non-block) sparsity setting, a rich collection of performance guarantees exists for various algorithms under

different noise models. In particular, a distinction is made between adversarial and random noise models. In the former case, nothing is known about \mathbf{w} except that it is bounded, $\|\mathbf{w}\|_2 \leq \varepsilon$; in particular, \mathbf{w} might be chosen so as to maximally harm a given estimation algorithm. Consequently, guarantees in this case are relatively weak, ensuring only that the error in \mathbf{x} is on the order of ε [2], [4], [5]. By contrast, when the noise is random, estimation performance is considerably improved for most noise realizations [4], [20], [21].

It is natural to seek an extension of these results to the block sparsity model. In the absence of noise, successful recovery of a block sparse parameter \mathbf{x} from measurements $\mathbf{y} = \mathbf{D}\mathbf{x}$ has been demonstrated in the past for both BOMP and L-OPT [11], [19]. However, to the best of our knowledge, the only result providing analytical guarantees for a block sparse estimator under noise was given in [11], where the performance of L-OPT was analyzed under adversarial noise. The goal of this paper is to analyze the performance of the greedy algorithms BOMP and BTH under both adversarial and random noise models. As we will see, despite the fact that these greedy algorithms are simpler and more efficient to implement, their performance is close to the optimal achievable results.

Specifically, we first analyze the adversarial noise model, and show that both BOMP and BTH achieve an error on the order of ε when the noise is bounded by $\|\mathbf{w}\|_2 \leq \varepsilon$. These results generalize previous guarantees in several ways: First, when each block contains one element, we recover the non-block sparsity guarantee of Donoho et al. [5]. Second, when the noise bound ε equals 0, we obtain the noise-free guarantees of Eldar et al. [19].

We next turn to the random noise model, and examine in particular the case in which \mathbf{w} is white Gaussian noise. We derive the Cramér–Rao bound (CRB) for estimating \mathbf{x} from its measurements, and show that this bound equals the error of the “oracle estimator” which knows the locations of the nonzero blocks of \mathbf{x} . However, while the oracle estimator relies on information which is unavailable in practice, the CRB is known to be achievable by the maximum likelihood (ML) technique at high SNR. Unfortunately, the ML approach is NP-complete, and thus can probably not be implemented efficiently. Nevertheless, we proceed to show that both BOMP and BTH come within a nearly constant factor of the CRB at high SNR, for dictionaries satisfying suitable requirements. Once again, when each block contains one element, we can recover previously known guarantees for non-block sparsity [21] from our results. Furthermore, we show that in typical block sparse situations, the performance guarantees of block algorithms is substantially better than that of non-block techniques.

The rest of this paper is organized as follows. The block sparse setting is defined in Section II, and the BOMP and BTH techniques are described in Section III. The adversarial noise model is then analyzed in Section IV. The treatment of random noise begins with the derivation of the CRB in Section V, while performance guarantees for this case appear in Section VI. Finally, the guarantees and the CRB are compared with the actual performance of BOMP and BTH in a numerical study in Section VII.

II. PROBLEM SETTING

A. Notation

The following notation is used throughout the paper. Matrices and vectors are denoted by boldface uppercase letters \mathbf{M} and boldface lowercase letters \mathbf{v} , respectively. The ℓ_2 norm of a vector \mathbf{v} is $\|\mathbf{v}\|_2$ and the spectral norm of a matrix \mathbf{M} is $\|\mathbf{M}\|$. The expectation of a random vector \mathbf{v} will be denoted $\mathbb{E}\{\mathbf{v}\}$ or, occasionally, $\mathbb{E}_{\mathbf{x}}\{\mathbf{v}\}$, where the subscript is intended to emphasize the fact that the expectation is a function of the deterministic quantity \mathbf{x} . The adjoint and the Moore–Penrose pseudoinverse of a matrix \mathbf{M} are denoted, respectively, by \mathbf{M}^* and \mathbf{M}^\dagger , while the column space of \mathbf{M} is $\mathcal{R}(\mathbf{M})$. We denote by $\mathbf{v}[i]$ the i th d -element block of a vector \mathbf{v} of length $N = Md$. Thus

$$\mathbf{v}[i] \triangleq [v_{(i-1)d+1}, v_{(i-1)d+2}, \dots, v_{id}]^T, \quad 1 \leq i \leq M. \quad (1)$$

Consequently, we may write

$$\mathbf{v} = [\mathbf{v}^T[1], \dots, \mathbf{v}^T[M]]^T. \quad (2)$$

Similarly, given a matrix \mathbf{M} having N columns, the submatrix $\mathbf{M}[i]$ contains the columns $(i-1)d+1, (i-1)d+2, \dots, id$ of \mathbf{M} , i.e., those columns of \mathbf{M} which correspond to the i th block. The support $\text{supp}(\mathbf{v})$ of \mathbf{v} is defined as the set of indices of nonzero blocks of \mathbf{v} ; formally

$$\text{supp}(\mathbf{v}) \triangleq \{i : \mathbf{v}[i] \neq \mathbf{0}\}. \quad (3)$$

Given an index set I , the vector \mathbf{v}_I is constructed as the subvector of \mathbf{v} containing the blocks indexed by I ; in other words, if $I = \{i_1, \dots, i_p\}$, then

$$\mathbf{v}_I = [\mathbf{v}^T[i_1], \dots, \mathbf{v}^T[i_p]]^T. \quad (4)$$

Likewise, the submatrix \mathbf{M}_I contains the column blocks indexed by I , so that

$$\mathbf{M}_I = [\mathbf{M}[i_1], \dots, \mathbf{M}[i_p]]. \quad (5)$$

To uniquely define \mathbf{v}_I and \mathbf{M}_I , we will assume as a convention that the elements of I are sorted, i.e., $i_1 < i_2 < \dots < i_p$.

B. Problem Definition

Let $\mathbf{x} \in \mathbb{C}^N$ be a deterministic block-sparse vector, i.e., \mathbf{x} consists of M blocks $\mathbf{x}[1], \dots, \mathbf{x}[M]$ of size d , of which at most k are nonzero [19]. The maximum support size k is assumed to be known. The block sparsity restriction can then be written as

$$\mathbf{x} \in \mathfrak{X} \triangleq \{\mathbf{v} \in \mathbb{R}^N : |\text{supp}(\mathbf{v})| \leq k\}. \quad (6)$$

For convenience, let $S \triangleq \text{supp}(\mathbf{x})$ be the support of the parameter \mathbf{x} , and let $s = |S|$. Note the distinction between k and s : It is known that at most k blocks are nonzero, but the actual number of nonzero blocks s is unknown and may be smaller than k . In the sequel, it will be useful to define

$$\begin{aligned} |x_{\max}| &\triangleq \max_{i \in S} \|\mathbf{x}[i]\|_2, \\ |x_{\min}| &\triangleq \min_{i \in S} \|\mathbf{x}[i]\|_2. \end{aligned} \quad (7)$$

The block sparse model differs from the more common non-block sparsity setting: in the latter, it is assumed that a small number of entries (rather than blocks) in the vector \mathbf{x} are nonzero. To emphasize this difference, we will occasionally refer to the non-block sparsity model as “ordinary” or “scalar” sparsity.

We are given noisy observations

$$\mathbf{y} = \mathbf{D}\mathbf{x} + \mathbf{w} \quad (8)$$

where $\mathbf{D} \in \mathbb{C}^{L \times N}$ is a known, deterministic dictionary, and \mathbf{w} is a noise vector. Our goal is to estimate \mathbf{x} from the measurements \mathbf{y} . It will be convenient to denote the i th column (or “atom”) of \mathbf{D} as \mathbf{d}_i . Thus we have

$$\mathbf{D} = \underbrace{[\mathbf{d}_1, \dots, \mathbf{d}_d]}_{\mathbf{D}[1]}, \underbrace{[\mathbf{d}_{d+1}, \dots, \mathbf{d}_{2d}]}_{\mathbf{D}[2]}, \dots, \underbrace{[\mathbf{d}_{N-d+1}, \dots, \mathbf{d}_N]}_{\mathbf{D}[M]}. \quad (9)$$

We assume for simplicity that the dictionary atoms are normalized, $\|\mathbf{d}_i\|_2 = 1$. We also assume that the measurement system is underdetermined, i.e., the number of measurements L is less than the number of parameters N ; thus, we must utilize the structure \mathfrak{X} , for otherwise we have no hope of recovering \mathbf{x} from its measurements. Finally, we require that for any index set I of size $|I| \leq k$, the subdictionary \mathbf{D}_I has full column rank. This latter assumption is needed to ensure that after a support set is chosen, one may estimate \mathbf{x} using standard techniques for inverting an overcomplete set of linear equations, e.g., the least-squares approach.

We will provide performance guarantees for two separate noise models. First, we consider the adversarial setting, in which the noise is unknown but bounded,

$$\|\mathbf{w}\|_2 \leq \varepsilon \quad (10)$$

for a known constant $\varepsilon > 0$. In this case the goal is to provide performance guarantees which hold for all values of \mathbf{w} satisfying (10). Second, we treat additive white Gaussian noise, in which

$$\mathbf{w} \sim N(\mathbf{0}, \sigma^2 \mathbf{I}). \quad (11)$$

In this case \mathbf{w} is unbounded, and the goal will be to provide guarantees which hold with high probability.

Following [19], we define the block coherence of \mathbf{D} as

$$\mu_B \triangleq \max_{i \neq j} \frac{1}{d} \|\mathbf{D}^*[i]\mathbf{D}[j]\|. \quad (12)$$

We also define the sub-coherence

$$\nu = \max_{1 \leq \ell \leq M} \max_{(\ell-1)d+1 \leq i \neq j \leq \ell d} |\mathbf{d}_i^* \mathbf{d}_j|. \quad (13)$$

The block coherence and sub-coherence are generalizations of the concept of the coherence, which is defined as

$$\mu = \max_{1 \leq i \neq j \leq N} |\mathbf{d}_i^* \mathbf{d}_j| \quad (14)$$

and applies to dictionaries regardless of whether they have a block structure.

III. TECHNIQUES FOR BLOCK-SPARSE ESTIMATION

For reference and in order to fix notation, we now describe the two greedy algorithms for which we provide performance guarantees.

a) *Block-Thresholding (BTH)*: We propose the following straightforward extension of the well-known thresholding algorithm. Given a measurement vector $\mathbf{y} \in \mathbb{C}^L$, perform the following steps:

- 1) Compute the correlations

$$\rho_i = \|\mathbf{D}^*[i]\mathbf{y}\|_2, \quad i = 1, \dots, M. \quad (15)$$

- 2) Find the k largest correlations and denote their indices by i_1, \dots, i_k . In other words, find a set of indices $\hat{S} = \{i_1, \dots, i_k\}$ such that $\rho_i \geq \rho_j$ for all $i \in \hat{S}$ and $j \notin \hat{S}$.
- 3) The reconstructed signal is given by

$$\hat{\mathbf{x}}_{\text{BTH}} = \arg \min_{\tilde{\mathbf{x}}: \text{supp}(\tilde{\mathbf{x}}) = \hat{S}} \|\mathbf{y} - \mathbf{D}\tilde{\mathbf{x}}\|_2. \quad (16)$$

b) *Block Orthogonal Matching Pursuit (BOMP)*: The BOMP algorithm, based on the OMP algorithm [22], was first proposed in [19].

Given a measurement vector $\mathbf{y} \in \mathbb{C}^L$, perform the following steps:

- 1) Define $\mathbf{r}^0 = \mathbf{y}$.
- 2) For each $\ell = 1, \dots, k$, do the following:

- a) Set

$$i_\ell = \arg \max_i \|\mathbf{D}^*[i]\mathbf{r}^{\ell-1}\|_2. \quad (17)$$

- b) Set

$$\mathbf{x}^\ell = \arg \min_{\tilde{\mathbf{x}}: \text{supp}(\tilde{\mathbf{x}}) \subseteq \{i_1, \dots, i_\ell\}} \|\mathbf{y} - \mathbf{D}\tilde{\mathbf{x}}\|_2. \quad (18)$$

- c) Set $\mathbf{r}^\ell = \mathbf{y} - \mathbf{D}\mathbf{x}^\ell$.

- 3) The estimate is given by $\hat{\mathbf{x}}_{\text{BOMP}} = \mathbf{x}^k$.

c) *Oracle Estimator*: We will find it useful to analyze the oracle estimator, which is defined as the least-squares solution within the true support set, i.e.,

$$\hat{\mathbf{x}}_{\text{or}} = \arg \min_{\tilde{\mathbf{x}}: \text{supp}(\tilde{\mathbf{x}}) \subseteq S} \|\mathbf{x} - \tilde{\mathbf{x}}\|_2^2. \quad (19)$$

Using the notation introduced above, we have

$$\begin{aligned} (\hat{\mathbf{x}}_{\text{or}})_S &= (\mathbf{D}_S^* \mathbf{D}_S)^{-1} \mathbf{D}_S^* \mathbf{y}, \\ (\hat{\mathbf{x}}_{\text{or}})_{S^c} &= \mathbf{0} \end{aligned} \quad (20)$$

where $S^c = \{1, \dots, M\} \setminus S$ is the complement of the support set S . Note that the term “oracle estimator” is somewhat misleading, since $\hat{\mathbf{x}}_{\text{or}}$ relies on knowledge of the true support set S , and is therefore not a true estimator.

IV. GUARANTEES FOR ADVERSARIAL NOISE

We begin by stating our performance guarantees in the case of adversarial noise. The proofs of these results are quite technical and can be found in Appendix A.

Theorem 1. *Consider the setting of Section II with adversarial noise (10). Suppose that*

$$(1 - (d-1)\nu)|x_{\min}| > 2\varepsilon \sqrt{1 + (d-1)\nu} + (2k-1)d\mu_B |x_{\max}|. \quad (21)$$

Then, the BTH algorithm correctly identifies all elements of the support of \mathbf{x} , and its error is bounded by

$$\|\hat{\mathbf{x}}_{\text{BTH}} - \mathbf{x}\|_2^2 \leq \frac{\varepsilon^2}{1 - (d-1)\nu - (k-1)d\mu_B}. \quad (22)$$

Theorem 2. Consider the setting of Section II with adversarial noise (10). Suppose that

$$(1 - (d-1)\nu)|x_{\min}| > 2\varepsilon\sqrt{1 + (d-1)\nu} + (2k-1)d\mu_B|x_{\min}|. \quad (23)$$

Then, the BOMP algorithm identifies all elements of $\text{supp}(\mathbf{x})$, and its error is bounded by

$$\|\widehat{\mathbf{x}}_{\text{BOMP}} - \mathbf{x}\|_2^2 \leq \frac{\varepsilon^2}{1 - (d-1)\nu - (k-1)d\mu_B}. \quad (24)$$

The following remarks should be made concerning Theorems 1 and 2.

- *Scalar sparsity:* The scalar sparsity setting, in which \mathbf{x} has no more than k nonzero elements, can be recovered by choosing $d = 1$. In this case, BOMP and BTH reduce to their scalar versions, which are called OMP and thresholding, respectively, and the block-coherence μ_B equals the coherence μ of (14). Theorems 1 and 2 then coincide with the well-known results of Donoho et al. [5] for performance of scalar sparse signals under adversarial noise. As an example (and for future reference), the OMP performance guarantee is given below.

Corollary 1 (Donoho et al. [5]). Let $\mathbf{y} = \mathbf{D}\mathbf{x} + \mathbf{w}$ be a measurement vector of a signal \mathbf{x} having sparsity $\|\mathbf{x}\|_0 \leq k$. Suppose that the coherence μ of the dictionary \mathbf{D} satisfies

$$|x_{\min}|(1 - (2k-1)\mu) > 2\varepsilon. \quad (25)$$

Then, OMP recovers the correct support pattern of \mathbf{x} and achieves an error bounded by

$$\|\widehat{\mathbf{x}}_{\text{OMP}} - \mathbf{x}\|_2^2 \leq \frac{\varepsilon^2}{1 - (k-1)\mu}. \quad (26)$$

Note that in the case of ordinary sparsity, $d = 1$, and therefore $|x_{\min}|$ can be defined simply as the magnitude of the smallest nonzero element in \mathbf{x} .

- *Benefits and limitations of block sparsity:* It is interesting to compare the achievable performance guarantees when one utilizes the block-sparse structure, as opposed to merely using ordinary (scalar) sparsity information. For concreteness, we focus in this discussion on a comparison between OMP and BOMP, but identical conclusions can be drawn by comparing the thresholding algorithm with its block-sparse version BTH.

Consider a block sparse signal \mathbf{x} as defined in Section II. Such a signal can also be viewed as a scalar sparse signal of length $N = Md$, having no more than sd nonzero elements. It is readily shown that the coherence μ satisfies $\nu \leq \mu$ and $\mu_B \leq \mu$ [19]. Consequently,

$$\frac{\varepsilon^2}{1 - (d-1)\nu - (k-1)d\mu_B} \leq \frac{\varepsilon^2}{1 - (sd-1)\mu} \quad (27)$$

which implies that if the conditions for the performance guarantees of both BOMP and OMP hold, then the performance guarantee (24) for BOMP will be at least as good as that of OMP (26). Moreover, in typical block-sparse settings, both ν and μ_B will be substantially smaller than μ [19], and the guarantees for BOMP will then be considerably better.

These results notwithstanding, it should be noted that BOMP should not automatically be preferred over OMP in any setting. This is because the condition (23) of Theorem 2

can sometimes be weaker than that of OMP. Specifically, the factor $2\varepsilon\sqrt{1 + (d-1)\nu}$ in (23) is larger than the analogous term 2ε in (25).¹ This implies that if the sub-coherence ν is large, block sparse algorithms will not perform as well as their scalar counterparts. Such a result is to be expected: Highly correlated dictionary blocks may cause noise amplification, and in such cases, it may be preferable to separately correlate each atom with the measurements, rather than relying on the combined correlation of the entire block. Indeed, it would be quite surprising if a partition of *any* dictionary \mathbf{D} into arbitrary blocks could be shown to perform as well as a scalar sparsity algorithm, since the former adds a restriction on the possible support patterns of the vector \mathbf{x} . The lesson to be learned from this analysis is that block sparsity techniques are effective when the dictionary can be separated into blocks whose elements are orthogonal or nearly orthogonal.

- *Noiseless case:* The situation in which $\mathbf{y} = \mathbf{D}\mathbf{x}$, i.e., no noise is present in the system, has been previously analyzed in the context of block sparsity in [19]. This setting can be recovered by choosing the noise bound $\varepsilon = 0$. In this case, the condition (24) simplifies to

$$(d-1)\nu + (2k-1)d\mu_B < 1 \quad (28)$$

and Theorem 2 then amounts to a guarantee for perfect recovery of \mathbf{x} if (28) holds. This result for the noise-free setting has been previously demonstrated in [19, Thm. 3].

Similarly, by substituting $\varepsilon = 0$ into Theorem 1, one obtains a perfect recovery condition for BTH in the noiseless setting. Specifically, if the condition

$$(d-1)\nu \frac{|x_{\max}|}{|x_{\min}|} + (2k-1)d\mu_B < 1 \quad (29)$$

is satisfied, then BTH correctly recovers \mathbf{x} from its noiseless measurements $\mathbf{y} = \mathbf{D}\mathbf{x}$.

Since BTH is a much simpler algorithm than BOMP, it is not surprising that the necessary condition (29) for BTH is somewhat stronger than the corresponding condition (28) for BOMP. This difference between the conditions is indicative of the different strategies employed by the two techniques, and will be further discussed in Section VI.

- *Severity of the error:* As in the scalar sparsity scenario, the presence of adversarial noise severely limits the ability of any algorithm to perform denoising. This is evident from Theorems 1 and 2, which guarantee only that the distance between the estimates and the true value of \mathbf{x} is on the order of the noise magnitude ε . Given our detailed knowledge of the structure of the signal \mathbf{x} , one would expect more powerful denoising capabilities for typical noise realizations. Consequently, in the remainder of this paper, we adopt the assumption of random noise, which cannot align itself so as to maximally interfere with the recovery algorithms.

V. THE CRAMÉR–RAO BOUND

A central goal in assessing the quality of an estimator is to check its proximity to the best possible performance in the

¹The remaining terms in (23) are always no worse than the corresponding terms in (25).

given setting. To this end, it is common practice to compute the CRB for unbiased estimators [23], i.e., those techniques $\hat{\mathbf{x}}$ for which the bias $\mathbf{b}(\mathbf{x}) \triangleq \mathbb{E}_{\mathbf{x}}\{\hat{\mathbf{x}}\} - \mathbf{x}$ equals zero. The CRB is a lower bound on the mean-squared error $\text{MSE}(\hat{\mathbf{x}}, \mathbf{x}) = \mathbb{E}_{\mathbf{x}}\{\|\hat{\mathbf{x}} - \mathbf{x}\|_2^2\}$ for any unbiased estimator $\hat{\mathbf{x}}$.

To utilize the information inherent in the block sparsity structure, we apply the constrained CRB [24]–[27] to the present setting. In the constrained estimation scenario, one often seeks estimators which are unbiased for all parameter values in the constraint set [24], [25]. However, as we will see below, this requirement is too strict in the block sparse setting. Indeed, in Theorem 3 we show that it is not possible to construct *any* method which is unbiased for all feasible parameter values. Consequently, a weaker, local definition of unbiasedness is called for, which we refer to as \mathfrak{X} -unbiasedness [27].

Intuitively, an estimator $\hat{\mathbf{x}}$ is said to be \mathfrak{X} -unbiased at a point $\mathbf{x} \in \mathfrak{X}$ if $\mathbb{E}_{\mathbf{x}}\{\hat{\mathbf{x}}\} = \mathbf{x}$ holds at the point \mathbf{x} and at all points $\tilde{\mathbf{x}}$ in \mathfrak{X} which are sufficiently close to \mathbf{x} . To formally define \mathfrak{X} -unbiasedness, we first recall the concept of a feasible direction. A vector $\mathbf{v} \in \mathbb{C}^N$ is said to be a feasible direction at \mathbf{x} if, for any sufficiently small α , we have $\mathbf{x} + \alpha\mathbf{v} \in \mathfrak{X}$. We then say that $\hat{\mathbf{x}}$ is \mathfrak{X} -unbiased at \mathbf{x} if $\mathbb{E}_{\mathbf{x}}\{\hat{\mathbf{x}}\} = \mathbf{x}$ and if

$$\left. \frac{\partial \mathbf{b}(\mathbf{x} + \alpha\mathbf{v})}{\partial \alpha} \right|_{\alpha=0} = 0 \quad (30)$$

for any feasible direction \mathbf{v} . In other words, the bias is zero at \mathbf{x} and remains unchanged, up to a first-order approximation, when moving away from \mathbf{x} along feasible directions. This definition yields the following result, whose proof can be found in Appendix B.

Theorem 3 (Cramér–Rao bound for block-sparse signals). *Consider the setting of Section II in which the block sparse parameter vector \mathbf{x} is to be estimated from measurements corrupted by Gaussian noise (11).*

- (a) *Suppose \mathbf{x} contains fewer than k nonzero blocks, i.e., $s < k$. Then, no finite-variance estimator is \mathfrak{X} -unbiased at \mathbf{x} .*
- (b) *Suppose \mathbf{x} contains precisely k nonzero blocks, i.e., $s = k$. Then, any estimator which is \mathfrak{X} -unbiased at \mathbf{x} satisfies*

$$\text{MSE}(\hat{\mathbf{x}}, \mathbf{x}) \geq \sigma^2 \text{Tr}((\mathbf{D}_S^* \mathbf{D}_S)^{-1}). \quad (31)$$

We recall that both the MSE and the CRB are functions of the unknown vector \mathbf{x} , as is generally the case when estimating a deterministic parameter. It follows immediately from Theorem 3 that no finite-variance estimator can satisfy $\mathbb{E}_{\mathbf{x}}\{\hat{\mathbf{x}}\} = \mathbf{x}$ for all $\mathbf{x} \in \mathfrak{X}$, which explains why we previously avoided this simpler definition of unbiasedness in the constrained setting. Instead, restricting attention to a local unbiasedness requirement led to a finite CRB for almost all parameter values in \mathbf{x} : specifically, those parameters whose support is maximal, $|\text{supp}(\mathbf{x})| \triangleq s = k$.

For maximal-support values of \mathbf{x} , it is not difficult to show that the CRB (31) coincides with the MSE of the oracle estimator (20). In this case it is possible to get a sense for the value of the bound, as follows. From (44) of Lemma 1 (see Appendix A), we have that none of the eigenvalues of

$(\mathbf{D}_S^* \mathbf{D}_S)^{-1}$ are larger than $1/(1 - (d-1)\nu - (k-1)d\mu_B)$. Thus

$$\sigma^2 \text{Tr}((\mathbf{D}_S^* \mathbf{D}_S)^{-1}) \leq \frac{1}{1 - (d-1)\nu - (k-1)d\mu_B} kd\sigma^2. \quad (32)$$

In other words, when the block coherence and sub-coherence of \mathbf{D} are low, the bound of Theorem 3 will be close to $kd\sigma^2$. This value is typically much lower than the total noise variance $\mathbb{E}\{\|\mathbf{w}\|_2^2\} = L\sigma^2$. Thus, at least according to the CRB, it is possible to achieve substantial denoising in the presence of random noise. This stands in contrast to the rather disappointing guarantees presented for adversarial noise in the previous section. We may thus hope that the performance will be improved when considering random noise.

As opposed to the oracle estimator, which cannot be implemented in practice, it is well-known that the CRB can be asymptotically achieved at high SNR by the maximum likelihood (ML) estimator [23]. However, in the present setting, computing the ML estimator is NP-hard, and thus impractical. Consequently, it is of interest to determine whether there exist *efficient* techniques which come close to the performance bound (31), at least for high SNR values. As we will show in the next section, this question is answered in the affirmative: greedy block sparsity techniques do indeed approach the CRB for sufficiently high SNR.

VI. GUARANTEES FOR GAUSSIAN NOISE

In this section, we analyze the performance of block sparse algorithms when the noise \mathbf{w} is a Gaussian random variable having mean zero and covariance $\sigma^2 \mathbf{I}$. Our main performance guarantees are summarized in Theorems 4 and 5. The proofs of these theorems are found in Appendix C.

Theorem 4. *Consider the setting of Section II with additive white Gaussian noise $\mathbf{w} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$. Suppose it is known that*

$$(1 - (d-1)\nu)|x_{\min}| - (2k-1)d\mu_B|x_{\max}| \geq 2\sigma\sqrt{2\alpha d(1 + (d-1)\nu)\log N} \quad (33)$$

for some constant $\alpha \geq 1/(2d\log N)$. Then, with probability exceeding

$$1 - \frac{0.8d(2\alpha d\log N)^{d/2-1}}{N^{\alpha d-1}} \quad (34)$$

the BTH algorithm identifies the correct support of \mathbf{x} and achieves an error bounded by

$$\|\hat{\mathbf{x}}_{\text{BTH}} - \mathbf{x}\|_2^2 \leq \frac{2\alpha(1 + (d-1)\nu)}{(1 - (d-1)\nu - (k-1)d\mu_B)^2} dk\sigma^2 \log N. \quad (35)$$

Theorem 5. *Consider the setting of Section II with additive white Gaussian noise $\mathbf{w} \sim N(\mathbf{0}, \sigma^2 \mathbf{I})$. Suppose it is known that*

$$(1 - (d-1)\nu)|x_{\min}| - (2k-1)d\mu_B|x_{\min}| \geq 2\sigma\sqrt{2\alpha d(1 + (d-1)\nu)\log N} \quad (36)$$

for some constant $\alpha \geq 1/(2d \log N)$. Then, with probability exceeding (34), the BOMP algorithm identifies the correct support of \mathbf{x} and achieves an error bounded by

$$\|\hat{\mathbf{x}}_{\text{BOMP}} - \mathbf{x}\|_2^2 \leq \frac{2\alpha(1 + (d-1)\nu)}{(1 - (d-1)\nu - (k-1)d\mu_B)^2} dk\sigma^2 \log N. \quad (37)$$

We now provide some insights into the performance of block-sparse algorithms under random noise.

- *Random noise vs. adversarial noise:* As noted in Section IV, performance guarantees in the case of adversarial noise can ensure a recovery error on the order of the total noise magnitude. This is a result of the fact that the noise could, in principle, be concentrated in a single nonzero component of \mathbf{x} , whereupon it would be indistinguishable from the signal. However, for random noise, such an event is highly unlikely. Consequently, Theorems 4 and 5 provide much tighter performance guarantees: both theorems demonstrate that, with high probability, the estimation error is on the order of $dk\sigma^2 \log N$, i.e., within a constant times $\log N$ of the CRB presented in Section V. Since the noise variance $\mathbb{E}\{\|\mathbf{w}\|^2\}$ is given by $N\sigma^2$, and since typically $dk \log N \ll N$, we conclude that the block sparse algorithms have successfully removed a large portion of the noise, owing to the utilization of the union-of-subspaces structure.

- *BOMP vs. BTH:* Comparing Theorems 4 and 5 leads to an important insight concerning the advantage of the more sophisticated BOMP algorithm over its simpler counterpart. Indeed, the guarantee for BOMP requires condition (36), which basically states that $|x_{\min}|$ must be larger than a constant multiplied by the standard deviation of the noise. By contrast, for the BTH guarantee one requires the stronger condition (33), which can be interpreted as requiring $|x_{\min}|$ to be larger than a small constant times $|x_{\max}|$, plus another constant times the noise standard deviation.

To explain this difference, recall from Section III that the BTH approach relies on a single support-identification stage in which the blocks most highly correlated with the measurements are chosen as the estimated support set \hat{S} . Thus, for BTH to correctly identify the support, each block in S must be sufficiently large in magnitude to overcome interference from the noise and from the remaining blocks. Condition (33) can therefore be interpreted as a requirement that the magnitude $|x_{\min}|$ of the smallest nonzero block must be larger than the sum of the interference from the large nonzero blocks (the $|x_{\max}|$ term) and the noise. By contrast, the BOMP algorithm iteratively identifies support elements, maintaining a residual vector \mathbf{r}^ℓ containing the components of the measurement vector which have yet to be identified. Thus, BOMP requires only the ability to separately isolate each nonzero block, and hence its weaker condition (36), which necessitates only that $|x_{\min}|$ be larger than the noise.

Finally, it should be noted that when BTH and BOMP both identify the correct support set, the estimates of the two algorithms coincide, explaining the identical bounds on their performance. The conclusion from this analysis is that BOMP should be preferred if a wide dynamic range of block magnitudes is possible, but that when all blocks have roughly

the same size, the simpler and more efficient BTH technique can be used.

- *Scalar sparsity:* It is interesting to note that known results for scalar sparsity algorithms can be recovered from our block sparsity guarantees, by substituting $d = 1$ into Theorems 4 and 5. For example, consider the BOMP guarantee (Theorem 5). In the scalar case, this algorithm is known as OMP, and its performance guarantee can be written as follows.

Corollary 2. *Let $\mathbf{y} = \mathbf{D}\mathbf{x} + \mathbf{w}$ be a measurement vector of a signal \mathbf{x} having sparsity $\|\mathbf{x}\|_0 \leq k$. Suppose the coherence μ of \mathbf{D} satisfies*

$$|x_{\min}|(1 - (2k-1)\mu) \geq 2\sigma\sqrt{2\alpha \log N} \quad (38)$$

for some $\alpha > 1$. Then, with probability exceeding

$$1 - \frac{0.8/\sqrt{2}}{N^{\alpha-1}\sqrt{\alpha \log N}} \quad (39)$$

the OMP algorithm recovers the correct support of \mathbf{x} , and achieves an error bounded by

$$\|\hat{\mathbf{x}}_{\text{OMP}} - \mathbf{x}\|_2^2 \leq \frac{2\alpha}{(1 - (k-1)\mu)^2} k\sigma^2 \log N. \quad (40)$$

Corollary 2 is nearly identical to [21, Thm. 4], with the only difference being that the constant $0.8/\sqrt{2} \approx 0.566$ in (39) is replaced in [21] with the slightly better constant $1/\sqrt{\pi} \approx 0.564$. This slight discrepancy can be resolved if the more accurate version (88a) of Lemma 4 is used in the proof of Theorem 5, but the resulting expression becomes much more cumbersome in the block sparse case.

- *Block sparsity vs. scalar sparsity:* A legitimate question is whether the incorporation of the block sparsity structure substantially assists estimation algorithms. In other words, do the performance guarantees of the block algorithms BOMP and BTH compare favorably with the results achievable on identical signals using scalar sparsity algorithms, such as OMP and thresholding? This question is examined numerically in the next section.

VII. NUMERICAL EXPERIMENTS

From a practical point of view, it is important to determine whether the use of block sparse algorithms contributes significantly to the performance of estimation algorithms. After all, any block sparse signal containing k nonzero blocks of size d can also be viewed as a sparse signal containing kd nonzero elements. Is there a significant benefit in using the block algorithms rather than the ordinary scalar versions?

There are two possible approaches to answering this question. First, one may compare the performance achieved in practice by block sparse and scalar sparse algorithms. This requires a complete specification of the problem setting, including a choice of the parameter value \mathbf{x} , which is unknown in practice. Alternatively, one can compare the performance guarantees for block sparse techniques, which were derived in Section VI, to the previously known guarantees for scalar approaches [28]. The performance guarantees apply to all parameter values having a specified sparsity level, and are therefore more general. However, there may be a gap between

Problem Dimensions				Coherence		OMP		Block-OMP		Cramér-Rao
Blocks M	Block size d	Measurements L	Sparsity k	μ	μ_B	Guarantee/ σ^2	σ_{\max}	Guarantee/ σ^2	σ_{\max}	CRB/ σ^2
1200	5	3000	1	0.10	0.026	301.0	0.033	37.0	0.160	5.0
1200	5	3000	2	0.10	0.026	—	—	98.8	0.110	10.0
1200	5	3000	3	0.10	0.026	—	—	204.4	0.063	15.1
1200	5	3000	4	0.10	0.026	—	—	417.0	0.010	20.1
1200	5	3000	5	0.10	0.026	—	—	—	—	25.2
1200	5	3000	3	0.10	0.026	—	—	204.4	0.063	15.1
600	10	3000	3	0.10	0.015	—	—	364.3	0.049	30.2
300	20	3000	3	0.10	0.010	—	—	879.1	0.008	60.8
200	30	3000	3	0.10	0.007	—	—	—	—	91.8
1200	5	3000	1	0.10	0.026	301.0	0.033	37.0	0.160	5.0
1200	5	1000	1	0.17	0.043	—	—	37.0	0.144	5.0
1200	5	500	1	0.25	0.060	—	—	37.0	0.128	5.0
1200	5	100	1	0.51	0.133	—	—	37.0	0.062	5.0
1200	5	50	1	0.71	0.165	—	—	37.0	0.032	5.0
1200	5	20	1	0.90	0.197	—	—	37.0	0.003	5.0
1200	5	10	1	0.98	0.200	—	—	—	—	5.0

TABLE I
PERFORMANCE GUARANTEES FOR OMP AND BLOCK-OMP

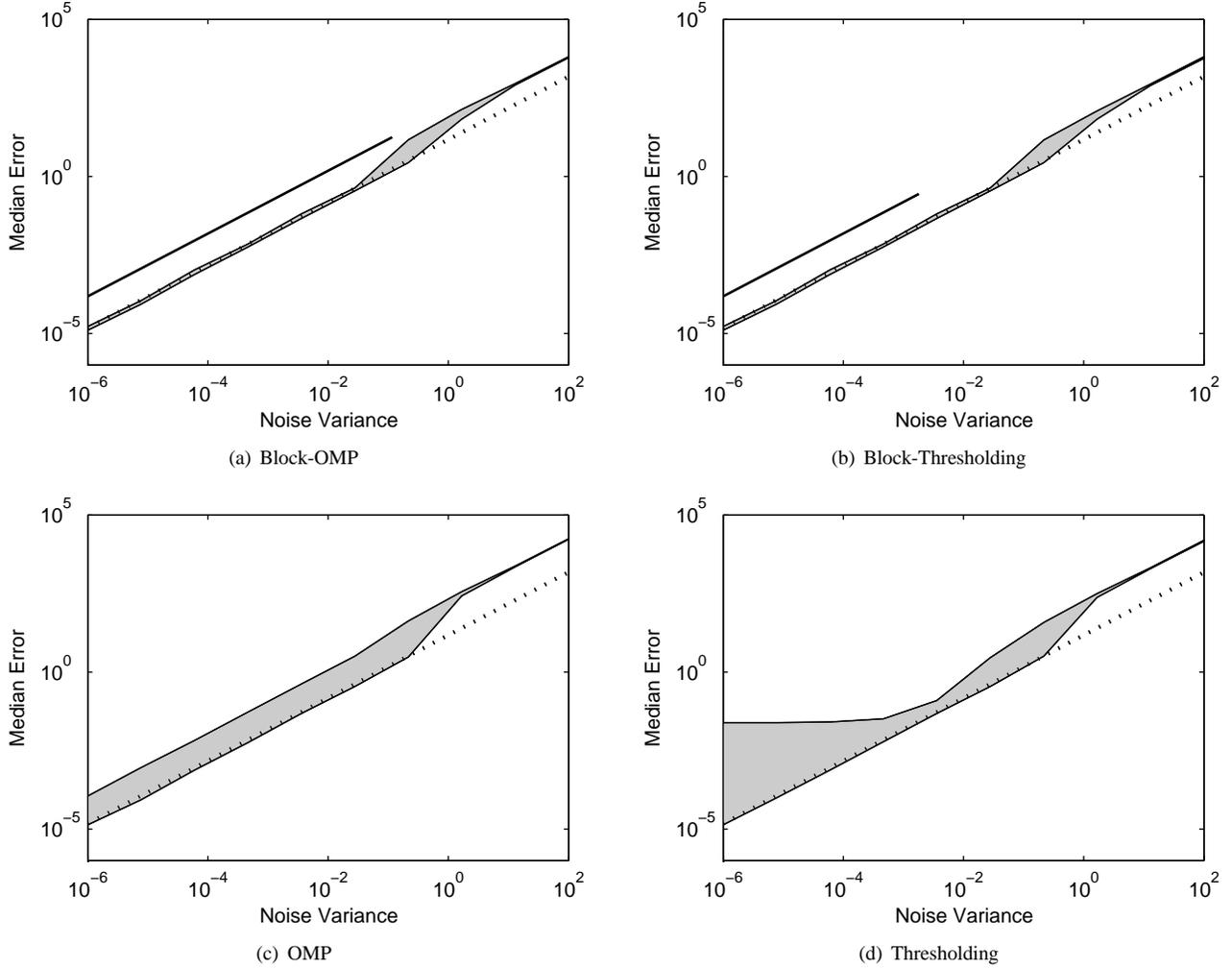


Fig. 1. Median squared error as a function of the noise variance for block and scalar sparse estimation algorithms. The shaded region indicates the range of errors encountered for different parameter values. The dotted line plots the CRB. The thick solid line in Figs. 1(a) and 1(b) indicates the performance guarantees for the block sparse algorithms; no guarantee can be made for the scalar sparsity techniques in Figs. 1(c) and 1(d).

the guarantee and the performance observed in practice. In order to take advantage of both approaches, in the following we compare both the actual performance and the guarantees of the various algorithms discussed in this paper.

In our experiments, we used dictionaries containing orthonormal blocks. Such dictionaries were constructed by first generating a random $L \times N$ matrix containing IID, zero-mean Gaussian random variables, and then performing a Gram–Schmidt procedure separately on the columns of each block. As a first experiment, we generated a variety of such dictionaries, and computed their coherence μ and block coherence μ_B . (The sub-coherence of dictionaries generated in this manner is necessarily $\nu = 0$.) These values were used to compute performance guarantees for BOMP (using Theorem 5) and for OMP (using Corollary 2). We assumed throughout that the minimum norm $|x_{\min}|$ among nonzero blocks equals 1 and that the minimum nonzero element equals $1/\sqrt{d}$. Some typical results are listed in Table I. To compute the guarantees in this table, the smallest value of α yielding a 99% probability of success was chosen. The resulting guarantee is listed in multiples of σ^2 . For example, a value of $\text{Guarantee}/\sigma^2 = 100$ means that $\|\hat{\mathbf{x}} - \mathbf{x}\|_2^2 \leq 100\sigma^2$ for 99% of the noise realizations. Also listed in Table I are the maximum noise standard deviations σ_{\max} for which the performance guarantees still hold. A dash (—) indicates that no guarantee can be made for the given setting even in the noise-free case.

It is evident from Table I that the block sparse algorithm BOMP is guaranteed to perform over a much wider range of problem settings than the scalar OMP approach. Furthermore, even when performance guarantees are provided for both techniques, those for BOMP are substantially stronger. To provide merely one striking example from Table I, note that 50 measurements suffice for BOMP to identify a signal composed of a single 5-element block among a set of 1200 possible blocks, whereas for OMP to identify such a signal at the same noise level, as many as 3000 measurements are required. The reason for this advantage is clear: the OMP algorithm must separately identify each nonzero component of the signal, and must therefore choose among a total of $\binom{1200}{5} \approx 2.1 \cdot 10^{13}$ possible support sets. This is obviously more challenging than identifying one nonzero block among a set of 1200 possibilities. Clearly, then, knowledge of a block-sparse structure can substantially improve performance if it is correctly utilized.

Table I also compares the performance guarantees with the CRB of Theorem 3. The CRB is listed for a random choice of support set S containing precisely k nonzero blocks; however, choosing different sets S only has a small effect on the value of the bound. The gap between these lower and upper bounds is not inconsiderable, and is typically on the order of a factor of 10. There are several reasons for this gap. First, the performance guarantees plotted above indicate an error which is obtained with 99% confidence, whereas the CRB is a bound on the MSE. By its very nature, the MSE averages out unusually disruptive noise realizations, and thus tends to be more optimistic. Second, different values of \mathbf{x} may yield significantly different performance; the performance guarantees apply to *all* values of \mathbf{x} , whereas the CRB is

plotted for a single, typical parameter value. Third, some loss of tightness undoubtedly results from the derivations of the theorems, i.e., there may still be room for improved bounds.

To measure the relative influence of these factors, we performed another experiment, in which the guarantees were compared with the actual performance of the various algorithms. To overcome the aforementioned pessimistic effect of a guarantee which holds with overwhelming probability, in this second experiment we computed guarantees with a 50% confidence level. In other words, these are assurances on the median of the distance between \mathbf{x} and its estimate, which captures the typical estimation error. We also computed the actual median error of the various algorithms for a variety of parameter values.

The details of this experiment are as follows. We constructed a 3000×6000 dictionary \mathbf{D} containing $M = 1200$ blocks of $d = 5$ atoms each, using the orthogonalization algorithm described above. The resulting coherence of \mathbf{D} was $\mu = 0.094$, the block coherence was $\mu_B = 0.026$, and since each block was orthonormal, the sub-coherence was $\nu = 0$. We then constructed a variety of block sparse vectors \mathbf{x} , each having $s = 3$ nonzero blocks, with $|x_{\min}| = 2\sqrt{d}$ and $|x_{\max}| = 3\sqrt{d}$. We chose the parameter vectors so as to cover as wide a range of scenarios as possible, within the aforementioned requirements. For example, some parameter vectors contained a block with a single nonzero component whose value was $|x_{\max}|$, while other vectors contained a block with each of the d elements receiving a value of $|x_{\max}|/\sqrt{d}$. Although it is clearly not feasible to cover the full range of possible parameter vectors, it is hoped that in this way some sense is given of the variability in performance for different parameter values. Indeed, as shown below, different parameters often yield widely differing estimation errors.

For each choice of a parameter vector, 20 noise realizations were generated and the resulting measurement vector \mathbf{y} was computed using (8). The BOMP, BTH, OMP, and thresholding algorithms were then applied to each of the measurement vectors. For every technique and each parameter vector, the median estimation error (among the noise realizations) was computed. The range of median estimation errors obtained for different choices of \mathbf{x} is plotted as a shaded area in Fig. 1.

In the present setting, neither of the scalar sparsity algorithms was capable of providing a performance guarantee. For BOMP and BTH, performance guarantees were available, and these are plotted as a solid line in Fig. 1. These guarantees are valid only up to a certain maximal noise variance, at which point the solid line in Fig. 1 stops. The results are also compared with the CRB of Theorem 3. It should be emphasized that the CRB is a bound on the MSE, rather than the median error, although in practice the differences between these two quantities appear to be quite small. It is also worth recalling that the CRB is a bound on unbiased estimators, while all of the techniques discussed herein are biased; nevertheless, it is evident that the CRB still provides a rough measure of the optimal performance of the proposed algorithms.

Several comments are in order concerning Fig. 1. First, the performance of both block sparse algorithms exhibits a

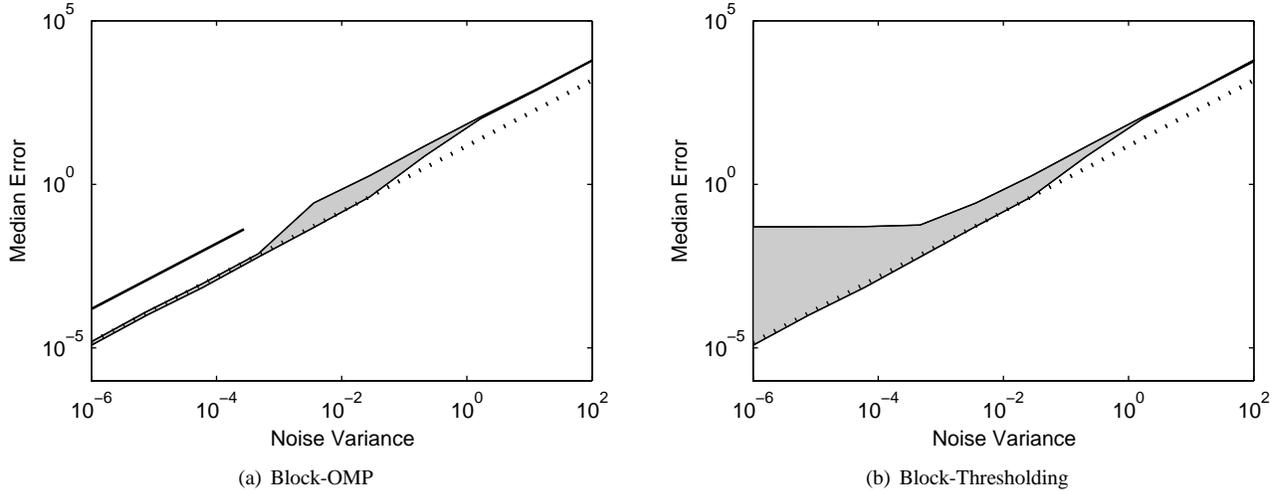


Fig. 2. Median squared error as a function of the noise variance for block sparse estimation algorithms. The shaded region indicates the range of errors encountered for different parameter values. The dotted line plots the CRB. The thick solid line in Fig. 2(a) indicates the performance guarantee for BOMP; no guarantee can be made for BTH. The deteriorated performance of BTH is a result of the existence of low-magnitude blocks.

transition: near-CRB performance for low noise levels deteriorates substantially when the noise level crosses a certain threshold. This behavior qualitatively matches the predictions of the performance guarantees, which ensure support recovery and near-CRB performance for sufficiently low noise levels. The threshold at which this transition occurs is identified fairly accurately for BOMP, and less so for BTH, although it is possible that there exist some (untested) parameter values for which the BTH transition occurs at lower noise levels. However, the numeric value of the performance guarantee is somewhat pessimistic: while the observed performance is close to the CRB for all parameter values, analytically one can guarantee only that the median error will not be larger than approximately 10 times the CRB. This result is most likely due to the various inequalities employed in the proofs of Theorems 4 and 5. Indeed, since the correct support is identified with high probability for most noise realizations, the BTH and BOMP algorithms will likely tend to coincide with the oracle estimator, whose error equals that of the CRB. The question of formally proving such a claim remains a topic for further research.

The advantages of the block sparse approach become evident when compared with scalar sparsity algorithms (Figs. 1(c) and 1(d)). For the scalar techniques, no performance guarantees can be made in the present setting. Unlike the block sparsity algorithms, the scalar approaches fail to recover the correct parameter vector even when the noise is negligible, and for some parameter values, their error does not converge to the CRB. The thresholding algorithm, in particular, ceases to improve (for some parameter values) as the noise is reduced, while the OMP approach, although significantly better than thresholding, does not converge to the CRB as do the block sparse techniques. This demonstrates the advantages of utilizing the fact that the signal is known to have a block-sparse structure.

The performance of BOMP (Fig. 1(a)) is quite similar to that of BTH (Fig. 1(b)) in the experiment above. This is not

surprising when one compares our problem setting with the guarantees of Section VI. Indeed, as we have seen, the primary difference between the BOMP and BTH algorithms is that the one-shot support estimation employed by BTH causes large-magnitude blocks to overshadow small-magnitude nonzero blocks. In the setting of Fig. 1, the range of magnitudes between $|x_{\max}| = 3\sqrt{d}$ and $|x_{\min}| = 2\sqrt{d}$ is not very large, and therefore BTH performs nearly as well as BOMP. The advantages of BOMP become readily apparent if one considers a wider dynamic range. This is illustrated in Fig. 2, in which the setup is identical to that of the previous experiment, except that parameter vectors having $|x_{\min}| = 0.1\sqrt{d}$ and $|x_{\max}| = \sqrt{d}$ were chosen, yielding a 10-fold dynamic range in the block magnitudes. In this case, while the guarantee for BOMP is hardly changed, the conditions for Theorem 4 no longer hold, so that nothing can be ensured concerning the BTH technique. Indeed, in Fig. 2 we see that BTH performs poorly for some parameter values even when the noise level is low, and its performance is no longer proportional to the CRB.

VIII. CONCLUSION

In this paper, we analyzed the performance of the greedy block algorithms BOMP and BTH under the adversarial and Gaussian noise models. In the adversarial setting $\|\mathbf{w}\|_2 \leq \varepsilon$, we showed that the estimation error equals a constant times the noise bound ε , which shows that performance in this case will not necessarily reduce the noise power. The situation is much better in the presence of random noise, where we saw that, under suitable conditions, greedy techniques obtain an error on the order of $dk\sigma^2 \log N$ with high probability; this is substantially lower than the input noise power $N\sigma^2$. Indeed, the BTH and BOMP algorithms come close to the CRB and the error of the oracle estimator.

There remain many open questions concerning the performance of block sparse techniques under random noise. For example, for scalar sparsity, performance guarantees for convex

relaxation techniques do not require assumptions on the SNR. An important challenge is to determine whether similar SNR-independent results can be demonstrated for block convex relaxation techniques such as L-OPT. Furthermore, it is well-known that scalar sparsity guarantees can be strengthened if the restricted isometry constants of the dictionary \mathbf{D} are known, as is the case, for example, when \mathbf{D} is chosen from an appropriate random ensemble. Thus, it is also of interest to provide guarantees for block techniques under random noise based on an extension of the RIP to the block sparse setting. One such extension has already been proposed in [11], and its application to the Gaussian noise model may provide tighter bounds for some performance algorithms.

APPENDIX A PROOFS FOR ADVERSARIAL NOISE

We begin by providing several lemmas which will prove useful for the analysis under both the adversarial and the Gaussian noise models.

Lemma 1. *Given a dictionary \mathbf{D} having block coherence μ_B and sub-coherence ν , we have*

$$\|\mathbf{D}^*[i]\mathbf{D}[j]\| \leq d\mu_B \quad \text{for all } i \neq j \quad (41)$$

and

$$\|\mathbf{D}[i]\|^2 = \|\mathbf{D}^*[i]\mathbf{D}[i]\| \leq 1 + (d-1)\nu. \quad (42)$$

If $1 - (d-1)\nu > 0$, then

$$\|(\mathbf{D}^*[i]\mathbf{D}[i])^{-1}\| \leq \frac{1}{1 - (d-1)\nu}. \quad (43)$$

Suppose $1 - (d-1)\nu - (k-1)d\mu_B > 0$ and let I be an index set with $|I| \leq k$. Then

$$\|(\mathbf{D}_I^*\mathbf{D}_I)^{-1}\| \leq \frac{1}{1 - (d-1)\nu - (k-1)d\mu_B}. \quad (44)$$

Proof: The bound (41) follows directly from the definition (12) of block coherence. To prove (42)–(43), observe that the diagonal elements of the matrix $\mathbf{D}^*[i]\mathbf{D}[i]$ equal 1, while the off-diagonal elements are bounded in magnitude by ν . Therefore, by the Gershgorin circle theorem [29], all eigenvalues of $\mathbf{D}^*[i]\mathbf{D}[i]$ are in the range $[1 - (d-1)\nu, 1 + (d-1)\nu]$, demonstrating (42). Furthermore, it follows that the eigenvalues of $(\mathbf{D}^*[i]\mathbf{D}[i])^{-1}$ are in the range $[(1 + (d-1)\nu)^{-1}, (1 - (d-1)\nu)^{-1}]$, leading to (43).

It remains to prove (44). To this end, let $|I| = \ell \leq k$ and write $\mathbf{D}_I^*\mathbf{D}_I$ as

$$\mathbf{D}_I^*\mathbf{D}_I = \begin{pmatrix} \mathbf{M}[1,1] & \mathbf{M}[1,2] & \cdots & \mathbf{M}[1,\ell] \\ \mathbf{M}[2,1] & \mathbf{M}[2,2] & \cdots & \mathbf{M}[2,\ell] \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{M}[\ell,1] & \mathbf{M}[\ell,2] & \cdots & \mathbf{M}[\ell,\ell] \end{pmatrix} \quad (45)$$

where each $\mathbf{M}[i,j]$ is a $d \times d$ matrix containing the correlations between two blocks of dictionary atoms. From the definition of block coherence, we have

$$\|\mathbf{M}[i,j]\| \leq d\mu_B, \quad \text{for all } i \neq j. \quad (46)$$

By a generalization of the Gershgorin circle theorem [30, Thm. 2], it follows that all eigenvalues λ of $\mathbf{D}_I^*\mathbf{D}_I$ satisfy

$$\begin{aligned} \|\mathbf{M}[i,i] - \lambda\mathbf{I}\| &\leq \sum_{j \neq i} \|\mathbf{M}[i,j]\| \leq (\ell-1)d\mu_B \\ &\leq (k-1)d\mu_B. \end{aligned} \quad (47)$$

Now, from the definition of sub-coherence, the off-diagonal elements of $\mathbf{M}[i,i]$ are no larger in magnitude than ν , while the diagonal elements of $\mathbf{M}[i,i]$ all equal 1. Therefore, by the Gershgorin circle theorem, given an arbitrary constant λ , all eigenvalues of the $d \times d$ matrix $\mathbf{M}[i,i] - \lambda\mathbf{I}$ are in the range $[1 - \lambda - (d-1)\nu, 1 - \lambda + (d-1)\nu]$. Consequently

$$\|\mathbf{M}[i,i] - \lambda\mathbf{I}\| \geq 1 - \lambda - (d-1)\nu. \quad (48)$$

Combining with (47) and rearranging, we conclude that all eigenvalues of $\mathbf{D}_I^*\mathbf{D}_I$ satisfy

$$\lambda \geq 1 - (d-1)\nu - (k-1)d\mu_B. \quad (49)$$

Consequently, the eigenvalues of $(\mathbf{D}_I^*\mathbf{D}_I)^{-1}$ are no larger than $(1 - (d-1)\nu - (k-1)d\mu_B)^{-1}$, establishing (44). ■

Lemma 2. *Consider the setting of Section II, and suppose it is known that*

$$\max_{1 \leq j \leq M} \|\mathbf{D}^*[j]\mathbf{w}\|_2 < \tau \quad (50)$$

for a given value $\tau > 0$. If the dictionary \mathbf{D} satisfies

$$(1 - (d-1)\nu) |x_{\max}| > 2\tau + (2s-1)d\mu_B |x_{\max}| \quad (51)$$

then

$$\max_{j \in S} \|\mathbf{D}^*[j]\mathbf{y}\|_2 > \max_{j \notin S} \|\mathbf{D}^*[j]\mathbf{y}\|_2 \quad (52)$$

where $S = \text{supp}(\mathbf{x})$.

If (51) is replaced by the stronger condition

$$(1 - (d-1)\nu) |x_{\min}| > 2\tau + (2s-1)d\mu_B |x_{\max}| \quad (53)$$

then

$$\min_{j \in S} \|\mathbf{D}^*[j]\mathbf{y}\|_2 > \max_{j \notin S} \|\mathbf{D}^*[j]\mathbf{y}\|_2. \quad (54)$$

Proof: The proof is an extension of [21, Lemma 3] to the block-sparse case, and is ultimately inspired by [5]. We first note that

$$\begin{aligned} \max_{j \notin S} \|\mathbf{D}^*[j]\mathbf{y}\|_2 &= \max_{j \notin S} \left\| \mathbf{D}^*[j]\mathbf{w} + \sum_{i \in S} \mathbf{D}^*[j]\mathbf{D}[i]\mathbf{x}[i] \right\|_2 \\ &\leq \max_{j \notin S} \|\mathbf{D}^*[j]\mathbf{w}\|_2 + \max_{j \notin S} \sum_{i \in S} \|\mathbf{D}^*[j]\mathbf{D}[i]\| |x_{\max}|. \end{aligned} \quad (55)$$

By (50), the first term in (55) is smaller than τ . Together with (41), we obtain

$$\max_{j \notin S} \|\mathbf{D}^*[j]\mathbf{y}\|_2 < \tau + sd\mu_B |x_{\max}| \leq \tau + kd\mu_B |x_{\max}|. \quad (56)$$

On the other hand,

$$\begin{aligned} \max_{j \in S} \|\mathbf{D}^*[j]\mathbf{y}\|_2 &= \max_{j \in S} \left\| \mathbf{D}^*[j]\mathbf{w} + \sum_{i \in S} \mathbf{D}^*[j]\mathbf{D}[i]\mathbf{x}[i] \right\|_2 \\ &\geq \max_{j \in S} \|\mathbf{D}^*[j]\mathbf{D}[j]\mathbf{x}[j]\|_2 \\ &\quad - \max_{j \in S} \left\| \mathbf{D}^*[j]\mathbf{w} + \sum_{i \in S \setminus \{j\}} \mathbf{D}^*[j]\mathbf{D}[i]\mathbf{x}[i] \right\|_2. \end{aligned} \quad (57)$$

As we have seen in the proof of Lemma 1, the eigenvalues of $\mathbf{D}^*[j]\mathbf{D}[j]$ are bounded in the range $[1-(d-1)\nu, 1+(d-1)\nu]$. Consequently

$$\begin{aligned} \max_{j \in S} \|\mathbf{D}^*[j]\mathbf{D}[j]\mathbf{x}[j]\|_2 &\geq \max_{j \in S} (1-(d-1)\nu) \|\mathbf{x}[j]\|_2 \\ &= (1-(d-1)\nu) |x_{\max}|. \end{aligned} \quad (58)$$

Combining this result with (57), we have

$$\begin{aligned} \max_{j \in S} \|\mathbf{D}^*[j]\mathbf{y}\|_2 &\geq (1-(d-1)\nu) |x_{\max}| \\ &\quad - \max_{j \in S} \sum_{i \in S \setminus \{j\}} \|\mathbf{D}^*[j]\mathbf{D}[i]\mathbf{x}[i]\|_2 - \max_{j \in S} \|\mathbf{D}^*[j]\mathbf{w}\|_2. \end{aligned} \quad (59)$$

Together with (50) and (41), this implies that

$$\begin{aligned} \max_{j \in S} \|\mathbf{D}^*[j]\mathbf{y}\|_2 &> (1-(d-1)\nu) |x_{\max}| - (k-1) |x_{\max}| d\mu_B - \tau \\ &= (1-(d-1)\nu) |x_{\max}| - (2k-1) |x_{\max}| d\mu_B - 2\tau \\ &\quad + k |x_{\max}| d\mu_B + \tau. \end{aligned} \quad (60)$$

Merging the results (56) and (60) yields

$$\begin{aligned} \max_{j \in S} \|\mathbf{D}^*[j]\mathbf{y}\|_2 &> \max_{j \notin S} \|\mathbf{D}^*[j]\mathbf{y}\|_2 \\ &\quad + (1-(d-1)\nu) |x_{\max}| - (2k-1) |x_{\max}| d\mu_B - 2\tau. \end{aligned} \quad (61)$$

Consequently, if (51) holds, then (52) follows, as required.

In a similar fashion, observe that

$$\begin{aligned} \min_{j \in S} \|\mathbf{D}^*[j]\mathbf{y}\|_2 &= \min_{j \in S} \left\| \sum_{i \in S} \mathbf{D}^*[j]\mathbf{D}[i]\mathbf{x}[i] + \mathbf{D}^*[j]\mathbf{w} \right\|_2 \\ &\geq \min_{j \in S} \|\mathbf{D}^*[j]\mathbf{D}[j]\mathbf{x}[j]\|_2 \\ &\quad - \max_{j \in S} \sum_{i \in S \setminus \{j\}} \|\mathbf{D}^*[j]\mathbf{D}[i]\mathbf{x}[i]\|_2 - \|\mathbf{D}^*[j]\mathbf{w}\|_2. \end{aligned} \quad (62)$$

As noted previously, all eigenvalues of $\mathbf{D}^*[j]\mathbf{D}[j]$ are larger than or equal to $1-(d-1)\nu$, and therefore

$$\min_{j \in S} \|\mathbf{D}^*[j]\mathbf{D}[j]\mathbf{x}[j]\|_2 \geq (1-(d-1)\nu) |x_{\min}|. \quad (63)$$

Furthermore, using (41) we have, for $i \neq j$,

$$\|\mathbf{D}^*[j]\mathbf{D}[i]\mathbf{x}[i]\|_2 \leq \|\mathbf{D}^*[j]\mathbf{D}[i]\| |x_{\max}| \leq d\mu_B |x_{\max}|. \quad (64)$$

Substituting (50), (63), and (64) into (62) provides us with

$$\begin{aligned} \min_{j \in S} \|\mathbf{D}^*[j]\mathbf{y}\|_2 &> (1-(d-1)\nu) |x_{\min}| - (k-1) d\mu_B |x_{\max}| - \tau \\ &= (1-(d-1)\nu) |x_{\min}| - (2k-1) d\mu_B |x_{\max}| - 2\tau \\ &\quad + k d\mu_B |x_{\max}| + \tau. \end{aligned} \quad (65)$$

Finally, using (56) we obtain

$$\begin{aligned} \min_{j \in S} \|\mathbf{D}^*[j]\mathbf{y}\|_2 &> \max_{j \notin S} \|\mathbf{D}^*[j]\mathbf{y}\|_2 \\ &\quad + (1-(d-1)\nu) |x_{\min}| - (2k-1) d\mu_B |x_{\max}| - 2\tau. \end{aligned} \quad (66)$$

Therefore, if the condition (53) is satisfied, then (54) holds, completing the proof. \blacksquare

We are now ready to prove Theorems 1 and 2.

Proof of Theorem 1: Using (10) and (42), we have for all j

$$\|\mathbf{D}^*[j]\mathbf{w}\|_2 \leq \|\mathbf{D}[j]\| \cdot \|\mathbf{w}\|_2 \leq \varepsilon \sqrt{1+(d-1)\nu}. \quad (67)$$

Thus, (50) holds with $\tau = \varepsilon \sqrt{1+(d-1)\nu}$.

In light of (21), the condition (53) for the second part of Lemma 2 holds, and therefore, by Lemma 2, we conclude that (54) holds. It follows that all blocks $\mathbf{D}[i]$ with $i \in S$ are more highly correlated than the off-support blocks $\mathbf{D}[i]$, $i \notin S$. Thus, the estimated support \hat{S} contains the true support set S (with the possible addition of superfluous indices if $s < k$). It follows from the definition (16) of $\hat{\mathbf{x}}_{\text{BTH}}$ that $(\hat{\mathbf{x}}_{\text{BTH}})_{\hat{S}} = \mathbf{D}_{\hat{S}}^\dagger \mathbf{y}$, and thus

$$\begin{aligned} \|\mathbf{x} - \hat{\mathbf{x}}_{\text{BTH}}\|_2^2 &= \|\mathbf{x}_{\hat{S}} - (\hat{\mathbf{x}}_{\text{BTH}})_{\hat{S}}\|_2^2 \\ &= \|\mathbf{D}_{\hat{S}}^\dagger \mathbf{D}_{\hat{S}} \mathbf{x}_{\hat{S}} - \mathbf{D}_{\hat{S}}^\dagger \mathbf{y}\|_2^2 \\ &\leq \|\mathbf{D}_{\hat{S}}^\dagger\|^2 \cdot \|\mathbf{y} - \mathbf{D}_{\hat{S}} \mathbf{x}\|_2^2 \\ &= \|\mathbf{D}_{\hat{S}}^\dagger\|^2 \cdot \|\mathbf{w}\|_2^2 \end{aligned} \quad (68)$$

where we have used the fact that $\mathbf{D}_{\hat{S}}^\dagger \mathbf{D}_{\hat{S}} = \mathbf{I}$, which follows from our assumption that \mathbf{D}_I has full row rank for any set I of size s (see Section II).

Since $|x_{\min}| \leq |x_{\max}|$, it follows from (21) that

$$1-(d-1)\nu > (2k-1) d\mu_B. \quad (69)$$

Therefore, we may apply (44), yielding

$$\begin{aligned} \|\mathbf{D}_{\hat{S}}^\dagger\|^2 &= \|(\mathbf{D}_{\hat{S}}^* \mathbf{D}_{\hat{S}})^{-1}\| \\ &\leq \frac{1}{1-(d-1)\nu - (k-1) d\mu_B}. \end{aligned} \quad (70)$$

Combining this result with (68) and using (10), we obtain (22), as required. \blacksquare

Proof of Theorem 2: As shown in the proof of Theorem 1, it follows from (10) that (50) holds with $\tau = \varepsilon \sqrt{1+(d-1)\nu}$. From (23) we then have

$$(1-(d-1)\nu) |x_{\min}| > 2\tau + (2k-1) d\mu_B |x_{\min}|. \quad (71)$$

Since $|x_{\max}| \geq |x_{\min}|$, this implies the condition (51) for the first part of Lemma 2. Thus, by Lemma 2, the dictionary block most highly correlated with \mathbf{y} is a block within the support S

of \mathbf{x} . In other words, the first iteration in the BOMP algorithm correctly identifies an element within the support S .

The proof continues by induction. Assume we have reached the ℓ th iteration with $2 \leq \ell \leq s$ and that all previous iterations have correctly identified elements of S . In other words, using the notation of Section III, we have $i_1, \dots, i_{\ell-1} \in S$.

By definition, we now have

$$\mathbf{r}^\ell = \mathbf{y} - \mathbf{D}\mathbf{x}^{\ell-1} = \mathbf{D}\tilde{\mathbf{x}}^{\ell-1} + \mathbf{w} \quad (72)$$

where $\tilde{\mathbf{x}}^{\ell-1} \triangleq \mathbf{x} - \mathbf{x}^{\ell-1}$ is the estimation error after $\ell-1$ iterations. Since $\text{supp}(\mathbf{x}) = S$ and, by induction, $\text{supp}(\mathbf{x}^{\ell-1}) \subset S$, we have $\text{supp}(\tilde{\mathbf{x}}^{\ell-1}) \subset S$. Furthermore, $\ell-1 < s$, so that $\text{supp}(\mathbf{x}^{\ell-1})$ contains less than s elements, and is thus a strict subset of S . It follows that at least one nonzero block in $\tilde{\mathbf{x}}^{\ell-1}$ is equal to the corresponding block in \mathbf{x} . Therefore

$$\max_j \|\tilde{\mathbf{x}}^{\ell-1}[j]\|_2 \geq |x_{\min}|. \quad (73)$$

To summarize, by (72), \mathbf{r}^ℓ can be thought of as a noisy measurement of the block sparse vector $\tilde{\mathbf{x}}^{\ell-1}$, which contains a block whose norm is at least $|x_{\min}|$. Using (73) and (23), we find that the condition (51) holds for this modified estimation problem. Consequently, by Lemma 2, we have

$$\max_{j \in S} \|\mathbf{D}^*[j]\mathbf{r}^{\ell-1}\|_2 > \max_{j \notin S} \|\mathbf{D}^*[j]\mathbf{r}^{\ell-1}\|_2. \quad (74)$$

Therefore, by (17), the ℓ th iteration of the BOMP algorithm will choose an index i_ℓ belonging to the correct support set S , as long as $\ell \leq s$.

Since the BOMP algorithm never chooses the same support element twice, we conclude that precisely the s elements of S will be identified in the first s iterations. If $s < k$, then the remaining iterations will identify some additional elements not in S , so that ultimately the estimated support set $\hat{S} = \{i_1, \dots, i_k\}$ will satisfy $\hat{S} \supseteq S$. The estimate $\hat{\mathbf{x}}_{\text{BOMP}}$ therefore satisfies $(\hat{\mathbf{x}}_{\text{BOMP}})_{\hat{S}} = \mathbf{D}_{\hat{S}}^\dagger \mathbf{y}$. Following the procedure (68)–(70) in the proof of Theorem 1, we obtain in an identical manner the required result (24). ■

APPENDIX B PROOF OF THEOREM 3

To compute the CRB, we must first determine the Fisher information matrix $\mathbf{J}(\mathbf{x})$ for estimating \mathbf{x} from \mathbf{y} of (8). This can be done using a standard formula [23, p. 85] and yields

$$\mathbf{J}(\mathbf{x}) = \frac{1}{\sigma^2} \mathbf{D}^* \mathbf{D}. \quad (75)$$

We now identify, for each $\mathbf{x} \in \mathfrak{X}$, an orthonormal basis for the feasible direction subspace, which is defined as the smallest subspace of \mathbb{C}^N containing all feasible directions at \mathbf{x} . To this end, denote by \mathbf{e}_i the i th column of the $N \times N$ identity matrix. Consider first points $\mathbf{x} \in \mathfrak{X}$ for which $s < k$. In other words, these are parameter values whose support S contains fewer than k elements. For such values of \mathbf{x} , we have, for any ε and any $1 \leq i \leq N$,

$$|\text{supp}(\mathbf{x} + \varepsilon \mathbf{e}_i)| \leq |S| + 1 < k + 1 \leq k \quad (76)$$

and therefore $\mathbf{x} + \varepsilon \mathbf{e}_i \in \mathfrak{X}$ for any ε and for any i . Consequently, the set of feasible directions at \mathbf{x} includes

$\{\mathbf{e}_1, \dots, \mathbf{e}_N\}$, and the feasible direction subspace is therefore \mathbb{C}^N itself. Thus, for values \mathbf{x} containing fewer than k nonzero blocks, a convenient choice of a basis for the feasible direction subspace consists of the columns of the identity matrix.

Next, consider maximal-support parameter values, i.e., vectors \mathbf{x} for which $s = k$. It is now no longer possible to add any vector \mathbf{e}_i to \mathbf{x} without violating the constraints. Indeed, it is not difficult to see that the only feasible directions are linear combinations of the unit vectors \mathbf{e}_i for which i belongs to one of the blocks in S . These unit vectors can thus be chosen as a basis for the feasible direction subspace.

Let $\mathbf{U}(\mathbf{x})$ be a matrix whose columns comprise the chosen orthonormal basis for the feasible direction subspace at \mathbf{x} . Note that the dimensions of $\mathbf{U}(\mathbf{x})$ change with \mathbf{x} ; specifically, $\mathbf{U}(\mathbf{x}) = \mathbf{I}_{N \times N}$ when $|S| < k$, and $\mathbf{U}(\mathbf{x})$ is an $N \times sd$ matrix otherwise. A necessary condition for a finite-variance \mathfrak{X} -unbiased estimator to exist at a point \mathbf{x} is [27, Thm. 1]

$$\mathcal{R}(\mathbf{U}(\mathbf{x})\mathbf{U}^*(\mathbf{x})) \subseteq \mathcal{R}(\mathbf{U}(\mathbf{x})\mathbf{U}^*(\mathbf{x})\mathbf{J}(\mathbf{x})\mathbf{U}(\mathbf{x})\mathbf{U}^*(\mathbf{x})). \quad (77)$$

When $s < k$, we have $\mathbf{U}(\mathbf{x}) = \mathbf{I}$. In this case, using (75), the condition (77) becomes

$$\mathbb{C}^N \subseteq \mathcal{R}(\mathbf{J}(\mathbf{x})) = \mathcal{R}(\mathbf{D}^* \mathbf{D}). \quad (78)$$

Since the dimensions of \mathbf{D} are $L \times N$ with $L < N$, the rank of $\mathbf{D}^* \mathbf{D}$ is at most L , and thus $\mathcal{R}(\mathbf{D}^* \mathbf{D})$ cannot include the entire space \mathbb{C}^N . We conclude that in this case, (77) does not hold, and therefore no \mathfrak{X} -unbiased estimator exists at points \mathbf{x} for which $|S| < s$, proving part (a) of the theorem.

Let us now turn to maximal-support parameter values \mathbf{x} . As we have seen above, in this case the matrix $\mathbf{U}(\mathbf{x})$ consists of the columns \mathbf{e}_i for which i is an element of a block within the support of \mathbf{x} . Therefore, the product $\mathbf{D}\mathbf{U}(\mathbf{x})$ selects those atoms of \mathbf{D} belonging to blocks within S , i.e., $\mathbf{D}\mathbf{U}(\mathbf{x}) = \mathbf{D}_S$. Using (75), this leads to

$$\mathbf{U}^*(\mathbf{x})\mathbf{J}(\mathbf{x})\mathbf{U}(\mathbf{x}) = \frac{1}{\sigma^2} \mathbf{D}_S^* \mathbf{D}_S \quad (79)$$

which is invertible by assumption (see Section II). It follows that the condition (77) holds for maximal-support parameters \mathbf{x} . One can therefore apply [27, Thm. 1], which states that for such values of \mathbf{x} ,

$$\text{MSE}(\hat{\mathbf{x}}, \mathbf{x}) \geq \text{Tr} \left(\mathbf{U}(\mathbf{x}) (\mathbf{U}^*(\mathbf{x})\mathbf{J}(\mathbf{x})\mathbf{U}(\mathbf{x}))^\dagger \mathbf{U}^*(\mathbf{x}) \right). \quad (80)$$

Combining with (79) and using the fact that $\mathbf{U}^*(\mathbf{x})\mathbf{U}(\mathbf{x}) = \mathbf{I}$, we obtain (31), proving part (b) of the theorem.

APPENDIX C PROOFS FOR GAUSSIAN NOISE

We begin with two lemmas which prove some useful properties of the Gaussian distribution. The first of these is a generalization of a result due to Šidák [31].

Lemma 3. *Let $\mathbf{v}_1, \dots, \mathbf{v}_M$ be a set of M jointly Gaussian random vectors. Suppose that $\mathbb{E}\{\mathbf{v}_i\} = \mathbf{0}$ for all i , but that the covariances of the vectors are unspecified and that the vectors*

are not necessarily independent. We then have

$$\begin{aligned} & \Pr\{\|\mathbf{v}_1\|_2 \leq c_1, \|\mathbf{v}_2\|_2 \leq c_2, \dots, \|\mathbf{v}_M\|_2 \leq c_M\} \\ & \geq \Pr\{\|\mathbf{v}_1\|_2 \leq c_1\} \cdot \Pr\{\|\mathbf{v}_2\|_2 \leq c_2\} \cdots \\ & \quad \cdots \Pr\{\|\mathbf{v}_M\|_2 \leq c_M\}. \end{aligned} \quad (81)$$

Proof: We will demonstrate that

$$\begin{aligned} & \Pr\{\|\mathbf{v}_1\|_2 \leq c_1, \|\mathbf{v}_2\|_2 \leq c_2, \dots, \|\mathbf{v}_M\|_2 \leq c_M\} \\ & \geq \Pr\{\|\mathbf{v}_1\|_2 \leq c_1\} \Pr\{\|\mathbf{v}_2\|_2 \leq c_2, \dots, \|\mathbf{v}_M\|_2 \leq c_M\}. \end{aligned} \quad (82)$$

The result then follows by induction. For simplicity of notation, we will prove that (82) holds for the case $M = 2$; the general result can be shown in the same manner.

Denote by $f(\mathbf{v}_1|\mathbf{v}_2)$ the pdf of \mathbf{v}_1 conditioned on \mathbf{v}_2 . Observe that, for a deterministic value \mathbf{w} , the pdf $f(\mathbf{v}_1|\mathbf{w})$ defines a Gaussian random vector whose mean depends linearly on \mathbf{w} , but whose covariance is constant in \mathbf{w} . Therefore, using a result due to Anderson [32], it follows that

$$\Pr\{\|\mathbf{v}_1\|_2 \leq c_1 | \mathbf{v}_2 = \alpha \mathbf{w}\} = \int_{\|\mathbf{u}_1\|_2 \leq c_1} f(\mathbf{u}_1 | \alpha \mathbf{w}) d\mathbf{u} \quad (83)$$

is a non-increasing function of α .

Next, denoting by $f(\mathbf{v}_2)$ the marginal pdf of \mathbf{v}_2 , we have

$$\begin{aligned} a(c_1, c_2) & \triangleq \Pr\{\|\mathbf{v}_1\|_2 \leq c_1 | \|\mathbf{v}_2\|_2 \leq c_2\} \\ & = \frac{\int_{\|\mathbf{u}\|_2 \leq c_1} \int_{\|\mathbf{w}\|_2 \leq c_2} f(\mathbf{u}|\mathbf{w}) f(\mathbf{w}) d\mathbf{w} d\mathbf{u}}{\Pr\{\|\mathbf{v}_2\|_2 \leq c_2\}} \\ & = \frac{\int_{\|\mathbf{w}\|_2 \leq c_2} \Pr\{\|\mathbf{v}_1\|_2 \leq c_1 | \mathbf{v}_2 = \mathbf{w}\} f(\mathbf{w}) d\mathbf{w}}{\int_{\|\mathbf{w}\|_2 \leq c_2} f(\mathbf{w}) d\mathbf{w}}. \end{aligned} \quad (84)$$

Thus, the function $a(c_1, c_2)$ is a weighted average of expressions of the form $\Pr\{\|\mathbf{v}_1\|_2 \leq c_1 | \mathbf{v}_2 = \mathbf{w}\}$ for values of \mathbf{w} satisfying $\|\mathbf{w}\|_2 \leq c_2$. However, as we have shown, $\Pr\{\|\mathbf{v}_1\|_2 \leq c_1 | \mathbf{v}_2 = \mathbf{w}\}$ is non-increasing in $\|\mathbf{w}\|_2$. Consequently, $a(c_1, c_2)$ is non-increasing in c_2 .

On the other hand, observe that as $c_2 \rightarrow \infty$, the probability of the event $\|\mathbf{v}_2\|_2 \leq c_2$ converges to 1. Thus we have

$$\lim_{c_2 \rightarrow \infty} a(c_1, c_2) = \Pr\{\|\mathbf{v}_1\|_2 \leq c_1\}. \quad (85)$$

Combined with the fact that $a(c_1, c_2)$ is non-increasing in c_2 , we find that

$$a(c_1, c_2) \geq \Pr\{\|\mathbf{v}_1\|_2 \leq c_1\} \quad \text{for all } c_1, c_2. \quad (86)$$

Using the definition of $a(c_1, c_2)$ and applying Bayes's rule, we obtain

$$\begin{aligned} & \Pr\{\|\mathbf{v}_1\|_2 \leq c_1, \|\mathbf{v}_2\|_2 \leq c_2\} \\ & \geq \Pr\{\|\mathbf{v}_1\|_2 \leq c_1\} \Pr\{\|\mathbf{v}_2\|_2 \leq c_2\} \end{aligned} \quad (87)$$

and thus complete the proof. \blacksquare

Our next lemma bounds the tail probability of the chi-squared distribution.

Lemma 4. Let \mathbf{u} be a d -dimensional Gaussian random vector having mean zero and covariance \mathbf{I} . Then, for any $t \geq 1$, we have

$$\Pr\{\|\mathbf{u}\|_2^2 \geq t^2\} \leq \frac{(d-2)!! [d/2]}{2^{d/2-1} \Gamma(d/2)} t^{d-2} e^{-t^2/2} \quad (88a)$$

$$\leq 0.8dt^{d-2} e^{-t^2/2} \quad (88b)$$

where $\Gamma(z) \triangleq \int_0^\infty t^{z-1} e^{-t} dt$ is the Gamma function and

$$n!! \triangleq \prod_{0 \leq i < n/2} (n-2i) \quad (89)$$

is the double factorial operator.

Of the two bounds provided in (88), the first is somewhat tighter, but obviously more cumbersome. For analytical tractability, we will use the latter bound in the sequel.

Proof of Lemma 4: The expression $\|\mathbf{u}\|_2^2$ is distributed as a chi-squared random variable with d degrees of freedom. Therefore, its tail probability is given by [33, §16.3]

$$\Pr\{\|\mathbf{u}\|_2^2 \geq t^2\} = \frac{\Gamma(d/2, t^2/2)}{\Gamma(d/2)} \quad (90)$$

where $\Gamma(a, z)$ is the incomplete Gamma function $\Gamma(a, z) \triangleq \int_z^\infty t^{a-1} e^{-t} dt$. It follows from the series expansion of $\Gamma(a, z)$ that [34, §6.5.32]

$$\begin{aligned} \Gamma\left(\frac{d}{2}, \frac{t^2}{2}\right) & \leq \frac{e^{-t^2/2}}{2^{d/2-1} t^2} [t^d + (d-2)t^{d-2} \\ & \quad + (d-2)(d-4)t^{d-4} + \cdots + (d-2)!! t^m] \end{aligned} \quad (91)$$

where $m = 1$ when d is odd and $m = 2$ when d is even. Note that (91) holds with equality for even d , but the inequality is strict for odd d . Since $t \geq 1$, we can enlarge each of the terms in the square brackets in (91) by replacing it with $(d-2)!! t^d$. The total number of terms in brackets is $\lceil d/2 \rceil$, yielding

$$\Gamma\left(\frac{d}{2}, \frac{t^2}{2}\right) \leq \frac{e^{-t^2/2}}{2^{d/2-1}} t^{d-2} (d-2)!! \left[\frac{d}{2}\right]. \quad (92)$$

Substituting into (90) demonstrates (88a).

To prove (88b), we distinguish between even and odd values of d . Assume first that d is even and denote $d = 2p$. We then have

$$\Gamma(d/2) = \Gamma(p) = (p-1)! \quad (93)$$

and

$$(d-2)!! = (2p-2)!! = 2^{p-1} (p-1)!. \quad (94)$$

Substituting these values into (88a) and simplifying yields

$$\Pr\{\|\mathbf{u}\|_2^2 \geq t^2\} \leq \frac{d}{2} t^{d-2} e^{-t^2/2} \quad (95)$$

which clearly satisfies (88b).

Similarly, assume that d is odd and write $d = 2p + 1$. Substituting the formula

$$\Gamma(d/2) = \Gamma(p + 1/2) = \frac{(2p-1)!! \sqrt{\pi}}{2^p} \quad (96)$$

into (88a), we obtain

$$\Pr\{\|\mathbf{u}\|_2^2 \geq t^2\} \leq \sqrt{\frac{2}{\pi}} \frac{d+1}{2} t^{d-2} e^{-t^2/2}. \quad (97)$$

It is easily verified that

$$\sqrt{\frac{2}{\pi}} \frac{d+1}{2} \leq 0.8d \quad \text{for all } d \geq 1. \quad (98)$$

Substituting back into (97) yields the required result. \blacksquare

Our next result applies more specifically to the block sparse estimation setting. Following [4], [21], we consider the event

$$B = \left\{ \max_{1 \leq i \leq M} \|\mathbf{D}^*[i]\mathbf{w}\|_2^2 \leq \tau^2 \right\} \quad (99)$$

where

$$\tau^2 = 2d\sigma\alpha(1 + (d-1)\nu) \log N \quad (100)$$

for a given $\alpha > 1/(2d \log N)$. We then have the following lemma.

Lemma 5. *Under the setting of Section II, assume that \mathbf{w} is a Gaussian random vector with mean zero and covariance $\sigma^2 \mathbf{I}$. Then, the probability of the event B of (99) is bounded by*

$$\Pr\{B\} \geq 1 - \frac{0.8(2\alpha d \log N)^{d/2-1}}{N^{\alpha d-1}}. \quad (101)$$

Proof: Observe that $\mathbf{D}^*[i]\mathbf{w}$ is a d -dimensional Gaussian random vector with mean zero and covariance $\sigma^2 \mathbf{D}^*[i]\mathbf{D}[i]$. Therefore, the random vector

$$\mathbf{u} = \frac{1}{\sigma} (\mathbf{D}^*[i]\mathbf{D}[i])^{-1/2} \mathbf{D}^*[i]\mathbf{w} \quad (102)$$

is a d -dimensional Gaussian random vector with mean zero and covariance \mathbf{I} . We thus have

$$\begin{aligned} \Pr\{\|\mathbf{D}^*[i]\mathbf{w}\|_2^2 \leq \tau^2\} &= \Pr\left\{\sigma^2 \|(\mathbf{D}^*[i]\mathbf{D}[i])^{1/2} \mathbf{u}\|_2^2 \leq \tau^2\right\} \\ &\geq \Pr\left\{\sigma^2 \|\mathbf{D}^*[i]\mathbf{D}[i]\| \cdot \|\mathbf{u}\|_2^2 \leq \tau^2\right\} \\ &\geq \Pr\left\{\|\mathbf{u}\|_2^2 \leq \frac{\tau^2}{\sigma^2(1+(d-1)\nu)}\right\} \end{aligned} \quad (103)$$

where, in the last step, we used (42). Using Lemma 4 and substituting the value (100) of τ^2 , we obtain

$$\Pr\{\|\mathbf{D}^*[i]\mathbf{w}\|_2^2 \leq \tau^2\} \geq 1 - \eta \quad (104)$$

where

$$\begin{aligned} \eta &\triangleq 1 - 0.8d(2\alpha d \log N)^{d/2-1} \exp(-d\alpha \log N) \\ &= 1 - \frac{0.8d(2\alpha d \log N)^{d/2-1}}{N^{\alpha d}}. \end{aligned} \quad (105)$$

Using Lemma 3, we have

$$\begin{aligned} \Pr\{B\} &\geq \prod_{i=1}^M \Pr\{\|\mathbf{D}^*[i]\mathbf{w}\|_2^2 \leq \tau^2\} \\ &= (1 - \eta)^M. \end{aligned} \quad (106)$$

When $\eta > 1$, the bound (101) is meaningless and the theorem holds vacuously. Otherwise, when $\eta \leq 1$, we have

$$\Pr\{B\} \geq 1 - M\eta \quad (107)$$

where we used the fact that $(1 - \eta)^M \geq 1 - M\eta$ whenever $\eta \leq 1$ and $M \geq 1$. Substituting the value of η from (105) and recalling that $N = Md$ yields the required result. \blacksquare

We are now ready to prove Theorems 4 and 5.

Proof of Theorem 4: By Lemma 5, the event B of (99) occurs with probability exceeding (34). Furthermore, using (33), it follows from Lemma 2 that under the event B , all blocks in the correct support set S are more highly correlated with \mathbf{y} than the off-support blocks. Consequently, when B occurs, we have $S \subseteq \hat{S}$, where \hat{S} is the support estimated by the BTH algorithm. Note, however, that the estimated set \hat{S} will contain additional blocks not in S if $s < k$. It follows that

$$\begin{aligned} \|\mathbf{x} - \hat{\mathbf{x}}_{\text{BTH}}\|_2^2 &= \|\mathbf{x}_{\hat{S}} - (\hat{\mathbf{x}}_{\text{BTH}})_{\hat{S}}\|_2^2 \\ &= \|\mathbf{D}_{\hat{S}}^\dagger \mathbf{D}_{\hat{S}} \mathbf{x}_{\hat{S}} - \mathbf{D}_{\hat{S}}^\dagger \mathbf{D}_{\hat{S}} \mathbf{y}\|_2^2 \\ &\leq \|(\mathbf{D}_{\hat{S}}^* \mathbf{D}_{\hat{S}})^{-1}\|^2 \cdot \|\mathbf{D}_{\hat{S}}^* \mathbf{w}\|_2^2 \\ &\leq \|(\mathbf{D}_{\hat{S}}^* \mathbf{D}_{\hat{S}})^{-1}\|^2 \cdot \sum_{i \in \hat{S}} \|\mathbf{D}^*[i]\mathbf{w}\|_2^2 \end{aligned} \quad (108)$$

where we have used the fact that $\mathbf{D}_{\hat{S}}^\dagger \mathbf{D}_{\hat{S}} = \mathbf{I}$, which is a consequence of the assumption that $\mathbf{D}_{\hat{S}}$ has full row rank (see Section II). Using (44) and (99), we have that when B occurs

$$\|\mathbf{x} - \hat{\mathbf{x}}_{\text{BTH}}\|_2^2 \leq \frac{k\tau^2}{(1 - (d-1)\nu - (k-1)d\mu_B)^2}. \quad (109)$$

Substituting the value (100) of τ yields the required result (35). \blacksquare

Proof of Theorem 5: It follows from Lemma 5 that the event B occurs with probability exceeding (34). Our goal in this proof will thus be to show that, if B does occur, then the BOMP algorithm correctly identifies all elements of the support S of \mathbf{x} (although some off-support elements may be identified as well if $s < k$). The remainder of the proof will then follow the steps of the proof of Theorem 4.

To demonstrate that the correct support is recovered, we begin by analyzing the first iteration of the BOMP algorithm. This iteration chooses a block i_1 having maximal correlation $\|\mathbf{D}^*[i_1]\mathbf{y}\|_2$ with the measurements \mathbf{y} . Now, since $|x_{\max}| \geq |x_{\min}|$, the condition (36) implies (51), with τ given by (100). Consequently, by Lemma 2, under the event B we find that the first iteration of BOMP identifies an element i_1 in the correct support set S .

To show that the next $s-1$ iterations of the BOMP algorithm also identify support elements, we proceed by induction. Specifically, assume that $\ell-1 < s$ iterations have correctly identified elements $i_1, \dots, i_{\ell-1}$, all of which are in the support set S . As in the proof of Theorem 2, define the estimation error after $\ell-1$ iterations as $\tilde{\mathbf{x}}^{\ell-1} \triangleq \mathbf{x} - \mathbf{x}^{\ell-1}$. By the induction hypothesis, $\text{supp}(\tilde{\mathbf{x}}) \subset S$, and clearly $\text{supp}(\mathbf{x}) = S$. Thus $\text{supp}(\tilde{\mathbf{x}}) \subset S$, i.e., the support of $\tilde{\mathbf{x}}$ is a strict subset of S . Using the same arguments as in the proof of Theorem 2, we find that $\tilde{\mathbf{x}}^{\ell-1}$ contains a block whose norm is at least $|x_{\min}|$. Therefore, we can consider a modified estimation problem, in which \mathbf{r}^ℓ is a noisy measurement vector of the block sparse signal $\tilde{\mathbf{x}}^{\ell-1}$. Together with (36), this implies that (51) holds for the modified setting. Therefore, by (52), the block in \mathbf{r}^ℓ having maximal correlation with the measurements is an element of S . Consequently, BOMP will correctly identify a support element in the ℓ th iteration. Since the BOMP algorithm never selects a previously chosen support element, we find by induction that

the support set S will be identified in full after s iterations. If $s < k$, then the remaining $k - s$ iterations will identify arbitrary off-support elements.

Denoting by \hat{S} the complete k -element support set identified by the BOMP approach, we thus have $S \subseteq \hat{S}$. Following the technique (108)–(109) used in the proof of Theorem 4 thus yields the required result (37). ■

REFERENCES

- [1] D. L. Donoho, "Compressed sensing," *IEEE Trans. Inf. Theory*, vol. 52, no. 4, pp. 1289–1306, 2006.
- [2] E. J. Candès, J. K. Romberg, and T. Tao, "Stable signal recovery from incomplete and inaccurate measurements," *Comm. Pure Appl. Math.*, vol. LIX, pp. 1207–1223, 2006.
- [3] A. M. Bruckstein, D. L. Donoho, and M. Elad, "From sparse solutions of systems of equations to sparse modeling of signals and images," *SIAM Review*, vol. 51, no. 1, pp. 34–81, Feb. 2009.
- [4] J. A. Tropp, "Just relax: Convex programming methods for identifying sparse signals in noise," *IEEE Trans. Inf. Theory*, vol. 52, no. 3, pp. 1030–1051, 2006.
- [5] D. L. Donoho, M. Elad, and V. N. Temlyakov, "Stable recovery of sparse overcomplete representations in the presence of noise," *IEEE Trans. Inf. Theory*, vol. 52, no. 1, pp. 6–18, 2006.
- [6] M. Elad and M. Aharon, "Image denoising via sparse and redundant representations over learned dictionaries," *IEEE Trans. Image Process.*, vol. 15, no. 12, pp. 3736–3745, Dec. 2006.
- [7] M. M. Bronstein, A. M. Bronstein, M. Zibulevsky, and Y. Y. Zeevi, "Blind deconvolution of images using optimal sparse representations," *IEEE Trans. Image Process.*, vol. 14, no. 6, pp. 726–736, Jun. 2005.
- [8] M. J. Fadili, J.-L. Starck, and F. Murtagh, "Inpainting and zooming using sparse representations," *The Computer Journal*, vol. 52, no. 1, pp. 64–79, 2009.
- [9] Y. M. Lu and M. N. Do, "A theory for sampling signals from a union of subspaces," *IEEE Trans. Signal Process.*, vol. 56, no. 6, pp. 2334–2345, Jun. 2008.
- [10] T. Blumensath and M. E. Davies, "Sampling theorems for signals from the union of finite-dimensional linear subspaces," *IEEE Trans. Inf. Theory*, vol. 55, no. 4, pp. 1872–1882, Apr. 2009.
- [11] Y. C. Eldar and M. Mishali, "Robust recovery of signals from a structured union of subspaces," *IEEE Trans. Inf. Theory*, vol. 55, no. 11, pp. 5302–5316, Nov. 2009.
- [12] K. Gedalyahu and Y. C. Eldar, "Time delay estimation from low rate samples: A union of subspaces approach," *IEEE Trans. Signal Process.*, vol. 58, no. 6, pp. 3017–3031, Jun. 2010.
- [13] K. Gedalyahu, R. Tur, and Y. C. Eldar, "Multichannel sampling of pulse streams at the rate of innovation," *IEEE Trans. Signal Process.*, submitted. [Online]. Available: <http://arxiv.org/abs/1004.5070>
- [14] M. Mishali and Y. C. Eldar, "From theory to practice: Sub-Nyquist sampling of sparse wideband analog signals," *IEEE J. Sel. Topics Signal Process.*, vol. 4, no. 2, pp. 375–391, Apr. 2010.
- [15] —, "Blind multi-band signal reconstruction: Compressed sensing for analog signals," *IEEE Trans. Signal Process.*, vol. 57, no. 3, pp. 993–1009, Mar. 2009.
- [16] M. Mishali, Y. C. Eldar, O. Dounaevsky, and E. Shoshan, "Xampling: Analog to digital at sub-Nyquist rates," *IET J. Circuits, Devices and Systems*, 2010, to appear. [Online]. Available: <http://arxiv.org/abs/0912.2495>
- [17] M. Yuan and Y. Lin, "Model selection and estimation in regression with grouped variables," *J. Royal Statist. Soc. B*, vol. 68, pp. 49–67, 2006.
- [18] M. Stojnic, F. Parvaresh, and B. Hassibi, "On the reconstruction of block-sparse signals with an optimal number of measurements," *IEEE Trans. Signal Process.*, vol. 57, no. 8, pp. 3075–3085, Aug. 2009.
- [19] Y. C. Eldar, P. Kuppinger, and H. Bölcskei, "Block-sparse signals: Uncertainty relations and efficient recovery," *IEEE Trans. Signal Process.*, vol. 58, no. 6, pp. 3042–3054, Jun. 2010.
- [20] E. Candès and T. Tao, "The Dantzig selector: Statistical estimation when p is much larger than n ," *Ann. Statist.*, vol. 35, no. 6, pp. 2313–2351, 2007, with discussion.
- [21] Z. Ben-Haim, Y. C. Eldar, and M. Elad, "Coherence-based performance guarantees for estimating a sparse vector under random noise," *IEEE Trans. Signal Process.*, 2010, to appear. [Online]. Available: <http://arxiv.org/abs/0903.4579>
- [22] Y. C. Pati, R. Rezaifar, and P. S. Krishnaprasad, "Orthogonal matching pursuit: Recursive function approximation with applications to wavelet decomposition," in *Proc. 27th Asilomar Conf. Signals, Systems, and Computers*, Nov. 1993, pp. 40–44.
- [23] S. M. Kay, *Fundamentals of Statistical Signal Processing: Estimation Theory*. Englewood Cliffs, NJ: Prentice Hall, 1993.
- [24] J. D. Gorman and A. O. Hero, "Lower bounds for parametric estimation with constraints," *IEEE Trans. Inf. Theory*, vol. 26, no. 6, pp. 1285–1301, Nov. 1990.
- [25] P. Stoica and B. C. Ng, "On the Cramér–Rao bound under parametric constraints," *IEEE Signal Process. Lett.*, vol. 5, no. 7, pp. 177–179, 1998.
- [26] Z. Ben-Haim and Y. C. Eldar, "On the constrained Cramér–Rao bound with a singular Fisher information matrix," *IEEE Signal Process. Lett.*, vol. 16, no. 6, pp. 453–456, Jun. 2009.
- [27] —, "The Cramér–Rao bound for estimating a sparse parameter vector," *IEEE Trans. Signal Process.*, vol. 58, no. 6, pp. 3384–3389, Jun. 2010.
- [28] Z. Ben-Haim, Y. C. Eldar, and M. Elad, "Coherence-based near-oracle performance guarantees for sparse estimation under Gaussian noise," in *Proc. Int. Conf. Acoustics, Speech, and Signal Processing (ICASSP 2010)*, Dallas, TX, Mar. 2010, pp. 3590–3593.
- [29] G. H. Golub and C. F. Van Loan, *Matrix Computations*, 3rd ed. Baltimore, MD: Johns Hopkins University Press, 1996.
- [30] D. G. Feingold and R. S. Varga, "Block diagonally dominant matrices and generalizations of the Gerschgorin circle theorem," *Pacific J. Math.*, vol. 12, no. 4, pp. 1241–1250, 1962.
- [31] Z. Šidák, "Rectangular confidence regions for the means of multivariate normal distributions," *J. Amer. Statist. Assoc.*, vol. 62, no. 318, pp. 626–633, Jun. 1967.
- [32] T. W. Anderson, "The integral of a symmetric unimodal function over a symmetric convex set and some probability inequalities," *Proc. Am. Math. Soc.*, vol. 6, pp. 170–176, 1955.
- [33] A. Stuart and J. K. Ord, *Kendall's Advanced Theory of Statistics*, 6th ed. London: Edward Arnold, 1994, vol. 1.
- [34] M. Abramowitz and I. A. Stegun, *Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables*. New York: Dover, 1964.