

# Performance Analysis for Sparse Support Recovery

Gongguo Tang, *Student Member, IEEE*, Arye Nehorai, *Fellow, IEEE*,

## Abstract

In this paper, the performance of estimating the common support for jointly sparse signals based on their projections onto lower-dimensional space is analyzed. Support recovery is formulated as a multiple-hypothesis testing problem and both upper and lower bounds on the probability of error are derived for general measurement matrices, by using the Chernoff bound and Fano's inequality, respectively. The form of the upper bound shows that the performance is determined by a single quantity that is a measure of the incoherence of the measurement matrix, while the lower bound reveals the importance of the total measurement gain. To demonstrate its immediate applicability, the lower bound is applied to derive the minimal number of samples needed for accurate direction of arrival (DOA) estimation for an algorithm based on sparse representation. When applied to Gaussian measurement ensembles, these bounds give necessary and sufficient conditions to guarantee a vanishing probability of error for majority realizations of the measurement matrix. Our results offer surprising insights into the sparse signal reconstruction based on their projections. For example, as far as support recovery is concerned, the well-known bound in compressive sensing is generally not sufficient if the Gaussian ensemble is used. Our study provides an alternative performance measure, one that is natural and important in practice, for signal recovery in compressive sensing as well as other application areas taking advantage of signal sparsity.

## Index Terms

support recovery, jointly sparse signals, compressive sensing, multiple hypothesis testing, probability of error, Chernoff bound, Fano's inequality

This work was supported by the Department of Defense under the Air Force Office of Scientific Research MURI Grant FA9550-05-1-0443, and ONR Grant N000140810849.

## I. INTRODUCTION

Support recovery for jointly sparse signals concerns accurately estimating the non-zero component locations shared by a set of sparse signals based on a limited number of noisy linear observations. More specifically, suppose  $\{x(t) \in \mathbb{F}^N, t = 1, 2, \dots, T\}$ ,  $\mathbb{F} = \mathbb{R}$  or  $\mathbb{C}$ , is a sequence of jointly sparse signals (possibly under a sparsity-inducing basis  $\Phi$  instead of the canonical domain) with a common support  $S$ , which is the index set indicating the non-vanishing signal coordinates. The observation model is linear:

$$y(t) = Ax(t) + w(t) \quad t = 1, 2, \dots, T. \quad (1)$$

In (1),  $A \in \mathbb{F}^{M \times N}$  is the measurement matrix,  $y(t) \in \mathbb{F}^M$  the noisy data vector, and  $w(t) \in \mathbb{F}^M$  an additive noise. In most cases, the sparsity level  $K \triangleq |S|$  and the number of observations  $M$  is far less than  $N$ , the dimension of the ambient space. This problem arises naturally in several signal processing areas such as compressive sensing [1]–[5], source localization [6]–[9], sparse approximation and signal denoising [10].

Compressive Sensing [1]–[5] is a major field that motivates our study for support recovery. In the classical setting of compressive sensing, only one snapshot is considered; *i.e.*,  $T = 1$ . The goal is to recover a long vector  $x := x(1)$  with a small fraction of non-zero coordinates from the much shorter observation vector  $y := y(1)$ . Since most natural signals are compressible under some basis and can be well approximated by their  $K$ -sparse representation [11], this scheme, if properly justified, will reduce the necessary sampling rate beyond the limit set by Nyquist and Shannon [4], [5]. Surprisingly, if  $M = O(K \log(\frac{N}{K})) \ll N$  and the measurement matrix is generated randomly from, for example, a Gaussian distribution, we can recover  $x$  exactly in the noise-free setting by solving a linear programming task. Besides, various methods have been designed for the noise case [12]–[16]. Along with these algorithms, rigorous theoretical analysis is provided to guarantee their effectiveness in terms of, for example, various  $l_p$ -norms between the estimator  $\hat{x}$  and the true value of  $x$  [12]–[16]. However, these results offer no guarantee that we can recover the support of a sparse signal correctly.

The accurate recovery of signal support is crucial to compressive sensing both in theory and in practice. Since for signal recovery it is necessary to have  $K \leq M$ , signal component values can be computed by solving a least-square problem once its support is obtained. Therefore, support recovery is a stronger theoretical criterion than various  $l_p$ -norms. In practice, the success of

compressive sensing in a variety of applications relies on its ability for correct support recovery because the non-zero component indices usually have significant physical meanings. The support of temporally or spatially sparse signals reveals the timing or location for important events such as anomalies. The non-zero indices in the Fourier domain indicate the harmonics existing in a signal [17], which is critical for tasks such as spectrum sensing for cognitive radios [18]. In compressive DNA arrays for bio-sensing, the existence of certain target agents in the tested solution is reflected by the locations of non-vanishing coordinates, while the magnitudes are determined by their concentrations [19]–[22]. For compressive radar imaging, the sparsity constraints are usually imposed on the discretized time–frequency domain. The distance and velocity of an object have a direct correspondence to its locations in the time-frequency domain. The magnitude determined by coefficients of reflection is of less physical significance [23]–[25]. In sparse linear regression [26], the recovered parameter support corresponds to the few factors that explain the data. In all these applications, the support is physically more significant than the component values.

Our study of sparse support recovery is also motivated by the recent reformulation of the source localization problem as one of sparse spectrum estimation. In [6], the authors transform the process of source localization using sensory arrays into the task of estimating the spectrum of a sparse signal by discretizing the parameter manifold. This method exhibits super-resolution in the estimation of direction of arrival (DOA) compared with traditional techniques such as beamforming [27], Capon [28], and MUSIC [29], [30]. Since the basic model employed in [6] applies to several other important problems in signal processing (see [31] and references therein), the principle is readily applicable to those cases. This idea is later generalized and extended to other source localization settings in [7]–[9]. For source localization, the support of the sparse signal reveals the DOA of sources. Therefore, the recovery algorithm’s ability to effect exact support recovery is key to the effectiveness of the method. We also note that usually multiple temporal snapshots are collected, which results in a jointly sparse signal sets as in (1). In addition, since  $M$  is the number of sensors while  $T$  is the number of temporal samples, it is far more expensive to increase  $M$  than  $T$ . The same comments apply to several other examples in the compressive sensing applications discussed in the previous paragraph, especially the compressive DNA arrays, spectrum sensing for cognitive radios, and compressive sensing radar imaging.

The signal recovery problem with joint sparsity constraint [32]–[35], also termed the multiple measurement vector (MMV) problem [36]–[40], has been considered in a line of previous works.

Several algorithms, among them Simultaneous Orthogonal Matching Pursuit (SOMP) [33], [36], [39]; convex relaxation [40];  $\ell_1$ -minimization [37], [38]; and M-FOCUSS [36], are proposed and analyzed, either numerically or theoretically. These algorithms are multiple-dimension extensions of their one-dimension counterparts. Most performance measures of the algorithms are concerned with bounds on various norms of the difference between the true signals and their estimates or their closely related variants. The performance bounds usually involve the mutual coherence between the measurement matrix  $A$  and the basis matrix  $\Phi$  under which the measured signals  $x(t)$  have a jointly sparse representation. However, with joint sparsity constraints, a natural measure of performance would be the model (1)'s potential for correctly identifying the true common support, and hence the algorithm's ability to achieve this potential. As part of their research, J. Chen and X. Huo derived, in a noiseless setting, sufficient conditions on the uniqueness of solutions to (1) under  $\ell_0$  and  $\ell_1$  minimization. In [36], S. Cotter *et. al.* numerically compared the probabilities of correctly identifying the common support by basic matching pursuit, orthogonal matching pursuit, FOCUSS, and regularized FOCUSS in the multiple-measurement setting with a range of SNRs and different numbers of snapshots.

The contribution of our work is threefold. First, we introduce a hypothesis-testing framework to study the performance for multiple support recovery. We employ well-known tools in statistics and information theory such as the Chernoff bound and Fano's inequality to derive both upper and lower bounds for the probability of error. The upper bound we derive is for the *optimal* decision rule, in contrast to a performance analysis for specific sub-optimal reconstruction algorithms [12]–[16]. Hence, the bound can be viewed as a measure of the measurement system's ability to correctly identify the true support. Our bounds isolate important quantities that are crucial for system performance. Since our analysis is based on measurement matrices with as few assumptions as possible, the results can be used as a guidance in system design. Second, we apply these performance bounds to other more specific situations and derive necessary and sufficient conditions in terms of the system parameters to guarantee a vanishing probability of error. In particular, we study necessary conditions for accurate source localization by the mechanism proposed in [6]. By restricting our attention to Gaussian measurement matrices, we derive a result parallel to those for classical compressive sensing [1], [2], namely, the number of measurements that are sufficient for signal reconstruction. Even if we adopt the probability of error as the performance criterion, we get the same bound on  $M$  as in [1], [2]. However, our

result suggests that generally it is impossible to obtain the true support accurately with only one snapshot. We also obtain a necessary condition that shows that the  $\ln \frac{N}{K}$  term cannot be dropped in compressive sensing. Last but not least, in the course of studying the performance bounds we explore the eigenvalue structure of a fundamental matrix in support recovery hypothesis testing for both general measurement matrices and the Gaussian measurement ensemble. These results are of independent interest.

The paper is organized as follows. In Section II, we introduce the mathematical model and briefly review the fundamental ideas in hypothesis testing. Section III is devoted to the derivation of upper bounds on the probability of error for general measurement matrices. We first derive an upper bound on the probability of error for the binary support recovery problem by employing the well-known Chernoff bound in detection theory [41] and extend it to multiple support recovery. We also study the effect of noise on system performance. In Section IV, an information theoretic lower bound is given by using the Fano's inequality, and a necessary condition is shown for the DOA problem considered in [6]. We focus on the Gaussian ensemble in Section V. Necessary and sufficient conditions for system parameters for accurate support recovery are given and their implications discussed. The paper is concluded in Section VI.

## II. NOTATIONS, MODELS, AND PRELIMINARIES

### A. Notations

We first introduce some notations used throughout this paper. Suppose  $x \in \mathbb{F}^N$  is a column vector. We denote by  $S = \text{supp}(x) \subseteq \{1, \dots, N\}$  the support of  $x$ , which is defined as the set of indices corresponding to the non-zero components of  $x$ . For a matrix  $X$ ,  $S = \text{supp}(X)$  denotes the index set of non-zero rows of  $X$ . Here the underlying field  $\mathbb{F}$  can be assumed as  $\mathbb{R}$  or  $\mathbb{C}$ . We consider both real and complex cases simultaneously. For this purpose, we denote a constant  $\kappa = 1/2$  or  $1$  for the real or complex case, respectively.

Suppose  $S$  is an index set. We denote by  $|S|$  the number of elements in  $S$ . For any column vector  $x \in \mathbb{F}^N$ ,  $x^S \in \mathbb{F}^{|S|}$  is the vector in  $\mathbb{F}^{|S|}$  formed by the components of  $x$  indicated by the index set  $S$ ; for any matrix  $B$ ,  $B^S$  denotes the submatrix formed by picking the rows of  $B$  corresponding to indices in  $S$ , while  $B_S$  is the submatrix with columns from  $B$  indicated by  $S$ . If  $I$  and  $J$  are two index sets, then  $B_J^I = (B^I)_J$ , the submatrix of  $B$  with rows indicated by  $I$  and columns indicated by  $J$ .

Transpose of a vector or matrix is denoted by  $'$  while conjugate transpose by  $^\dagger$ .  $A \otimes B$  represents the Kronecker product of two matrices. For a vector  $v$ ,  $\text{diag}(v)$  is the diagonal matrix with the elements of  $v$  in the diagonal. The identity matrix of dimension  $M$  is  $I_M$ . The trace of matrix  $A$  is given by  $\text{tr}(A)$ , the determinant by  $|A|$ , and the rank by  $\text{rank}(A)$ . Though the notation for determinant is inconsistent with that for cardinality for an index set, the exact meaning can always be understood from the context.

Bold symbols are reserved for *random* vectors and matrices. We use  $\mathbb{P}$  to denote the probability of an event and  $\mathbb{E}$  the expectation. The underlying probability space can be inferred from the context. Gaussian distribution for a random vector in field  $\mathbb{F}$  with mean  $\mu$  and covariance matrix  $\Sigma$  is represented by  $\mathbb{FN}(\mu, \Sigma)$ . Matrix variate Gaussian distribution [42] for  $\mathbf{Y} \in \mathbb{F}^{M \times T}$  with mean  $\Theta \in \mathbb{F}^{M \times T}$  and covariance matrix  $\Sigma \otimes \Psi$ , where  $\Sigma \in \mathbb{F}^{M \times M}$  and  $\Psi \in \mathbb{F}^{T \times T}$ , is denoted by  $\mathbb{FN}_{M,T}(\Theta, \Sigma \otimes \Psi)$

Suppose  $\{f_n\}_{n=1}^\infty, \{g_n\}_{n=1}^\infty$  are two positive sequences,  $f_n = o(g_n)$  means that  $\lim_{n \rightarrow \infty} \frac{f_n}{g_n} = 0$ . An alternative notation in this case is  $g_n \gg f_n$ . We use  $f_n = O(g_n)$  to denote that there exists an  $N \in \mathbb{N}$  and  $C > 0$  independent of  $N$  such that  $f_n \leq Cg_n$  for  $n \geq N$ . Similarly,  $f_n = \Omega(g_n)$  means  $f_n \geq Cg_n$  for  $n \geq N$ . These simple but expedient notations introduced by G. H. Hardy greatly simplify derivations [43].

## B. Models

Next, we introduce our mathematical model. Suppose  $\mathbf{x}(t) \in \mathbb{F}^N, t = 1, \dots, T$  are jointly sparse signals with common support; that is, only a few components of  $\mathbf{x}(t)$  are non-zero and the indices corresponding to these non-zero components are the same for all  $t = 1, \dots, T$ . The common support  $S = \text{supp}(\mathbf{x}(t))$  has known size  $K = |S|$ . We assume that the vectors  $\mathbf{x}^S(t), t = 1, \dots, T$  formed by the non-zero components of  $\mathbf{x}(t)$  follow *i.i.d.*  $\mathbb{FN}(0, I_K)$ . The measurement model is as follows:

$$\mathbf{y}(t) = A\mathbf{x}(t) + \mathbf{w}(t), t = 1, 2, \dots, T, \quad (2)$$

where  $A$  is the measurement matrix and  $\mathbf{y}(t) \in \mathbb{F}^M$  the measurements. The additive noises  $\mathbf{w}(t) \in \mathbb{F}^N$  are assumed to follow *i.i.d.*  $\mathbb{FN}(0, \sigma^2 I_M)$ . Note that assuming unit variance for signals loses no generality since only the ratio of signal variance to noise variance appears in all subsequence analyses. In this sense, we view  $1/\sigma^2$  as the signal-to-noise ratio (SNR).

Let  $\mathbf{X} = \begin{bmatrix} \mathbf{x}(1) & \mathbf{x}(2) & \cdots & \mathbf{x}(T) \end{bmatrix}$  and  $\mathbf{Y}$ ,  $\mathbf{W}$  be defined in a similar manner. Then we write the model in the more compact matrix form:

$$\mathbf{Y} = \mathbf{A}\mathbf{X} + \mathbf{W}. \quad (3)$$

We start our analysis for general measurement matrix  $A$ . For an arbitrary measurement matrix  $A \in \mathbb{F}^{M \times N}$ , if every  $M \times M$  submatrix of  $A$  is non-singular, we then call  $A$  a *non-degenerate* measurement matrix. In this case, the corresponding linear system  $Ax = b$  is said to have the *Unique Representation Property (URP)*, the implication of which is discussed in [12]. While most of our results apply to general non-degenerate measurement matrices, we need to impose more structure on the measurement matrices in order to obtain more profound results. In particular, we will consider Gaussian measurement matrix  $\mathbf{A}$  whose elements  $A_{mn}$  are generated from *i.i.d.*  $\mathbb{FN}(0, 1)$ . However, since our performance analysis is carried out by conditioning on a particular realization of  $\mathbf{A}$ , we still use non-bold  $A$  except in Section V. The role played by the variance of  $A_{mn}$  is indistinguishable from that of a signal variance and hence can be combined to  $1/\sigma^2$ , the SNR, by the note in the previous paragraph.

We now consider two hypothesis-testing problems. The first one is a binary support recovery problem:

$$\begin{cases} H_0 : \text{supp}(\mathbf{X}) = S_0 \\ H_1 : \text{supp}(\mathbf{X}) = S_1 \end{cases}. \quad (4)$$

The results we obtain for binary support recovery (4) offer insight into our second problem: the multiple support recovery. In the multiple support recovery problem we choose one among  $C_N^K$  distinct candidate supports of  $\mathbf{X}$ , which is a multiple-hypothesis testing problem:

$$\begin{cases} H_0 : \text{supp}(\mathbf{X}) = S_0 \\ H_1 : \text{supp}(\mathbf{X}) = S_1 \\ \vdots \\ H_{L-1} : \text{supp}(\mathbf{X}) = S_{L-1} \end{cases}. \quad (5)$$

### C. Preliminaries for Hypothesis Testing

We now briefly introduce the fundamentals of hypothesis testing. The following discussion is based mainly on [41]. In a simple binary hypothesis test, the goal is to determine which of two

candidate distributions is the true one that generates the data matrix (or vector)  $\mathbf{Y}$ :

$$\begin{cases} H_0 : \mathbf{Y} \sim p(\mathbf{Y}|H_0) \\ H_1 : \mathbf{Y} \sim p(\mathbf{Y}|H_1) \end{cases}. \quad (6)$$

There are two types of errors when one makes a choice based on the observed data  $\mathbf{Y}$ . A *false alarm* corresponds to choosing  $H_1$  when  $H_0$  is true, while a *miss* happens by choosing  $H_0$  when  $H_1$  is true. The probabilities of these two types of errors are called the probability of a false alarm and the probability of a miss, which are denoted by

$$P_F = \mathbb{P}(\text{Choose } H_1 | H_0), \quad (7)$$

$$P_M = \mathbb{P}(\text{Choose } H_0 | H_1), \quad (8)$$

respectively. Depending on whether one knows the prior probabilities  $\mathbb{P}(H_0)$  and  $\mathbb{P}(H_1)$  and assigns losses to errors, different criteria can be employed to derive the optimal decision rule. In this paper we adopt the probability of error with equal prior probabilities of  $H_0$  and  $H_1$  as the decision criterion; that is, we try to find the optimal decision rule by minimizing

$$P_{\text{err}} = P_F \mathbb{P}(H_0) + P_M \mathbb{P}(H_1) = \frac{1}{2} P_F + \frac{1}{2} P_D.$$

The optimal decision rule is then given by the *likelihood ratio test*:

$$\ell(\mathbf{Y}) = \ln \frac{p(\mathbf{Y}|H_1)}{p(\mathbf{Y}|H_0)} \underset{H_0}{\overset{H_1}{\gtrless}} 0. \quad (9)$$

The probability of error for the optimal decision rule, namely, the likelihood ratio test (9), is an indication of the underlying system's maximal possible performance. In many cases of interest, the simple binary hypothesis testing problem (6) is derived from a signal-generation system. For example, in a digital communication system, hypotheses  $H_0$  and  $H_1$  correspond to the transmitter sending digit 0 and 1, respectively, and the distributions of the observed data under the hypotheses are determined by the modulation method of the system. Therefore, the minimal probability of error achieved by the likelihood ratio test is a measure of the performance of the modulation method. For the problem addressed in this paper, the minimal probability of error reflects the measurement matrix's ability to distinguish different signal supports.

The Chernoff bound [41] is a well-known tight upper bound on the probability of error. In many cases, the optimum test can be derived and implemented efficiently. However, it is usually extremely difficult to calculate an exact expression for the performance. Even if such an



expression can be derived, it is too complicated to be of practical use. For that reason, sometimes a simple bound turns out to be more useful in many problems of practical importance. The Chernoff bound, based on the moment generating function of the test statistic  $\ell(\mathbf{Y})$  (9), is one of the most famous ones, blending efficiency and simplicity.

Define  $\mu(s)$  as the logarithm of the moment generating function of  $\ell(Y)$ ; *i.e.*,

$$\begin{aligned}\mu(s) &\triangleq \ln \int_{-\infty}^{\infty} e^{s\ell(\mathbf{Y})} p(\mathbf{Y}|\mathbf{H}_0) d\mathbf{Y} \\ &= \ln \int_{-\infty}^{\infty} [p(\mathbf{Y}|\mathbf{H}_1)]^s [p(\mathbf{Y}|\mathbf{H}_0)]^{1-s} d\mathbf{Y}.\end{aligned}\quad (10)$$

Then the Chernoff bound states that

$$P_F \leq \exp[\mu(s_m)] \leq \exp[\mu(s)], \quad (11)$$

$$P_M \leq \exp[\mu(s_m)] \leq \exp[\mu(s)], \quad (12)$$

and

$$P_{\text{err}} \leq \frac{1}{2} \exp[\mu(s_m)] \leq \frac{1}{2} \exp[\mu(s)], \quad (13)$$

where  $0 \leq s \leq 1$  and  $s_m = \operatorname{argmin}_{0 \leq s \leq 1} \mu(s)$ . We use these bounds to study the performance of the support recovery problem.

We next extend to multiple-hypothesis testing the key elements of the binary hypothesis testing. The goal in a simple multiple-hypothesis testing problem is to make a choice among  $L$  distributions based on the observations:

$$\left\{ \begin{array}{ll} \mathbf{H}_0 : & \mathbf{Y} \sim p(\mathbf{Y}|\mathbf{H}_0) \\ \mathbf{H}_1 : & \mathbf{Y} \sim p(\mathbf{Y}|\mathbf{H}_1) \\ & \vdots \\ \mathbf{H}_{L-1} : & \mathbf{Y} \sim p(\mathbf{Y}|\mathbf{H}_{L-1}) \end{array} \right. . \quad (14)$$

Using the total probability of error as a decision criterion and assuming equal prior probabilities for all hypotheses, we obtain the optimal decision rule given by

$$\mathbf{H}^* = \operatorname{argmax}_{0 \leq i \leq L-1} p(\mathbf{Y}|\mathbf{H}_i). \quad (15)$$

The total probability of error is defined as

$$\begin{aligned}P_{\text{err}} &= \sum_{i=0}^{L-1} \mathbb{P}(\mathbf{H}^* \neq \mathbf{H}_i | \mathbf{H}_i) \mathbb{P}(\mathbf{H}_i) \\ &= \frac{1}{L} \sum_{i=0}^{L-1} \mathbb{P}(\mathbf{H}^* \neq \mathbf{H}_i | \mathbf{H}_i).\end{aligned}\quad (16)$$

The union bound implies that

$$\begin{aligned}\mathbb{P}(H^* \neq H_i | H_i) &= \mathbb{P}\left\{\bigcup_{j \neq i} [p(\mathbf{Y} | H_j) > p(\mathbf{Y} | H_i)] \middle| H_i\right\} \\ &\leq \sum_{j \neq i} \mathbb{P}\{p(\mathbf{Y} | H_j) > p(\mathbf{Y} | H_i) | H_i\}.\end{aligned}$$

Therefore, we obtain an upper bound on  $P_{\text{err}}$  in multiple-hypothesis testing case by the Chernoff bound (13) for binary hypothesis testing:

$$\begin{aligned}P_{\text{err}} &\leq \frac{1}{2L} \sum_{i \neq j} \left\{ \mathbb{P}\left(\ln \frac{p(\mathbf{Y} | H_j)}{p(\mathbf{Y} | H_i)} > 0 \middle| H_i\right) \right. \\ &\quad \left. + \mathbb{P}\left(\ln \frac{p(\mathbf{Y} | H_j)}{p(\mathbf{Y} | H_i)} < 0 \middle| H_j\right) \right\} \\ &= \frac{1}{2L} \sum_{i=0}^{L-1} \sum_{\substack{j=0 \\ j \neq i}}^{L-1} 2P_{\text{err}}(H_i, H_j) \\ &\leq \frac{1}{2L} \sum_{i=0}^{L-1} \sum_{\substack{j=0 \\ j \neq i}}^{L-1} \exp[\mu(s; H_i, H_j)], \quad 0 \leq s \leq 1, \tag{17}\end{aligned}$$

where  $P_{\text{err}}(H_i, H_j)$  is the probability of error in the binary hypothesis testing for  $H_i$  and  $H_j$ , and  $\exp[\mu(s; H_i, H_j)]$  is the corresponding moment-generating function. Hence, we obtain an upper bound for multiple-hypothesis testing from that for binary hypothesis testing.

### III. UPPER BOUND ON PROBABILITY OF ERROR FOR NON-DEGENERATE MEASUREMENT MATRICES

In this section, we apply the general theory for hypothesis testing, the Chernoff bound on the probability of error in particular, to the support recovery problems (4) and (5). We first study binary support recovery, which lays the foundation for the general support recovery problem. The result is based on a theorem that counts the eigenvalues of a particular matrix defined by the columns of the measurement matrix corresponding to the two candidate supports. We then obtain a bound for the general support recovery problem via (17).

### A. Binary Support Recovery

Under model (3) and the assumptions pertaining to it, observations  $\mathbf{Y}$  follow a matrix variate Gaussian distribution [42] when the true support is  $S$ :

$$\mathbf{Y}|S \sim \mathbb{FN}_{M,T}(0, \Sigma_S \otimes \mathbf{I}_T), \quad (18)$$

with the probability density function (pdf) given by

$$p(\mathbf{Y}|S) = \frac{1}{(\pi/\kappa)^{\kappa MT} |\Sigma_S|^{\kappa T}} \exp \left[ -\kappa \text{tr} (\mathbf{Y}^\dagger \Sigma_S^{-1} \mathbf{Y}) \right], \quad (19)$$

where  $\Sigma_S = A_S A_S^\dagger + \sigma^2 \mathbf{I}_M$  is the common covariance matrix for each column of  $\mathbf{Y}$ . The binary support recovery problem (4) is equivalent to a linear Gaussian binary hypothesis testing problem:

$$\begin{cases} H_0 : \mathbf{Y} \sim \mathbb{FN}_{M,T}(0, \Sigma_{S_0} \otimes \mathbf{I}_T) \\ H_1 : \mathbf{Y} \sim \mathbb{FN}_{M,T}(0, \Sigma_{S_1} \otimes \mathbf{I}_T) \end{cases}. \quad (20)$$

From now on, for notation simplicity we will denote  $\Sigma_{S_i}$  by  $\Sigma_i$ . The optimal decision rule with minimal probability of error given by the likelihood ratio test  $\ell(\mathbf{Y})$  (9) reduces to

$$-\kappa \text{tr} [\mathbf{Y}^\dagger (\Sigma_1^{-1} - \Sigma_0^{-1}) \mathbf{Y}] - \kappa T \ln \frac{|\Sigma_1|}{|\Sigma_0|} \underset{H_0}{\overset{H_1}{\geq}} 0. \quad (21)$$

To analyze the performance of the likelihood ratio test (21), we first compute the log moment generating function of  $\ell(\mathbf{Y})$  according to (10):

$$\begin{aligned} \mu(s) &= \ln \int [p(\mathbf{Y}|H_1)]^s [p(\mathbf{Y}|H_0)]^{1-s} d\mathbf{Y} \\ &= \ln \left[ \frac{1}{(\pi/\kappa)^{\kappa MT} |\Sigma_1|^{\kappa s T} |\Sigma_0|^{\kappa(1-s)T}} \right. \\ &\quad \times \left. \int \exp \left\{ -\kappa \text{tr} [\mathbf{Y}^\dagger (s\Sigma_1^{-1} + (1-s)\Sigma_0^{-1}) \mathbf{Y}] \right\} d\mathbf{Y} \right] \\ &= \ln \left[ \frac{|s\Sigma_1^{-1} + (1-s)\Sigma_0^{-1}|^{-\kappa T}}{|\Sigma_1|^{\kappa s T} |\Sigma_0|^{\kappa(1-s)T}} \right] \\ &= -\kappa T \ln |sH^{1-s} + (1-s)H^{-s}|, \quad 0 \leq s \leq 1, \end{aligned} \quad (22)$$

where  $H = \Sigma_0^{1/2} \Sigma_1^{-1} \Sigma_0^{1/2}$ . The computation of the exact minimizer  $s_m = \arg\min_{0 \leq s \leq 1} \mu(s)$  is non-trivial and will lead to an expression of  $\mu(s_m)$  too complicated to handle. When  $|S_0| = |S_1|$

and the columns of  $A$  are not highly correlated, for example in the case of  $A$  with *i.i.d.* elements,  $s_m \approx \frac{1}{2}$ . We then take  $s = \frac{1}{2}$  in the Chernoff bounds (11), (12), and (13). Note that this is not an approximation and the bound we get is exact.

As positive definite Hermitian matrices,  $H$  and  $H^{-1}$  can be simultaneously diagonalized by an unitary matrix. Suppose that the eigenvalues of  $H$  are  $\lambda_1 \geq \dots \geq \lambda_{k_0} > 1 = \dots = 1 > \sigma_1 \geq \dots \geq \sigma_{k_1}$  and  $D = \text{diag}[\lambda_1, \dots, \lambda_{k_0}, 1, \dots, 1, \sigma_1, \dots, \sigma_{k_1}]$ . Then it is easy to show that

$$\begin{aligned} \mu(1/2) &= -\kappa T \ln \left| \frac{D^{1/2} + D^{-1/2}}{2} \right| \\ &= -\kappa T \left[ \sum_{j=1}^{k_0} \ln \left( \frac{\sqrt{\lambda_j} + 1/\sqrt{\lambda_j}}{2} \right) \right. \\ &\quad \left. + \sum_{j=1}^{k_1} \ln \left( \frac{\sqrt{\sigma_j} + 1/\sqrt{\sigma_j}}{2} \right) \right]. \end{aligned} \quad (24)$$

Therefore, it is necessary to count the numbers of eigenvalues of  $H$  that are greater than 1, equal to 1 and less than 1, *i.e.*, the values of  $k_0$  and  $k_1$  for general non-degenerate measurement matrix  $A$ . We have the following theorem on the eigenvalue structure of  $H$ :

**Proposition 1** *For any non-degenerate measurement matrix  $A$ , let  $H = \Sigma_0^{1/2} \Sigma_1^{-1} \Sigma_0^{1/2}$ ,  $k_i = |S_0 \cap S_1|$ ,  $k_0 = |S_0 \setminus S_1| = |S_0| - k_i$ ,  $k_1 = |S_1 \setminus S_0| = |S_1| - k_i$  and assume  $M \geq k_0 + k_1$ ; then  $k_0$  eigenvalues of matrix  $H$  are greater than 1,  $k_1$  less than 1, and  $M - (k_0 + k_1)$  equal to 1.*

**Proof:** See Appendix A.

For binary support recovery (4) with  $|S_0| = |S_1| = K$ , we have  $k_0 = k_1 \triangleq k_d$ . Employing the Chernoff bounds (13) and Proposition 1, we have

**Proposition 2** *The probability of error for the binary support recovery problem (4) is bounded by*

$$P_{\text{err}} \leq \frac{1}{2} \left[ \frac{\bar{\lambda}_{S_0, S_1} \bar{\lambda}_{S_1, S_0}}{16} \right]^{-\kappa k_d T / 2}, \quad (25)$$

where  $\bar{\lambda}_{S_i, S_j}$  is the geometric mean of the eigenvalues of  $H = \Sigma_i^{1/2} \Sigma_j^{-1} \Sigma_i^{1/2}$  that are greater than one.

**Proof:** According to (13) and (24), we have

$$\begin{aligned}
P_{\text{err}} &\leq \frac{1}{2} \exp \left[ \mu \left( \frac{1}{2} \right) \right] \\
&\leq \frac{1}{2} \left[ \prod_{j=1}^{k_d} \left( \frac{\sqrt{\lambda_j}}{2} \right) \prod_{j=1}^{k_d} \left( \frac{1/\sqrt{\sigma_j}}{2} \right) \right]^{-\kappa T} \\
&= \frac{1}{2} \left[ \frac{\left( \prod_{j=1}^{k_d} \lambda_j \right)^{1/k_d} \left( \prod_{j=1}^{k_d} \frac{1}{\sigma_j} \right)^{1/k_d}}{16} \right]^{-\kappa k_d T/2}.
\end{aligned}$$

Define  $\bar{\lambda}_{S_i, S_j}$  as the geometric mean of the eigenvalues of  $H = \Sigma_i^{1/2} \Sigma_j^{-1} \Sigma_i^{1/2}$  that are greater than one. Then obviously we have  $\bar{\lambda}_{S_0, S_1} = \left( \prod_{j=1}^{k_d} \lambda_j \right)^{1/k_d}$ . Since  $H^{-1}$  and  $\Sigma_1^{1/2} \Sigma_0^{-1} \Sigma_1^{1/2}$  have the same set of eigenvalues,  $1/\sigma_j, j = 1, \dots, k_d$  are the eigenvalues of  $\Sigma_1^{1/2} \Sigma_0^{-1} \Sigma_1^{1/2}$  that are greater than 1. We conclude that  $\bar{\lambda}_{S_0, S_1} = \left( \prod_{j=1}^{k_d} 1/\sigma_j \right)^{1/k_d}$ . ■

Note that  $\bar{\lambda}_{S_0, S_1} \bar{\lambda}_{S_1, S_0}$  completely determines the measurement system (3)'s performance in differentiating two different signal supports. It must be larger than the constant 16 for a vanishing bound when more temporal samples are taken. Once the threshold 16 is exceeded, taking more samples will drive the probability of error to 0 exponentially fast. From numerical simulations and our results on the Gaussian measurement matrix,  $\bar{\lambda}_{S_i, S_j}$  does not vary much when  $k_d$  changes, as long as the elements in the measurement matrix  $A$  are highly uncorrelated. Therefore, quite appealing to intuition, the larger the size  $k_d$  of the difference set between the two candidate supports, the smaller the probability of error.

### B. Multiple Support Recovery

Now we are ready to use the union bound (17) to study the probability of error for the multiple support recovery problem (5). We assume each candidate support  $S_i$  has known cardinality  $K$ , and we have  $L = \binom{N}{K}$  such supports. Our general approach is also applicable to cases for which we have some prior information on the structure of the signal's sparsity pattern, for example the setup in model-based compressive sensing [44]. In these cases, we usually have  $L \ll \binom{N}{K}$  supports, and a careful examination on the intersection pattern of these supports will give a better bound. However, in this paper we will not address this problem and will instead focus on the full support recovery problem with  $L = \binom{N}{K}$ . Defining  $\bar{\lambda} = \min_{i \neq j} \{\bar{\lambda}_{S_i, S_j}\}$ , we have the following theorem:

**Theorem 1** *If  $\bar{\lambda} > 4 [K(N-K)]^{\frac{1}{\kappa T}}$ , then the probability of error for the full support recovery problem (5) with  $|S_i| = K$  and  $L = \binom{N}{K}$  is bounded by*

$$P_{\text{err}} \leq \frac{1}{2} \frac{\frac{K(N-K)}{(\bar{\lambda}/4)^{\kappa T}}}{1 - \frac{K(N-K)}{(\bar{\lambda}/4)^{\kappa T}}}. \quad (26)$$

**Proof:** Combining the bound in Proposition 2 and Equation (17), we have

$$\begin{aligned} P_{\text{err}} &\leq \frac{1}{2L} \sum_{i=0}^{L-1} \sum_{\substack{j=1 \\ j \neq i}}^{L-1} \left[ \frac{\bar{\lambda}_{S_i, S_j} \bar{\lambda}_{S_j, S_i}}{16} \right]^{-\kappa k_d T / 2} \\ &\leq \frac{1}{2L} \sum_{i=0}^{L-1} \sum_{\substack{j=1 \\ j \neq i}}^{L-1} \left( \frac{\bar{\lambda}}{4} \right)^{-\kappa k_d T}. \end{aligned}$$

Here  $k_d$  depends on the supports  $S_i$  and  $S_j$ . For fixed  $S_i$ , the number of supports that have a difference set with  $S_i$  with cardinality  $k_d$  is  $\binom{K}{k_d} \binom{N-K}{k_d}$ . Therefore, using  $\binom{K}{k_d} \leq K^{k_d}$  and  $\binom{N-K}{k_d} \leq (N-K)^{k_d}$  and the summation formula for geometric series, we obtain

$$\begin{aligned} P_{\text{err}} &\leq \frac{1}{2L} \sum_{i=0}^{L-1} \sum_{k_d=1}^K \binom{K}{k_d} \binom{N-K}{k_d} \left( \frac{\bar{\lambda}}{4} \right)^{-\kappa k_d T} \\ &\leq \frac{1}{2} \sum_{k_d=1}^K \left[ \frac{K(N-K)}{(\bar{\lambda}/4)^{\kappa T}} \right]^{k_d} \\ &\leq \frac{1}{2} \frac{\frac{K(N-K)}{(\bar{\lambda}/4)^{\kappa T}}}{1 - \frac{K(N-K)}{(\bar{\lambda}/4)^{\kappa T}}}. \quad \blacksquare \end{aligned}$$

We make several comments here. First,  $\bar{\lambda}$  depends solely on the measurement matrix  $A$ . Compared with the results in [45], where the bounds involve the signal, we get more insight into what quantity of the measurement matrix is important in support recovery. This information is obtained by modelling the signals  $\mathbf{x}(t)$  as Gaussian random vectors. The quantity  $\bar{\lambda}$  effectively characterizes system (3)'s ability to distinguish different supports. In some sense,  $\bar{\lambda}$  plays the same role as the RIP constant [3]–[5] in classical compressive sensing. As we mentioned before, once we get the support, the corresponding component values can be reconstructed by solving a least square problem. Therefore,  $\bar{\lambda}$  should be closely related to the RIP condition.

Second, we observe that increasing the number of temporal samples plays two roles simultaneously in the measurement system. For one thing, it decreases the threshold  $4[K(N-K)]^{\frac{1}{\kappa T}}$

that  $\bar{\lambda}$  must exceed for a vanishing probability of error. However, since  $\lim_{T \rightarrow \infty} 4[K(N-K)]^{\frac{1}{\kappa T}} = 4$  for fixed  $K$  and  $N$ , increasing time samples can reduce the threshold only to a certain limit. For another, once the threshold is exceeded, the probability of error turns to 0 exponentially fast as  $T$  increases.

In addition, the final bound (26) is of the same order as the probability of error when  $k_d = 1$ . The probability of error  $P_{\text{err}}$  is dominated by the probability of error in cases for which the estimated support differs by only one index from the true support, which are the most difficult cases for the decision rule to make a choice. However, in practice we can imagine that these cases induce the least loss. Therefore, if we assign weights/costs to the errors based on  $k_d$ , then the weighted probability of error or average cost would be much lower. For example, we can choose the costs to exponentially decrease when  $k_d$  increases. Another possible choice of cost function is to assume zero cost when  $k_d$  is below a certain critical number. Our results can be easily extended to these scenarios.

Finally, note that our bound (26) applies to any non-degenerate matrix. In Section V, we apply the bound to Gaussian measurement matrices. The additional structure allows us to derive more profound results on the behavior of the bound.

### C. The Effect of Noise

In this subsection, we explore how the noise variance affects the probability of error, which is equivalent to analyzing the behavior of  $\bar{\lambda}_{S_i, S_j}$  and  $\bar{\lambda}$  as indicated in (25) and (26).

We now derive bounds on the eigenvalues of  $H$  that are greater than 1. The lower bound is expressed in terms of the QR decomposition of a submatrix of the measurement matrix with the noise variance  $\sigma^2$  isolated.

**Proposition 3** *For any non-degenerate measurement matrix  $A$ , let  $H = \Sigma_0^{1/2} \Sigma_1^{-1} \Sigma_0^{1/2}$  with  $\Sigma_i = A_{S_i} A_{S_i}^\dagger + \sigma^2 \mathbf{I}_M$ ,  $k_i = |S_0 \cap S_1|$ ,  $k_0 = |S_0 \setminus S_1| = |S_0| - k_i$ ,  $k_1 = |S_1 \setminus S_0| = |S_1| - k_i$  and assume  $M \geq k_0 + k_1$ . Then the sorted eigenvalues of  $H$  that are greater than 1 are*

- 1) *lower bounded by the corresponding eigenvalues of  $\mathbf{I}_{k_0} + \frac{1}{\sigma^2} R_{33} R_{33}^\dagger$ , where  $R_{33}$  is the  $k_0 \times k_0$  submatrix at the lower-right corner of the upper triangle matrix in the QR decomposition of  $\begin{bmatrix} A_{S_1 \setminus S_0} & A_{S_1 S_0} & A_{S_0 \setminus S_1} \end{bmatrix}$*
- 2) *upper bounded by the corresponding eigenvalues of  $\mathbf{I}_{k_0} + \frac{1}{\sigma^2} A_{S_0 \setminus S_1} A_{S_0 \setminus S_1}^\dagger$*

**Proof:** See Appendix B.

The importance of this proposition is twofold. First, by isolating the noise variance from the expression of matrix  $H$ , this theorem clearly shows that when noise variance decreases to zero, the relatively large eigenvalues of  $H$  will blow up, which results in increased performance in support recovery. Second, the bounds provide ways to analyze special measurement matrices, especially the Gaussian measurement ensemble discussed in the next section.

We have the following corollary:

**Corollary 1** *For support recovery problems (4) and (5) with support size  $K$ , suppose  $M \geq 2K$ ; then there exist constants  $c_1, c_2 > 0$  that depend only on the measurement matrix  $A$  such that*

$$1 + \frac{c_2}{\sigma^2} \geq \bar{\lambda} \geq 1 + \frac{c_1}{\sigma^2}. \quad (27)$$

*This corollary implies from (25) and (26) that*

$$\lim_{\sigma^2 \rightarrow 0} P_{\text{err}} = 0 \quad (28)$$

*and the speed of convergence is approximately  $(\sigma^2)^{\kappa k_d T}$  and  $(\sigma^2)^{\kappa T}$  for the binary and multiple cases, respectively.*

**Proof:** According to Proposition 3, for any fixed  $S_i, S_j$ , the eigenvalues of  $H = \Sigma_i^{1/2} \Sigma_j^{-1} \Sigma_i^{1/2}$  that are greater than 1 are lower bounded by those of  $\mathbf{I}_{k_d} + \frac{1}{\sigma^2} R_{33} R_{33}^\dagger$ ; hence we have

$$\begin{aligned} \bar{\lambda}_{S_i, S_j} &\geq \left| \mathbf{I}_{k_d} + \frac{1}{\sigma^2} R_{33} R_{33}^\dagger \right|^{1/k_d} \\ &\geq \left| \mathbf{I}_{k_d} + \frac{1}{\sigma^2} R_{33} R_{33}^\dagger \right|^{1/k_d} \\ &= \left[ \prod_{l=1}^{k_d} \left( 1 + \frac{1}{\sigma^2} r_{ll}^2 \right) \right]^{1/k_d} \\ &\geq 1 + \frac{1}{\sigma^2} \left( \prod_{l=1}^{k_d} r_{ll}^2 \right)^{1/k_d}, \end{aligned} \quad (29)$$

where  $r_{ll}$  is the  $l$ th diagonal element of  $R_{33}$ . For the second inequality we have used Fact 8. 11. 20 in [46]. Since  $A$  is non-degenerate,  $\begin{bmatrix} A_{S_j \setminus S_i} & A_{S_j S_i} & A_{S_i \setminus S_j} \end{bmatrix}$  is full rank and  $r_{ll}^2 > 0, 0 \leq l \leq$



$k_d$  for all  $S_i, S_j$ . Defining  $c_1$  as the minimal value of  $\left(\prod_{l=1}^{k_d} r_{ll}^2\right)^{1/k_d}$ 's over all possible support pairs  $S_i, S_j$ , we then have  $c_1 > 0$  and

$$\bar{\lambda} \geq 1 + \frac{c_1}{\sigma^2}.$$

On the other hand, the upper bound on the eigenvalues of  $H$  yields

$$\begin{aligned} \bar{\lambda}_{S_i, S_j} &\leq \left| \mathbf{I}_{k_d} + \frac{1}{\sigma^2} A_{S_i \setminus S_j} A_{S_i \setminus S_j}^\dagger \right|^{1/k} \\ &\leq 1 + \frac{1}{\sigma^2} \text{tr} \left( A_{S_i \setminus S_j} A_{S_i \setminus S_j}^\dagger \right) \\ &\leq 1 + \frac{1}{\sigma^2 k_d} \sum_{\substack{1 \leq m \leq M \\ n \in S_i \setminus S_j}} |A_{mn}|^2. \end{aligned} \quad (30)$$

Therefore, we have

$$\bar{\lambda} \leq 1 + \frac{c_2}{\sigma^2},$$

with  $c_2 = \max_{S: |S|=K} \sum_{\substack{1 \leq m \leq M \\ n \in S}} A_{mn}^2$ . All other statements in the theorem follows immediately from (25) and (26). ■

Corollary 1 suggests that in the limiting case where there is no noise,  $M \geq 2K$  is sufficient to recover a  $K$ -sparse signal. This fact has been observed in [3]. Our results also shows that the optimal decision rule, which is unfortunately inefficient, is robust to noise. Another extreme case is one for which the noise variance  $\sigma^2$  is very large. Then from  $\ln(1+x) \approx x, 0 < x \ll 1$ , the bounds in (25) and (26) are approximated by  $e^{-\kappa k_d T / \sigma^2}$  and  $e^{-\kappa T / \sigma^2}$ . Therefore, the performance improves only slightly by increasing the number of temporal samples.

The diagonal elements of  $R_{33}$ ,  $r_{ll}$ 's, have clear meanings. Since QR factorization is equivalent to the Gram-Schmidt orthogonalization procedure,  $r_{11}$  is the distance of the first column of  $A_{S_i/S_j}$  to the subspace spanned by the columns of  $A_{S_j}$ ;  $r_{22}$  is the distance of the second column of  $A_{S_i/S_j}$  to the subspace spanned by the columns of  $A_{S_j}$  plus the first column of  $A_{S_i/S_j}$ , and so on. Therefore,  $\bar{\lambda}_{S_i, S_j}$  is a measure of how well the columns of  $A_{S_i/S_j}$  can be expressed by the columns of  $A_{S_j}$ , or, put another way, a measure of the incoherence between the columns of  $A_{S_i}$  and  $A_{S_j}$ . Similarly,  $\bar{\lambda}$  is an indicator of the incoherence of the entire matrix  $A$  of order  $K$ .

To relate  $\bar{\lambda}$  with the incoherence, we consider the case for  $K = 1$  and  $\mathbb{F} = \mathbb{R}$ . By restricting our attention to matrices with unit columns, the above discussion implies that a better bound is achieved if the minimal distance of all pairs of column vectors of matrix  $A$  is maximized.

Finding such a matrix  $A$  is equivalent to finding a matrix with the inner product between columns as large as possible, since the distance between two unit vectors  $u$  and  $v$  is  $2 - 2| \langle u, v \rangle |$  where  $\langle u, v \rangle = u'v$  is the inner product between  $u$  and  $v$ . The restricted isometry propriety (RIP) [3], [4] constant  $\delta_{2K}$  is defined by

$$1 - \delta_{2K} \leq \frac{\|Ax\|_2^2}{\|x\|_2^2} \leq 1 + \delta_{2K}, \quad |\text{supp}(x)| = 2K. \quad (31)$$

A direct computation shows that the minimal value of  $\delta_2$  is the minimum among the absolute values of inner products between all pairs of columns of  $A$ . Hence, the requirements imposed on matrix  $A$  by minimizing the RIP constant  $\delta_2$  and maximizing our  $\bar{\lambda}$  coincide.

#### IV. AN INFORMATION THEORETICAL LOWER BOUND ON PROBABILITY OF ERROR

In this section, we derive an information theoretical lower bound on the probability of error for *any* decision rule in the multiple support recovery problem. The main tool is a variant of the well-known Fano's inequality [47]. In the variant, the average probability of error in a multiple-hypothesis testing problem is bounded in terms of the Kullback-Leibler divergence [48]. Suppose that we have a random vector  $\mathbf{Y}$  with  $L$  possible densities  $f_0, \dots, f_{L-1}$ . Denote the average of the Kullback-Leibler divergence between any pair of densities by

$$\beta = \frac{1}{L^2} \sum_{i,j} D_{KL}(f_i || f_j). \quad (32)$$

Then by Fano's inequality [49], [45], the probability of error (16) for *any* decision rule to identify the true density is lower bounded by

$$P_{\text{err}} \geq 1 - \frac{\beta + \ln 2}{\ln L}. \quad (33)$$

Since in the multiple support recovery problem (5), all the distributions involved are matrix variate Gaussian distributions with mean 0 and different variance, we now compute the Kullback-Leibler divergence between two matrix variate Gaussian distributions. Suppose  $f_i = \mathbb{FN}_{M,T}(0, \Sigma_i \otimes I_T)$ ,  $f_j = \mathbb{FN}_{M,T}(0, \Sigma_j \otimes I_T)$ , then the Kullback-Leibler divergence has closed form expression:

$$\begin{aligned} D_{KL}(f_i || f_j) &= E_{f_i} \ln \frac{f_i}{f_j} \\ &= \frac{1}{2} E_{f_i} \left[ -\kappa \text{tr} [\mathbf{Y}^\dagger (\Sigma_i^{-1} - \Sigma_j^{-1}) \mathbf{Y}] - \kappa T \ln \frac{|\Sigma_i|}{|\Sigma_j|} \right] \\ &= \frac{1}{2} \kappa T [\text{tr} (H_{i,j} - I_M) + 2 \ln \frac{|\Sigma_j|}{|\Sigma_i|}], \end{aligned}$$

where  $H_{i,j} = \Sigma_i^{1/2} \Sigma_j^{-1} \Sigma_i^{1/2}$ . Therefore, we obtain the average Kullback-Leibler divergence (32) as

$$\begin{aligned} \beta &= \frac{1}{L^2} \sum_{i,j} \frac{1}{2} \kappa T [2 \ln \frac{|\Sigma_j|}{|\Sigma_i|} + \text{tr}(H_{i,j}) - M] \\ &= \frac{\kappa T}{2L^2} \sum_{i,j} [\text{tr}(H_{i,j}) - M], \end{aligned}$$

where the  $\ln \frac{|\Sigma_j|}{|\Sigma_i|}$  terms all cancel out. Invoking Proposition 1, we get

$$\begin{aligned} \text{tr}(H_{i,j}) &= \sum \sigma(H_{i,j}) + \sum \lambda(H_{i,j}) + (M - 2k_d) \\ &\leq \sum [\lambda(H_{i,j}) - 1] + M, \end{aligned}$$

where  $\sigma(H_{i,j})$  and  $\lambda(H_{i,j})$  are the eigenvalues of  $H_{i,j}$  that are less than and greater than 1, respectively. For the last inequality, we drop  $\sigma(H_{i,j}) - 1 < 0$ . Define by  $\hat{\lambda}_{S_i, S_j} = \frac{1}{k_d} \sum \lambda(H_{i,j})$  the arithmetic mean of the eigenvalues of  $H_{i,j}$  that are greater than 1. The average Kullback-Leibler divergence is bounded by

$$\beta \leq \frac{\kappa T}{2L^2} \sum_{i=0}^{L-1} \sum_{k_d=1}^K \binom{K}{K-k_d} \binom{N-K}{k_d} k_d (\hat{\lambda}_{S_i, S_j} - 1).$$

Note that from Proposition 3, we have

$$\hat{\lambda}_{S_i, S_j} \leq \frac{1}{k_d} \text{tr} \left[ \mathbf{I}_M + \frac{1}{\sigma^2} A_{S_i \setminus S_j} A_{S_i \setminus S_j}^\dagger \right] \quad (34)$$

$$= 1 + \frac{1}{\sigma^2 k_d} \sum_{\substack{1 \leq m \leq M \\ n \in S_i \setminus S_j}} |A_{mn}|^2 \quad (35)$$

Therefore, we obtain

$$\beta \leq \frac{\kappa T}{2L^2 \sigma^2} \sum_{i=0}^{L-1} \sum_{k_d=1}^K \binom{K}{K-k_d} \binom{N-K}{k_d} \sum_{\substack{1 \leq m \leq M \\ n \in S_i \setminus S_j}} |A_{mn}|^2.$$

Due to the symmetry of the right-hand side, it must be of the form  $\alpha \sum_{\substack{1 \leq m \leq M \\ 1 \leq n \leq N}} |A_{mn}|^2 = \alpha \|A\|_F^2$ , where  $\|\cdot\|_F$  is the Frobenius norm. Setting all  $A_{mn} = 1$  gives  $\alpha = \frac{\kappa T K (N-K)}{2\sigma^2 N^2}$ , where we need to use the mean expression for hypergeometric distribution:

$$\sum_{k_d=1}^K \frac{\binom{K}{k_d} \binom{N-K}{k_d}}{\binom{N}{K}} k_d = \frac{K(N-K)}{N}.$$

Hence, we get

$$\beta \leq \frac{\kappa T K (N - K)}{2\sigma^2 N^2} \|A\|_F^2.$$

Therefore, the probability of error is lower bounded by

$$P_{\text{err}} \geq 1 - \frac{\frac{\kappa T K (N - K)}{2\sigma^2 N^2} \|A\|_F^2 + \ln 2}{\ln L}. \quad (36)$$

We conclude with the following theorem:

**Theorem 2** *For multiple support recovery problem (5), the probability of error for any decision rule is lower bounded by*

$$P_{\text{err}} \geq 1 - \frac{\kappa T \frac{K}{N} \left(1 - \frac{K}{N}\right) \|A\|_F^2}{2\sigma^2 \ln \binom{N}{K}} + o(1). \quad (37)$$

Each term in bound (37) has clear meanings. The Frobenius norm of measurement matrix  $\|A\|_F^2$  is the maximal possible total gain of system (2). Since the measured signal is  $K$ -sparse, only a fraction of the gain plays a role in the measurement, and its average over all possible  $K$ -sparse signals is  $\frac{K}{N} \left(1 - \frac{K}{N}\right) \|A\|_F^2$ . When  $K$  is small compared with  $N$ , we have that approximately  $\frac{K}{N}$  of  $\|A\|_F^2$  contributes to the measurements. The term  $\ln L = \ln \binom{N}{K}$  is the total uncertainty or entropy of the support variable  $S$ , since we impose a uniform prior on it. As long as  $K \leq \frac{N}{2}$ , increasing  $K$  will increase both the average gain exploited by the measurement system, namely,  $\frac{K}{N} \left(1 - \frac{K}{N}\right) \|A\|_F^2$ , and the entropy of the support variable  $S$ . The overall effect is, quite counterintuitively, a decrease of the lower bound in (37). We note that the bound decreases linearly with  $T$ , in contrast with (26)'s exponential decrease. This linear decrease will drive the bound to zero with finite samples. We commented previously that  $\frac{1}{\sigma^2}$  plays the role of the signal-to-noise ratio. Therefore, increasing the SNR will also force the bound to zero.

The lower bound (37) is loose in the sense that when  $T$ ,  $\|A\|_F^2$ , or the SNR is large enough the bound becomes negative, but when there is noise, perfect support recovery is generally impossible. We believe the limitation is inherited mainly from the Fano's inequality (33) rather than other intermediate inequalities. While the original Fano's inequality

$$H(P_{\text{err}}) + P_{\text{err}} \log(L - 1) \geq H(\mathbf{S}|\mathbf{Y})$$

is tight [47], the adoption of the average divergence (32) as an upper bound for the mutual information  $I(\mathbf{S}; \mathbf{Y})$  between the random support  $\mathbf{S}$  and the observation  $\mathbf{Y}$  results in an over-loose estimate (see the proof of (33) in [50]).

In the work of [6], each column of  $A$  is the array manifold vector function evaluated at a sample of the direction parameter. The implication of bound (37) for optimal design is that we should construct an array whose geometry leads to maximal  $\|A\|_F^2$ . However, under the narrowband signal assumption and narrowband array assumption [51], the array manifold vector for isotropic sensor arrays always has norm  $\sqrt{M}$  [52], which means that  $\|A\|_F^2 = MN$ . Hence in this case, the probability of error is always bounded by

$$P_{\text{err}} \geq 1 - \frac{\kappa T \frac{K}{N} \left(1 - \frac{K}{N}\right) MN}{2\sigma^2 \ln \binom{N}{K}} + o(1).$$

Therefore, we have the following theorem,

**Theorem 3** *Under the narrowband signal assumption and narrowband array assumption, for an isotropic sensor array in the DOA estimation scheme proposed in [6], suppose  $M \geq 2K$ . Then in order to let the probability of error  $P_{\text{err}} < \varepsilon$  with  $\varepsilon > 0$  for any decision rule, the number of measurements must satisfy the following:*

$$\begin{aligned} MT &\geq (1 - \varepsilon) \frac{2\sigma^2 \ln \binom{N}{K}}{K \left(1 - \frac{K}{N}\right)} + o(1) \\ &\geq (1 - \varepsilon) 2\sigma^2 \ln \frac{N}{K} + o(1). \end{aligned}$$

We comment that the same lower bound applies to Fourier measurement matrix (not normalized by  $1/\sqrt{M}$ ) due to the same line of argument. We will not explicitly present this result in the current paper.

Note that Theorem 3 is most strong when the number of sources is small, which is usually the case in the practice of radar and sonar. Our result shows that the number of samples is lower bounded by  $\ln N$  in this case. Note that  $N$  is the number of intervals we use to divide the whole range of DOA; hence, it is a measure of resolution. Therefore, the number of samples only needs to increase in the logarithm of  $N$ , which is very desirable. The symmetric roles played by  $M$  and  $T$  are also desirable since  $M$  is the number of sensors and is expensive to increase. As a consequence, we simply increase the number of samples to achieve a desired probability of error as long as  $M$  is greater than  $2K$ .

## V. SUPPORT RECOVERY FOR THE GAUSSIAN MEASUREMENT ENSEMBLE

In this section, we refine our results in previous sections from general non-degenerate measurement matrices to the Gaussian ensemble. Unless otherwise specified, we always assume that the

elements in a measurement matrix  $\mathbf{A}$  are *i.i.d.* samples from unit variance real or complex normal distributions. The Gaussian measurement ensemble is widely used and studied in compressive sensing [1]–[5]. The additional structure and the theoretical tools available enable us to derive deeper results in this case. We first show two corollaries on the eigenvalue structures of the Gaussian measurement matrix. Then we derive sufficient and necessary conditions in terms of  $M, N, K$  and  $T$  for the system to have a vanishing probability of error.

#### A. Eigenvalue Structure for a Gaussian Measurement Matrix

First, we observe that a Gaussian measurement matrix is non-degenerate with probability one, since any  $p \leq M$  random vectors  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p$  from  $\mathbb{FN}(0, \Sigma)$  with  $\Sigma \in \mathbb{R}^{M \times M}$  positive definite are linearly independent with probability one (refer to Theorem 3.2.1 in [42]). As a consequence, we have

**Corollary 2** *For Gaussian measurement matrix  $\mathbf{A}$ , let  $\mathbf{H} = \Sigma_0^{1/2} \Sigma_1^{-1} \Sigma_0^{1/2}$ ,  $k_i = |S_0 \cap S_1|$ ,  $k_0 = |S_0 \setminus S_1| = |S_0| - k_i$ ,  $k_1 = |S_1 \setminus S_0| = |S_1| - k_i$  and assume  $M \geq k_0 + k_1$ . Then with probability one,  $k_0$  eigenvalues of matrix  $\mathbf{H}$  are greater than 1,  $k_1$  less than 1, and  $M - (k_0 + k_1)$  equal to 1.*

We refine Proposition 3 based on the well-known QR factorization for Gaussian matrices [42], [53].

**Corollary 3** *With the same assumptions as in Corollary 2, then with probability one, the sorted eigenvalues of  $\mathbf{H}$  that are greater than 1 are*

- 1) *lower bounded by the corresponding ones of  $\mathbf{I}_{k_0} + \frac{1}{\sigma^2} \mathbf{R}_{33} \mathbf{R}_{33}^\dagger$ , where the elements of*

$\mathbf{R}_{33} = (r_{mn})_{k_0 \times k_0}$  *satisfy*

$$2\kappa r_{mn}^2 \sim \chi_{2\kappa(M-k_1-k_i-m+1)}^2, m = 1, \dots, k_0,$$

$$r_{mn} \sim \mathbb{FN}(0, 1), 1 \leq m < n \leq k_0.$$

- 2) *upper bounded by the corresponding ones of  $\mathbf{I}_{k_0} + \frac{1}{\sigma^2} \mathbf{A}_{S_0 \setminus S_1} \mathbf{A}_{S_0 \setminus S_1}^\dagger$ .*

Now with the distributions on the elements of the bounding matrices, we can give sharp estimates on  $\bar{\lambda}_{S_i, S_j}$  and  $\hat{\lambda}_{S_i, S_j}$ . In particular, we have the following theorem:

**Proposition 4** *For Gaussian measurement matrix  $\mathbf{A}$ , suppose  $S_i$  and  $S_j$  are a pair of distinct supports with the same size  $K$ . Then we have*

$$\begin{aligned} 1 + \frac{M}{\sigma^2} &\geq \mathbb{E} \bar{\lambda}_{S_i, S_j} \geq 1 + \frac{M - K - k_d}{\sigma^2} \\ 1 + \frac{M}{\sigma^2} &\geq \mathbb{E} \hat{\lambda}_{S_i, S_j} \geq 1 + \frac{M - K}{\sigma^2}. \end{aligned}$$

**Proof:** We copy the inequalities (29), (30) on  $\bar{\lambda}_{S_i, S_j}$  and  $\hat{\lambda}_{S_i, S_j}$  here:

$$1 + \frac{1}{\sigma^2 k_d} \sum_{\substack{1 \leq m \leq M \\ n \in S_i \setminus S_j}} |A_{mn}|^2 \geq \bar{\lambda}_{S_i, S_j} \geq 1 + \frac{1}{\sigma^2} \left( \prod_{m=1}^{k_d} |r_{mm}|^2 \right)^{1/k_d}.$$

Similarly, we can derive bounds for  $\hat{\lambda}_{S_i, S_j}$ :

$$1 + \frac{1}{\sigma^2 k_d} \sum_{\substack{1 \leq m \leq M \\ n \in S_i \setminus S_j}} |A_{mn}|^2 \geq \hat{\lambda}_{S_i, S_j} \geq 1 + \frac{1}{\sigma^2 k_d} \sum_{m < n} |r_{mn}|^2.$$

We notice that the upper bound in both cases is the same. The proof then reduces to the computation of three expectations, two of which are trivial:

$$\begin{aligned} \mathbb{E} \frac{1}{\sigma^2 k_d} \sum_{\substack{1 \leq m \leq M \\ n \in S_0 \setminus S_1}} |A_{mn}|^2 &= \frac{M}{\sigma^2}, \\ \mathbb{E} \frac{1}{\sigma^2 k_d} \sum_{m < n} |r_{mn}|^2 &= \frac{M - K}{\sigma^2}. \end{aligned}$$

Next, the independence of the  $r_{nn}$ 's and the convexity of exponential functions together with Jensen's inequality yield

$$\begin{aligned} &\mathbb{E} \frac{1}{\sigma^2} \left( \prod_{n=1}^{k_d} r_{nn}^2 \right)^{1/k_d} \\ &= \frac{1}{2\kappa\sigma^2} \mathbb{E} \exp \left[ \frac{1}{k_d} \sum_{n=1}^{k_d} \ln (2\kappa r_{nn}^2) \right] \\ &\geq \frac{1}{2\kappa\sigma^2} \exp \left[ \frac{1}{k_d} \sum_{n=1}^{k_d} \mathbb{E} \ln (2\kappa r_{nn}^2) \right]. \end{aligned}$$

Since  $(2\kappa r_{nn}^2) \sim \chi_{2\kappa(M-K-n+1)}^2$ , the expectation of logarithm is  $E \ln(2\kappa r_{nn}^2) = \ln 2 + \psi(\kappa(M-K-n+1))$ , where  $\psi(z) = \frac{\Gamma'(z)}{\Gamma(z)}$  is the digamma function. Note that  $\psi(z)$  is increasing and satisfies  $\psi(z+1) \geq \ln z$ . Therefore, we have

$$\begin{aligned}
& \mathbb{E} \frac{1}{\sigma^2} \left( \prod_{n=1}^{k_d} r_{nn}^2 \right)^{1/k_d} \\
& \geq \frac{1}{2\kappa\sigma^2} \exp \left[ \frac{1}{k_d} \sum_{n=1}^{k_d} \mathbb{E} \ln(2\kappa r_{nn}^2) \right] \\
& \geq \frac{1}{2\kappa\sigma^2} \exp \left[ \ln 2 + \frac{1}{k_d} \sum_{n=1}^{k_d} \psi(\kappa(M-K-n+1)) \right] \\
& \geq \frac{1}{\kappa\sigma^2} \exp[\psi(\kappa(M-K-k_d+1))] \\
& \geq \frac{1}{\kappa\sigma^2} \exp[\ln(\kappa(M-K-k_d))] \\
& \geq \frac{M-K-k_d}{\sigma^2}. \quad \blacksquare
\end{aligned}$$

The expected values for the critical quantities  $\bar{\lambda}_{S_i, S_j}$  and  $\hat{\lambda}_{S_i, S_j}$  lie between  $1 + \frac{M-2K}{\sigma^2}$  and  $1 + \frac{M}{\sigma^2}$ , linearly proportional to  $M$ . Note that in conventional compressive sensing, the variance of the elements of  $\mathbf{A}$  is usually taken to be  $\frac{1}{M}$ , which is equivalent to scaling the noise variance  $\sigma^2$  to  $M\sigma^2$  in our model. The resultant  $\bar{\lambda}_{S_i, S_j}$  and  $\hat{\lambda}_{S_i, S_j}$  are then centered between  $1 + \frac{1-2\frac{K}{M}}{\sigma^2}$  and  $1 + \frac{1}{\sigma^2}$ .

### B. Necessary Condition

One fundamental problem in compressive sensing is how many samples should the system take to guarantee a stable reconstruction. Although many sufficient conditions are available, non-trivial necessary conditions are rare. Besides, in previous works, stable reconstruction has been measured in the sense of  $l_p$  norms between the reconstructed signal and the true signal. In this section, we derive two necessary conditions on  $M$  and  $T$  in terms of  $N$  and  $K$  in order to guarantee respectively that, first,  $\mathbb{E}P_{\text{err}}$  turns to zeros and, second, for majority realizations of  $\mathbf{A}$ , the probability of error vanishes. More precisely, we have the following theorem:

**Theorem 4** *In the support recovery problem (5), for any  $\varepsilon, \delta > 0$ , a necessary condition of*



$\mathbb{E}P_{\text{err}} < \varepsilon$  is

$$\frac{\kappa MT}{\sigma^2} \geq (1 - \varepsilon) \frac{2 \ln \binom{N}{K}}{K \left(1 - \frac{K}{N}\right)} + o(1) \quad (38)$$

$$\geq (1 - \varepsilon) 2 \ln \frac{N}{K} + o(1), \quad (39)$$

and a necessary condition of  $\mathbb{P}\{P_{\text{err}}(\mathbf{A}) \leq \varepsilon\} \geq 1 - \delta$  is

$$\frac{\kappa MT}{\sigma^2} \geq (1 - \varepsilon - \delta) \frac{2 \ln \binom{N}{K}}{K \left(1 - \frac{K}{N}\right)} + o(1) \quad (40)$$

$$\geq (1 - \varepsilon - \delta) 2 \ln \frac{N}{K} + o(1). \quad (41)$$

**Proof:** Equation (37) and  $\mathbb{E}\|A\|_F^2 = \sum_{m,l} \mathbb{E}|A_{ml}|^2 = MN$  give

$$\mathbb{E}P_{\text{err}} \geq 1 - \frac{\kappa T \frac{K}{N} \left(1 - \frac{K}{N}\right) MN}{2\sigma^2 \ln \binom{N}{K}} + o(1).$$

Hence,  $\mathbb{E}P_{\text{err}} < \varepsilon$  entails

$$\begin{aligned} \frac{\kappa MT}{\sigma^2} &\geq (1 - \varepsilon) \frac{2 \ln \binom{N}{K}}{K \left(1 - \frac{K}{N}\right)} + o(1) \\ &\geq (1 - \varepsilon) 2 \ln \frac{N}{K} + o(1), \end{aligned}$$

Denote by  $E$  the event  $\{\mathbf{A} : P_{\text{err}}(\mathbf{A}) \leq \varepsilon\}$ ; then  $\mathbb{P}\{E^c\} \leq \delta$  and we have

$$\begin{aligned} \mathbb{E}P_{\text{err}} &= \int_E P_{\text{err}}(\mathbf{A}) + \int_{E^c} P_{\text{err}}(\mathbf{A}) \\ &\leq \varepsilon \mathbb{P}(E) + \mathbb{P}(E^c) \\ &\leq \varepsilon + \delta. \end{aligned}$$

Therefore, from the first part of the theorem, we obtain

$$\begin{aligned} \frac{\kappa MT}{\sigma^2} &\geq (1 - \varepsilon - \delta) \frac{2 \ln \binom{N}{K}}{K \left(1 - \frac{K}{N}\right)} + o(1) \\ &\geq (1 - \varepsilon - \delta) 2 \ln \frac{N}{K} + o(1). \quad \blacksquare \end{aligned}$$

Our result shows that as far as support recovery is concerned, one cannot avoid the  $\ln \frac{N}{K}$  term when only given one temporal sample. Worse, for conventional compressive sensing with a measurement matrix generated from a Gaussian random variable with variance  $1/M$ , the necessary condition becomes

$$\begin{aligned}
T &\geq \frac{2\sigma^2 \ln \binom{N}{K}}{\kappa K \left(1 - \frac{K}{N}\right)} + o(1) \\
&\geq \frac{2\sigma^2}{\kappa} \ln \frac{N}{K} + o(1),
\end{aligned}$$

which is independent of  $M$ . Therefore, it is impossible to have a vanishing  $\mathbb{E}P_{\text{err}}$  no matter how large an  $M$  one takes. Basically this situation arises because while taking more samples, one scales down the measurement gains  $A_{ml}$ , which effectively reduces the SNR and thus is not helpful in support recovery. As discussed below Theorem 3,  $\ln \binom{N}{K}$  is the uncertainty of the support variable  $S$ , and  $\ln \frac{N}{K}$  actually comes from it. Therefore, it is no surprise that the number of samples is determined by this quantity and cannot be made independent of it.

### C. Sufficient Condition

We derive a sufficient condition in parallel with sufficient conditions in compressive sensing. In compressive sensing, when only one temporal sample is available,  $M = \Omega \left( K \ln \frac{N}{K} \right)$  is enough for stable signal reconstruction for the majority of the realizations of measurement matrix  $A$  from a Gaussian ensemble with variance  $\frac{1}{M}$ . As shown in the previous subsection, if we take the probability of error for support recovery as a performance measure, it is impossible in this case to recover the support with vanishing probability of error. Therefore, we consider a Gaussian ensemble with unit variance. We first establish a lemma to estimate the lower tail of the distribution for  $\bar{\lambda}_{S_i, S_j}$ . We have shown that the  $\mathbb{E} \left( \bar{\lambda}_{S_i, S_j} \right)$  lie between  $1 + \frac{M-2K}{\sigma^2}$  and  $1 + \frac{M}{\sigma^2}$ . When  $\gamma$  is much less than  $1 + \frac{M-2K}{\sigma^2}$ , we expect that  $\mathbb{P} \left\{ \bar{\lambda}_{S_i, S_j} \leq \gamma \right\}$  decays quickly. More specifically, we have the following large deviation lemma:

**Lemma 1** *Suppose that  $\gamma = \frac{1}{3} \frac{M-2K}{\sigma^2}$ . Then there exists constant  $c > 0$  such that for  $M - 2K$  sufficiently large, we have*

$$\mathbb{P} \left\{ \bar{\lambda}_{S_i, S_j} \leq \gamma \right\} \leq \exp \left[ -c (M - 2K) \right].$$

This large deviation lemma together with the union bound gives us the following sufficient condition for support recovery:

**Theorem 5** *Suppose that*

$$M = \Omega \left( K \ln \frac{N}{K} \right) \tag{42}$$

and

$$\kappa T \ln \frac{M}{\sigma^2} \gg \ln [K(N - K)]. \quad (43)$$

Then given any realization of measurement matrix  $\mathbf{A}$  from a Gaussian ensemble, the optimal decision rule (15) for multiple support recovery problem (5) has a vanishing  $P_{\text{err}}$  with probability turning to one. In particular, if  $M = \Omega\left(K \ln \frac{N}{K}\right)$  and

$$T \gg \frac{\ln N}{\ln \ln N}, \quad (44)$$

then the probability of error turns to zero as  $N$  turns to infinity.

**Proof:** Denote  $\gamma = \frac{1}{3} \frac{M-2K}{\sigma^2}$ . Then according to the union bound, we have

$$\begin{aligned} & \mathbb{P}\{\bar{\lambda} \leq \gamma\} \\ &= \mathbb{P}\left\{\bigcup_{S_i \neq S_j} [\bar{\lambda}_{S_i, S_j} \leq \gamma]\right\} \\ &\leq \sum_{S_i \neq S_j} \mathbb{P}\{\bar{\lambda}_{S_i, S_j} \leq \gamma\}. \end{aligned}$$

Therefore, by application of Lemma 1,

$$\begin{aligned} & \mathbb{P}\{\bar{\lambda} \leq \gamma\} \\ &\leq \binom{N}{K}^2 K \exp\{-c(M - 2K)\} \\ &\leq \exp\left[-c(M - 2K) + 2K \log \frac{N}{K} + \log K\right]. \end{aligned}$$

Hence, as long as  $M = \Omega\left(K \log \frac{N}{K}\right)$ , we know that the exponent turns to  $-\infty$  as  $N \rightarrow \infty$ .

We now define  $E = \{\mathbf{A} : \bar{\lambda}(\mathbf{A}) > \gamma\}$ , where  $\mathbb{P}\{E\}$  approaches one as  $N$  turns to infinity. Now the upper bound (26) becomes

$$\begin{aligned} P_{\text{err}} &= O\left(\frac{K(N - K)}{\left(\frac{\bar{\lambda}}{12\sigma^2}\right)^{\kappa T}}\right) \\ &= O\left(\frac{K(N - K)}{\left(\frac{M}{\sigma^2}\right)^{\kappa T}}\right). \end{aligned}$$

Hence, if  $\kappa T \ln \frac{M}{\sigma^2} \gg \ln [K(N - K)]$ , we get a vanishing probability of error. In particular, under the assumption that  $M \geq \Omega\left(K \log \frac{N}{K}\right)$ , if  $T \gg \frac{\ln N}{\ln \ln N}$ , then  $\frac{\ln [K(N - K)]}{\ln [K \ln \frac{N}{K}]} \leq \frac{\ln N}{\ln \ln N}$  implies that  $K(N - K) \ll O\left(\left(\frac{K \ln \frac{N}{K}}{\sigma^2}\right)^{\kappa T}\right) = O\left(\frac{K(N - K)}{\left(\frac{M}{\sigma^2}\right)^{\kappa T}}\right)$  for suitably selected constants. ■

We now consider several special cases and explore the implications of the sufficient conditions. The discussions are heuristic in nature and their validity requires further checking.

If we set  $T = 1$ , then we need  $M$  to be much greater than  $N$  to guarantee a vanishing probability  $P_{\text{err}}$ . This restriction suggests that even if we have more observations than the original signal length  $N$ , in which case we can obtain the original sparse signal by solving a least square problem, we still might not be able to get the correct support because of the noise, if  $M$  is not sufficiently large compared to  $N$ . We discussed in the introduction that for many applications, the support of a signal has significant physical implications and its correct recovery is of crucial importance. Therefore, without multiple temporal samples, the scheme proposed by compressive sensing is questionable as far as support recovery is concerned. Worse, if we set the variance for the elements in  $\mathbf{A}$  to be  $1/M$  as in compressive sensing, which is equivalent to replacing  $\sigma^2$  with  $M\sigma^2$ , even increasing the number of temporal samples will not improve the probability of error significantly unless the noise variance is very small. Hence, using support recovery as a criterion, one cannot expect the compressive sensing scheme to work very well in the low SNR case. This conclusion is not a surprise, since we reduce the number of samples to achieve compression.

Another special case is when  $K = 1$ . In this case, the sufficient condition becomes  $M \geq \ln N$  and  $\kappa T \ln \frac{M}{\sigma^2} \gg \ln N$ . Now the number of total samples should satisfy  $MT \gg \frac{(\ln N)^2}{\ln \ln N}$  while the necessary condition states that  $MT = \Omega(\ln N)$ . The smallest gap between the necessary condition and sufficient condition is achieved when  $K = 1$ .

The result also exhibits several interesting properties in the general case. Compared with the necessary condition (40) and (41), the asymmetry in the sufficient condition is even more desirable in most cases because of the asymmetric cost associated with sensors and temporal samples. Once the threshold  $K \ln \frac{N}{K}$  of  $M$  is exceeded, we can achieve a desired probability of error by taking more temporal samples. If we were concerned only with the number of total samples, we would minimize  $MT$  subject to the constraints (42) and (43) to achieve a given level of probability of error. However, in applications for which timing is important, one has to increase sensors to reduce  $P_{\text{err}}$  to a certain limit.

The sufficient condition (42), (43), and (44) is separable in the following sense. We observe from the proof that the requirement  $M = \Omega\left(K \ln \frac{N}{K}\right)$  is used only to guarantee that the randomly generated measurement matrix is a good one in the sense that its incoherence  $\bar{\lambda}$  is sufficiently

large, as in the case of compressive sensing. It is in Lemma 1 that we use the Gaussian ensemble assumption. If another deterministic construction procedure (for attempts in this direction, see [54]) or random distribution gave measurement matrix with better incoherence  $\bar{\lambda}$ , it would be possible to reduce the orders for both  $M$  and  $T$ .

## VI. CONCLUSIONS

In this paper, we formulated the support recovery problems for jointly sparse signals as binary and multiple-hypothesis testings. Adopting the probability of error as the performance criterion, the optimal decision rules are given by the likelihood ratio test and the maximum *a posteriori* probability estimator. The latter reduces to the maximum likelihood estimator when equal prior probabilities are assigned to the supports. We then employed the Chernoff bound and Fano's inequality to derive bounds on the probability of error. We discussed the implications of these bounds at the end of Section III-B, Section III-C, Section IV, Section V-B, and Section V-C, in particular when they are applied to the DOA estimation problem considered in [6] and compressive sensing with a Gaussian measurement ensemble. We derived sufficient and necessary conditions to achieve a vanishing probability of error in both the mean and large probability senses. These conditions show the necessity of considering multiple temporal samples. The symmetric and asymmetric roles played by the spatial and temporal samples and their implications in system design were discussed. For compressive sensing, we demonstrated that it is impossible to obtain accurate signal support with only one temporal sample if the variance for the Gaussian measurement matrix scales with  $1/M$ .

This research on support recovery for jointly sparse signals is far from complete. Several questions remain to be answered. First, we notice an obvious gap between the necessary and sufficient conditions even in the simplest case with  $K = 1$ . Better techniques need to be introduced to refine the results. Second, as in the case for RIP, computation of the quantity  $\bar{\lambda}$  for an arbitrary measurement matrix is extremely difficult. Although we derive large deviation bounds on  $\bar{\lambda}$  and compute the expected value for  $\bar{\lambda}_{S_i, S_j}$  for the Gaussian ensemble, its behaviors in both the general and Gaussian cases require further study. Its relationship with RIP also needs to be clarified. Finally, our lower bound derived from Fano's inequality identifies only the role played by the total gain. The effect of the measurement matrix's incoherence is elusive. The answers to these questions will enhance our understanding of the measurement mechanism of

(2).

## APPENDIX A

### PROOF OF PROPOSITION 1

In this proof, we focus on the case for which both  $k_0 \neq 0$  and  $k_1 \neq 0$ . Other cases have similar and simpler proofs. The eigenvalues of  $H$  satisfy  $|\lambda I_M - H| = 0$ , which is equivalent to  $|\lambda \Sigma_1 - \Sigma_0| = 0$ . The substitution  $\mu = \lambda - 1$  defines

$$g(\mu) = |(\mu + 1) \Sigma_1 - \Sigma_0| = |\mu \Sigma_1 - (\Sigma_0 - \Sigma_1)|.$$

The following algebraic manipulation

$$\begin{aligned} G &\triangleq \Sigma_0 - \Sigma_1 \\ &= A_{S_0} A_{S_0}^\dagger - A_{S_1} A_{S_1}^\dagger \\ &= \left[ A_{S_0 \cap S_1} A_{S_0 \cap S_1}^\dagger + A_{S_0 \setminus S_1} A_{S_0 \setminus S_1}^\dagger \right] \\ &\quad - \left[ A_{S_0 \cap S_1} A_{S_0 \cap S_1}^\dagger + A_{S_1 \setminus S_0} A_{S_1 \setminus S_0}^\dagger \right] \\ &= A_{S_0 \setminus S_1} A_{S_0 \setminus S_1}^\dagger - A_{S_1 \setminus S_0} A_{S_1 \setminus S_0}^\dagger \end{aligned}$$

leads to

$$\begin{aligned} g(\mu) &= |\mu \Sigma_1 - G| \\ &= |\Sigma_1|^{\frac{1}{2}} \left| \mu I_M - \Sigma_1^{-\frac{1}{2}} G \Sigma_1^{-\frac{1}{2} \dagger} \right| |\Sigma_1|^{\frac{1}{2}}. \end{aligned}$$

Therefore, to prove the theorem, it suffices to show that  $\Sigma_1^{-\frac{1}{2}} G \Sigma_1^{-\frac{1}{2} \dagger}$  has  $k_0$  positive eigenvalues,  $k_1$  negative eigenvalues and  $M - (k_0 + k_1)$  zero eigenvalues or, put another way,  $\Sigma_1^{-\frac{1}{2}} G \Sigma_1^{-\frac{1}{2} \dagger}$  has inertia  $(k_0, k_1, M - (k_0 + k_1))$ . The Sylvester's law of inertia ([55], Theorem 4.5.8, p. 223) states that the inertia of a symmetric matrix is invariant under congruence transformations. Hence, we need only to show that  $G$  has inertia  $(k_0, k_1, M - (k_0 + k_1))$ . Clearly  $G = PQ^\dagger$  with  $P = \begin{bmatrix} A_{S_0 \setminus S_1} & A_{S_1 \setminus S_0} \end{bmatrix}$  and  $Q = \begin{bmatrix} A_{S_0 \setminus S_1} & -A_{S_1 \setminus S_0} \end{bmatrix}$ . To find the number of zero eigenvalues of  $G$ , we calculate the rank of  $G$ . The non-degenerateness of measurement matrix  $A$  implies that  $\text{rank}(P) = \text{rank}(Q) = k_0 + k_1$ . Therefore, from rank inequality ([55], Theorem 0.4.5, p.

13),

$$\begin{aligned} & \text{rank}(P) + \text{rank}(Q^\dagger) - (k_0 + k_1) \\ & \leq \text{rank}(PQ^\dagger) \\ & \leq \min \{ \text{rank}(P), \text{rank}(Q^\dagger) \}, \end{aligned}$$

and we conclude that  $\text{rank}(G) = k_0 + k_1$ .

To count the number of negative eigenvalues of  $G$ , we use the Jacobi-Sturm rule ([56], Theorem A.1.4, p. 320), which states that for an  $M \times M$  symmetric matrix whose  $j$ th leading principal minor has determinant  $d_j, j = 1, \dots, M$ , the number of nonnegative eigenvalues is equal to the number of sign changes of sequence  $1, d_1, \dots, d_M$ . We consider only the first  $k_0 + k_1$  leading principal minors, since higher order minors have determinant 0.

Suppose  $I = \{1, \dots, k_0 + k_1\}$  is an index set. Without loss of generality, we assume that  $P^I$  is nonsingular. Applying  $QL$  factorization (one variation of  $QR$  factorization, see [57]) to matrix  $P^I$ , we obtain  $P^I = OL$ , where  $O$  is an orthogonal matrix,  $OO^\dagger = I_{k_0+k_1}$ , and  $L = (l_{ij})_{(k_0+k_1) \times (k_0+k_1)}$  is a lower triangular matrix. The diagonal entries of  $L$  are nonzero since  $P^I$  is nonsingular. The partition of  $L$  into

$$L = \begin{bmatrix} L_1 & L_2 \end{bmatrix}$$

with  $L_1 \in \mathbb{F}^{(k_0+k_1) \times k_0}, L_2 \in \mathbb{F}^{(k_0+k_1) \times k_1}$ , and  $L_2 = \begin{bmatrix} 0 \\ L_3 \end{bmatrix}$  with  $L_3 \in \mathbb{F}^{k_1 \times k_1}$  implies

$$G_I^I = P^I(Q_I)^\dagger = O \begin{bmatrix} L_1 & L_2 \end{bmatrix} \begin{bmatrix} L_1^\dagger \\ -L_2^\dagger \end{bmatrix} O^\dagger.$$

Again using the invariance property of inertia under congruence transformation, we focus on the leading principal minors of  $U \triangleq \begin{bmatrix} L_1 & L_2 \end{bmatrix} \begin{bmatrix} L_1^\dagger \\ -L_2^\dagger \end{bmatrix}$ . Suppose  $J = \{1, \dots, j\}$ . For  $1 \leq j \leq k_0$ , from the lower triangularity of  $L$ , it is clear that

$$|(U_J^J)| = |(L_1)_J^J|^2 = \prod_{i=1}^j |l_{ii}|^2 > 0.$$

For  $k_0 + 1 \leq j \leq k_0 + k_1$ , suppose  $J_0 = \{1, \dots, k_0\}$  and  $J_1 = \{1, \dots, j - k_0\}$ . We then have

$$\begin{aligned} |U_J^J| &= \left| (L_1)_{J_0}^{J_0} \right|^2 \left| (L_3)_{J_1}^{J_1} \right| - \left| (L_3)_{J_1}^{J_1} \right|^\dagger \\ &= (-1)^{j-k_0} \left| (L_1)_{J_0}^{J_0} \right|^2 \left| (L_3)_{J_1}^{J_1} \right|^2 \\ &= (-1)^{j-k_0} \prod_{i=1}^j |l_{ii}|^2. \end{aligned}$$

Therefore, the sequence  $1, d_1, d_2, \dots, d_{k_0+k_1}$  has  $k_1$  sign changes, which implies that  $G_I^I$ —hence  $G$ —has  $k_1$  negative eigenvalues. Finally, we conclude that the theorem holds for  $H$ .

## APPENDIX B

### PROOF OF PROPOSITION 3

We first prove the first claim. From the proof of Proposition 1, it suffices to prove that the sorted positive eigenvalues of  $\Sigma_1^{-\frac{1}{2}} G \Sigma_1^{-\frac{1}{2} \dagger}$  are greater than those of  $\frac{1}{\sigma^2} R_{33} R_{33}^\dagger$ , where  $G = A_{S_0 \setminus S_1} A_{S_0 \setminus S_1}^\dagger - A_{S_1 \setminus S_0} A_{S_1 \setminus S_0}^\dagger$ . Since cyclic permutation of a matrix product does not change its eigenvalues, we restrict ourselves to  $\Sigma_1^{-1} G$ . Consider the  $QR$  decomposition

$$\begin{aligned} \begin{bmatrix} A_{S_1 \setminus S_0} & A_{S_1 S_0} & A_{S_0 \setminus S_1} \end{bmatrix} &= QR \\ &\triangleq \begin{bmatrix} Q_1 & Q_2 & Q_3 & Q_4 \end{bmatrix} \begin{bmatrix} R_{11} & R_{12} & R_{13} \\ 0 & R_{22} & R_{23} \\ 0 & 0 & R_{33} \\ 0 & 0 & 0 \end{bmatrix}, \end{aligned}$$

where  $Q \in \mathbb{F}^{M \times M}$  is an orthogonal matrix with partitions  $Q_1 \in \mathbb{F}^{M \times k_1}, Q_2 \in \mathbb{F}^{M \times k_i}, Q_3 \in \mathbb{F}^{M \times k_0}, R \in \mathbb{F}^{M \times (k_1 + k_i + k_0)}$  is an upper triangular matrix with partitions  $R_{11} \in \mathbb{F}^{k_1 \times k_1}, R_{22} \in \mathbb{F}^{k_i \times k_i}, R_{33} \in \mathbb{F}^{k_0 \times k_0}$ , and other submatrices have corresponding dimensions.



First, we note that

$$\begin{aligned}
& Q^\dagger G Q \\
&= \begin{bmatrix} R_{13} & R_{11} \\ R_{23} & 0 \\ R_{33} & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} R_{13}^\dagger & R_{23}^\dagger & R_{33}^\dagger & 0 \\ -R_{11}^\dagger & 0 & 0 & 0 \end{bmatrix} \\
&= \begin{bmatrix} \begin{bmatrix} R_{13} & R_{11} \\ R_{23} & 0 \end{bmatrix} \begin{bmatrix} R_{13}^\dagger & R_{23}^\dagger \\ -R_{11}^\dagger & 0 \end{bmatrix} \begin{bmatrix} R_{13}R_{33}^\dagger \\ R_{23}R_{33}^\dagger \end{bmatrix} & 0 \\ \begin{bmatrix} R_{33}R_{13}^\dagger & R_{33}R_{23}^\dagger \end{bmatrix} & R_{33}R_{33}^\dagger \\ 0 & 0 \end{bmatrix}.
\end{aligned}$$

Therefore, the last  $M - (k_1 + k_i + k_0)$  rows and columns of  $G$ —and hence of  $\Sigma_1^{-1}G$ —are zeros, which lead to the  $M - (k_1 + k_i + k_0)$  zero eigenvalues of  $\Sigma_1^{-\frac{1}{2}}G\Sigma_1^{-\frac{1}{2}\dagger}$ . We then drop these rows and columns in all matrices involved in subsequent analysis. In particular, the submatrix of  $Q^\dagger \Sigma_1 Q = Q^\dagger (\sigma^2 \mathbf{I}_M + A_{S_1} A_{S_1}^\dagger) Q$  without the last  $M - (k_1 + k_i + k_0)$  rows and columns is

$$\begin{aligned}
& \sigma^2 \mathbf{I}_M + \begin{bmatrix} R_{11} & R_{12} \\ 0 & R_{22} \\ 0 & 0 \end{bmatrix} \begin{bmatrix} R_{11}^\dagger & 0 & 0 \\ R_{12}^\dagger & R_{22}^\dagger & 0 \end{bmatrix} \\
&= \begin{bmatrix} \sigma^2 \mathbf{I}_{k_1+k_i} + \begin{bmatrix} R_{11} & R_{12} \\ 0 & R_{22} \end{bmatrix} \begin{bmatrix} R_{11}^\dagger & 0 \\ R_{12}^\dagger & R_{22}^\dagger \end{bmatrix} & 0 \\ 0 & \sigma^2 \mathbf{I}_{k_0} \end{bmatrix} \\
&\triangleq \begin{bmatrix} F & 0 \\ 0 & \sigma^2 \mathbf{I}_{k_0} \end{bmatrix}.
\end{aligned}$$

Define

$$\begin{aligned} & \begin{bmatrix} V & K^\dagger \\ K & R_{33}R_{33}^\dagger \end{bmatrix} \\ \triangleq & \begin{bmatrix} \begin{bmatrix} R_{13} & R_{11} \\ R_{23} & 0 \end{bmatrix} \begin{bmatrix} R_{13}^\dagger & R_{23}^\dagger \\ -R_{11}^\dagger & 0 \end{bmatrix} \begin{bmatrix} R_{13}R_{33}^\dagger \\ R_{23}R_{33}^\dagger \end{bmatrix} \\ \begin{bmatrix} R_{33}R_{13}^\dagger & R_{33}R_{23}^\dagger \end{bmatrix} & R_{33}R_{33}^\dagger \end{bmatrix}. \end{aligned}$$

Due to the invariance of eigenvalues with respect to orthogonal transformations and switching to the symmetrized version, we focus on

$$\begin{aligned} & \begin{bmatrix} F & 0 \\ 0 & \sigma^2 \mathbf{I}_{k_0} \end{bmatrix}^{-\frac{1}{2}} \begin{bmatrix} V & K^\dagger \\ K & R_{33}R_{33}^\dagger \end{bmatrix} \begin{bmatrix} F & 0 \\ 0 & \sigma^2 \mathbf{I}_{k_0} \end{bmatrix}^{-\frac{1}{2}\dagger} \\ = & \begin{bmatrix} F^{-\frac{1}{2}} V F^{-\frac{1}{2}\dagger} & F^{-\frac{1}{2}} K^\dagger \frac{1}{\sigma} \\ \frac{1}{\sigma} K F^{-\frac{1}{2}} & \frac{1}{\sigma^2} R_{33}R_{33}^\dagger \end{bmatrix}. \end{aligned}$$

Next we argue that the sorted positive eigenvalues of  $\begin{bmatrix} F^{-\frac{1}{2}} V F^{-\frac{1}{2}\dagger} & F^{-\frac{1}{2}} K^\dagger \frac{1}{\sigma} \\ \frac{1}{\sigma} K F^{-\frac{1}{2}} & \frac{1}{\sigma^2} R_{33}R_{33}^\dagger \end{bmatrix}$  are greater than the corresponding sorted eigenvalues of  $\frac{1}{\sigma^2} R_{33}R_{33}^\dagger$ .

For any  $\varepsilon > 0$ , we define a matrix  $M_{\varepsilon,N} = \begin{bmatrix} -N\mathbf{I}_{k_1+k_i} & 0 \\ 0 & \frac{1}{\sigma^2} R_{33}R_{33}^\dagger - \varepsilon \mathbf{I}_{k_0} \end{bmatrix}$ . Then we have

$$\begin{aligned} & \begin{bmatrix} F^{-\frac{1}{2}} V F^{-\frac{1}{2}\dagger} & F^{-\frac{1}{2}} K^\dagger \frac{1}{\sigma} \\ \frac{1}{\sigma} K F^{-\frac{1}{2}} & \frac{1}{\sigma^2} R_{33}R_{33}^\dagger \end{bmatrix} - M_{\varepsilon,N} \\ = & \begin{bmatrix} F^{-\frac{1}{2}} V F^{-\frac{1}{2}\dagger} + N\mathbf{I}_{k_1+k_i} & F^{-\frac{1}{2}} K^\dagger \frac{1}{\sigma} \\ \frac{1}{\sigma} K F^{-\frac{1}{2}} & \varepsilon \mathbf{I}_{k_0} \end{bmatrix} \end{aligned}$$

Note that  $\begin{bmatrix} F^{-\frac{1}{2}} V F^{-\frac{1}{2}\dagger} + N\mathbf{I}_{k_1+k_i} & F^{-\frac{1}{2}} K^\dagger \frac{1}{\sigma} \\ \frac{1}{\sigma} K F^{-\frac{1}{2}} & \varepsilon \mathbf{I}_{k_0} \end{bmatrix}$  is congruent to

$$\begin{bmatrix} F^{-\frac{1}{2}} V F^{-\frac{1}{2}\dagger} + N\mathbf{I}_{k_1+k_i} - \frac{1}{\varepsilon \sigma^2} F^{-\frac{1}{2}} K^\dagger K F^{-\frac{1}{2}} & 0 \\ 0 & \varepsilon \mathbf{I}_{k_0} \end{bmatrix}.$$

Clearly  $F^{-\frac{1}{2}}VF^{-\frac{1}{2}\dagger} + NI_{k_1+k_i} - \frac{1}{\varepsilon\sigma^2}F^{-\frac{1}{2}}K^\dagger KF^{-\frac{1}{2}}$  is positive definite when  $N$  is sufficiently large. Hence, when  $N$  is large enough, we obtain

$$\begin{bmatrix} F^{-\frac{1}{2}}VF^{-\frac{1}{2}\dagger} & F^{-\frac{1}{2}}K^\dagger\frac{1}{\sigma} \\ \frac{1}{\sigma}KF^{-\frac{1}{2}} & \frac{1}{\sigma^2}R_{33}R_{33}^\dagger \end{bmatrix} \succ M_{\varepsilon,N}.$$

Using Corollary 4.3.3 of [55], we conclude that the eigenvalues of  $\begin{bmatrix} F^{-\frac{1}{2}}VF^{-\frac{1}{2}\dagger} & F^{-\frac{1}{2}}K^\dagger\frac{1}{\sigma} \\ \frac{1}{\sigma}KF^{-\frac{1}{2}} & \frac{1}{\sigma^2}R_{33}R_{33}^\dagger \end{bmatrix}$  are greater than those of  $M_{\varepsilon,N}$  if sorted. From Proposition 1, we know that  $\begin{bmatrix} F^{-\frac{1}{2}}VF^{-\frac{1}{2}\dagger} + NI_{k_1+k_i} & F^{-\frac{1}{2}}K^\dagger\frac{1}{\sigma} \\ \frac{1}{\sigma}KF^{-\frac{1}{2}} & \varepsilon I_{k_0} \end{bmatrix}$  has exactly  $k_0$  positive eigenvalues, which are the only eigenvalues that could be greater than  $\lambda\left(\frac{1}{\sigma^2}R_{33}R_{33}^\dagger\right) - \varepsilon$ . Since  $\varepsilon$  is arbitrary, we finally conclude that the positive eigenvalues of  $\Sigma_1^{-1}G$  are greater than those of  $\frac{1}{\sigma^2}R_{33}R_{33}^\dagger$  if sorted in the same way.

For the second claim, we need some notations first. For any pair of symmetric (or Hermitian) matrices  $P$  and  $Q$ ,  $P \prec Q$  means that  $Q - P$  is positive definite and  $P \preceq Q$  means  $Q - P$  is nonnegative definite. Note that if  $P$  and  $Q$  are positive definite, then from Corollary 7.7.4 of [55]  $P \preceq Q$  if and only if  $Q^{-1} \preceq P^{-1}$ ; if  $P \preceq Q$  then the eigenvalues of  $P$  and  $Q$  satisfy  $\lambda_k(P) \leq \lambda_k(Q)$  if they are sorted in the same order (decreasing or increasing). Therefore,  $\sigma^2 I_M + A_{S_0 S_1} A_{S_0 S_1}^\dagger \preceq \sigma^2 I_M + A_{S_1} A_{S_1}^\dagger = \Sigma_1$  yields

$$\Sigma_0^{1/2} \Sigma_1^{-1} \Sigma_0^{1/2} \preceq \Sigma_0^{1/2} \left( \sigma^2 I_M + A_{S_0 S_1} A_{S_0 S_1}^\dagger \right)^{-1} \Sigma_0^{1/2}.$$

Since we are interested only in the eigenvalues, a cyclic permutation in the matrix product on the previous inequality's right-hand side gives us

$$\begin{aligned} & \left( \sigma^2 I_M + A_{S_0 S_1} A_{S_0 S_1}^\dagger \right)^{-\frac{1}{2}} \Sigma_0 \left( \sigma^2 I_M + A_{S_0 S_1} A_{S_0 S_1}^\dagger \right)^{-\frac{1}{2}} \\ &= I_M + \left( \sigma^2 I_M + A_{S_0 S_1} A_{S_0 S_1}^\dagger \right)^{-\frac{1}{2}} A_{S_0 \setminus S_1} \\ & \quad \times A_{S_0 \setminus S_1}^\dagger \left( \sigma^2 I_M + A_{S_0 S_1} A_{S_0 S_1}^\dagger \right)^{-\frac{1}{2}} \\ &\triangleq I_M + P. \end{aligned}$$

Until now we have shown that the sorted eigenvalues of  $H$  are less than the corresponding eigenvalues of  $I_M + P$ . In particular, the eigenvalues of  $H$  that are greater than 1 are upper bounded by 1 plus the eigenvalues of  $P$  if they are both sorted ascendantly. Since from the

definition of eigenvalues, the non-zero eigenvalues of  $AB$  and  $BA$  are the same for any matrices  $A$  and  $B$ , we conclude that the eigenvalues of  $H$  that are greater than 1 are less than those of

$$\begin{aligned} & \mathbf{I}_{k_0} + A_{S_0 \setminus S_1}^\dagger \left( \sigma^2 \mathbf{I}_M + A_{S_0 S_1} A_{S_0 S_1}^\dagger \right)^{-1} A_{S_0 \setminus S_1} \\ & \preceq \mathbf{I}_{k_0} + \frac{1}{\sigma^2} A_{S_0 \setminus S_1}^\dagger A_{S_0 \setminus S_1}. \end{aligned}$$

Therefore, the conclusion of the second part of the theorem holds. We comment here that usually it is not true that  $H \preceq \mathbf{I}_M + \frac{1}{\sigma^2} A_{S_0 \setminus S_1} A_{S_0 \setminus S_1}^\dagger$ . Only the inequality on eigenvalues hold.

■

## APPENDIX C

### PROOF OF LEMMA 1

For arbitrary fixed supports  $S_i, S_j$ , we have

$$\begin{aligned} \bar{\lambda}_{S_i, S_j} & \geq 1 + \frac{1}{2\kappa\sigma^2} \left( \prod_{l=1}^{k_d} 2\kappa r_{ll}^2 \right)^{1/k_d} \\ & \geq \frac{1}{2\kappa\sigma^2} \min_{1 \leq l \leq k_d} q_l, \end{aligned}$$

where  $2\kappa r_{ll}^2 \sim \chi_{2\kappa(M-K-l+1)}^2$  can be written as a sum of  $2\kappa(M-K-l+1)$  independent squared standard Gaussian random variables and  $q_l \sim \chi_{2\kappa(M-2K)}^2$  is obtained by dropping  $K-l+1$  of them. Therefore, using the union bound we obtain

$$\begin{aligned} & \mathbb{P} \{ \bar{\lambda}_{S_i, S_j} \leq \gamma \} \\ & \leq \mathbb{P} \left\{ \frac{1}{2\kappa\sigma^2} \min_{1 \leq l \leq k_d} q_l \leq \gamma \right\} \\ & \leq \mathbb{P} \left\{ \bigcup_{1 \leq l \leq k_d} [q_l \leq 2\kappa\sigma^2\gamma] \right\} \\ & \leq k_d \mathbb{P} \{ q_l \leq 2\kappa\sigma^2\gamma \}. \end{aligned}$$

Since  $\gamma = \frac{1}{3} \frac{M-2K}{\sigma^2}$  implies that  $2\kappa\sigma^2\gamma = \frac{2\kappa}{3} (M-2K) < 2\kappa(M-2K) - 2$ , the mode of  $\chi_{\kappa(M-2K)}^2$ , when  $M-2K$  is sufficiently large, we have

$$\begin{aligned} & \mathbb{P} \{ q_l \leq 2\kappa\sigma^2\gamma \} \\ & = \int_0^{2\kappa\sigma^2\gamma} \frac{(1/2)^{\kappa(M-2K)}}{\Gamma(\kappa(M-2K))} x^{\kappa(M-2K)-1} e^{-x/2} dx \\ & \leq \frac{[\kappa\sigma^2\gamma]^{\kappa(M-2K)}}{\Gamma(\kappa(M-2K))} e^{-\kappa\sigma^2\gamma}. \end{aligned}$$

The inequality  $\ln \Gamma(z) \geq (z - \frac{1}{2}) \ln z - z$  says that when  $M - 2K$  is large enough,

$$\begin{aligned}
& \mathbb{P} \{q_l \leq 2\kappa\sigma^2\gamma\} \\
& \leq \exp \{ \kappa(M - 2K) \log(\kappa\sigma^2\gamma) - \kappa\sigma^2\gamma \\
& \quad - \left[ \kappa(M - 2K) - \frac{1}{2} \right] \log[\kappa(M - 2K)] \\
& \quad + \kappa(M - 2K) \} \\
& \leq \exp \{-c(M - 2K)\},
\end{aligned}$$

where  $c < \kappa(\log 3 - 1)$ . Therefore, we have

$$\mathbb{P} \{ \bar{\lambda}_{S_i, S_j} \leq \gamma \} \leq K \exp \{-c(M - 2K)\}. \quad \blacksquare$$

## REFERENCES

- [1] E. Candès, J. Romberg, and T. Tao, “Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information,” *IEEE Trans. Inf. Theory*, vol. 52, no. 2, pp. 489–509, Feb. 2006.
- [2] D. Donoho, “Compressed sensing,” *IEEE Trans. Inf. Theory*, vol. 52, no. 4, pp. 1289–1306, Apr. 2006.
- [3] E. Candès and T. Tao, “Decoding by linear programming,” *IEEE Trans. Inf. Theory*, vol. 51, no. 12, pp. 4203–4215, Dec. 2005.
- [4] E. Candès and M. Wakin, “An introduction to compressive sampling,” *IEEE Signal Process. Mag.*, vol. 25, no. 2, pp. 21–30, Mar. 2008.
- [5] R. Baraniuk, “Compressive sensing [lecture notes],” *IEEE Signal Process. Mag.*, vol. 24, no. 4, pp. 118–121, Jul. 2007.
- [6] D. Malioutov, M. Cetin, and A. Willsky, “A sparse signal reconstruction perspective for source localization with sensor arrays,” *IEEE Trans. Signal Process.*, vol. 53, no. 8, pp. 3010–3022, Aug. 2005.
- [7] D. Model and M. Zibulevsky, “Signal reconstruction in sensor arrays using sparse representations,” *Signal Processing*, vol. 86, no. 3, pp. 624–638, Mar. 2006.
- [8] V. Cevher, M. Duarte, and R. G. Baraniuk, “Distributed target localization via spatial sparsity,” in *European Signal Processing Conference (EUSIPCO 2008)*, Lausanne, Switzerland, Aug. 2008.
- [9] V. Cevher, P. Indyk, C. Hegde, and R. G. Baraniuk, “Recovery of clustered sparse signals from compressive measurements,” in *Int. Conf. Sampling Theory and Applications (SAMPTA 2009)*, Marseille, France, May 2009, pp. 18–22.
- [10] S. S. Chen, D. L. Donoho, and M. A. Saunders, “Atomic decomposition by basis pursuit,” *SIAM J. Sci. Comp.*, vol. 20, no. 1, pp. 33–61, 1998.
- [11] S. Mallat, *A Wavelet Tour of Signal Processing*. San Diego, CA: Academic, 1998.
- [12] I. Gorodnitsky and B. Rao, “Sparse signal reconstruction from limited data using FOCUSS: A re-weighted minimum norm algorithm,” *IEEE Trans. Signal Process.*, vol. 45, no. 3, pp. 600–616, Mar. 1997.
- [13] E. Candès and T. Tao, “The Dantzig selector: Statistical estimation when  $p$  is much larger than  $n$ ,” *Ann. Statist.*, vol. 35, pp. 2313–2351, 2007.
- [14] J. Tropp and A. Gilbert, “Signal recovery from random measurements via orthogonal matching pursuit,” *IEEE Trans. Inf. Theory*, vol. 53, no. 12, pp. 4655–4666, Dec. 2007.

- [15] W. Dai and O. Milenkovic, "Subspace pursuit for compressive sensing signal reconstruction," *IEEE Trans. Inf. Theory*, vol. 55, no. 5, pp. 2230–2249, May 2009.
- [16] D. Needell and J. Tropp, "CoSaMP: Iterative signal recovery from incomplete and inaccurate samples," *Appl. Comput. Harmonic Anal.*, vol. 26, no. 3, pp. 301–321, May 2009.
- [17] P. Borgnat and P. Flandrin, "Time-frequency localization from sparsity constraints," in *Proc. IEEE Int. Conf. Acoustics, Speech and Signal Processing (ICASSP 2008)*, Las Vegas, NV, Apr. 2008, pp. 3785–3788.
- [18] Z. Tian and G. Giannakis, "Compressed sensing for wideband cognitive radios," in *Proc. IEEE Int. Conf. Acoustics, Speech and Signal Processing (ICASSP 2007)*, Honolulu, HI, Apr. 2007, pp. IV–1357–IV–1360.
- [19] M. A. Sheikh, S. Sarvotham, O. Milenkovic, and R. G. Baraniuk, "DNA array decoding from nonlinear measurements by belief propagation," in *Proc. IEEE Workshop Statistical Signal Processing (SSP 2007)*, Madison, WI, Aug. 2007, pp. 215–219.
- [20] H. Vikalo, F. Parvaresh, and B. Hassibi, "On recovery of sparse signals in compressed DNA microarrays," in *Proc. Asilomar Conf. Signals, Systems and Computers (ACSSC 2007)*, Pacific Grove, CA, Nov. 2007, pp. 693–697.
- [21] F. Parvaresh, H. Vikalo, S. Misra, and B. Hassibi, "Recovering sparse signals using sparse measurement matrices in compressed DNA microarrays," *IEEE J. Sel. Topics Signal Processing*, vol. 2, no. 3, pp. 275–285, Jun. 2008.
- [22] H. Vikalo, F. Parvaresh, S. Misra, and B. Hassibi, "Sparse measurements, compressed sampling, and DNA microarrays," in *Proc. IEEE Int. Conf. Acoustics, Speech and Signal Processing (ICASSP 2008)*, Las Vegas, NV, Apr. 2008, pp. 581–584.
- [23] R. Baraniuk and P. Steeghs, "Compressive radar imaging," in *IEEE Radar Conference*, Apr. 2007, pp. 128–133.
- [24] M. Herman and T. Strohmer, "Compressed sensing radar," in *IEEE Radar Conference*, May 2008, pp. 1–6.
- [25] M. A. Herman and T. Strohmer, "High-resolution radar via compressed sensing," *IEEE Trans. Signal Process.*, vol. 57, no. 6, pp. 2275–2284, Jun. 2009.
- [26] E. G. Larsson and Y. Selen, "Linear regression with a sparse parameter vector," *IEEE Trans. Signal Process.*, vol. 55, no. 2, pp. 451–460, Feb. 2007.
- [27] D. H. Johnson and D. E. Dudgeon, *Array Signal Processing: Concepts and Techniques*. Englewood Cliffs, NJ: Prentice Hall, 1993.
- [28] J. Capon, "High-resolution frequency-wavenumber spectrum analysis," *Proceedings of the IEEE*, vol. 57, no. 8, pp. 1408–1418, Aug. 1969.
- [29] R. Schmidt, "Multiple emitter location and signal parameter estimation," *IEEE Trans. Antennas Propag.*, vol. 34, no. 3, pp. 276–280, Mar. 1986.
- [30] G. Bienvenu and L. Kopp, "Adaptivity to background noise spatial coherence for high resolution passive methods," in *Proc. IEEE Int. Conf. Acoustics, Speech and Signal Processing (ICASSP 1980)*, vol. 5, Denver, CO, Apr. 1980, pp. 307–310.
- [31] P. Stoica and A. Nehorai, "MUSIC, maximum likelihood, and Cramer-Rao bound: Further results and comparisons," *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 38, no. 12, pp. 2140–2150, Dec. 1990.
- [32] M. Duarte, S. Sarvotham, D. Baron, M. Wakin, and R. Baraniuk, "Distributed compressed sensing of jointly sparse signals," in *Proc. Asilomar Conf. Signals, Systems and Computers (ACSSC 2005)*, Pacific Grove, CA, Nov. 2005, pp. 1537–1541.
- [33] M. Duarte, M. Wakin, D. Baron, and R. Baraniuk, "Universal distributed sensing via random projections," in *Int. Conf. Information Processing in Sensor Networks (IPSN 2006)*, Nashville, TN, Apr. 2006, pp. 177–185.
- [34] M. Fornasier and H. Rauhut, "Recovery algorithms for vector-valued data with joint sparsity constraints," *SIAM J. Numer. Anal.*, vol. 46, no. 2, pp. 577–613, 2008.
- [35] M. Mishali and Y. Eldar, "The continuous joint sparsity prior for sparse representations: Theory and applications," in *IEEE*

- Int. Workshop Computational Advances in Multi-Sensor Adaptive Processing (CAMPSPAP 2007)*, St. Thomas, U.S. Virgin Islands, Dec. 2007, pp. 125–128.
- [36] S. Cotter, B. Rao, K. Engan, and K. Kreutz-Delgado, “Sparse solutions to linear inverse problems with multiple measurement vectors,” *IEEE Trans. Signal Process.*, vol. 53, no. 7, pp. 2477–2488, Jul. 2005.
  - [37] J. Chen and X. Huo, “Sparse representations for multiple measurement vectors (MMV) in an over-complete dictionary,” in *Proc. IEEE Int. Conf. Acoustics, Speech and Signal Processing (ICASSP 2005)*, Philadelphia, PA, Mar. 2005, pp. 257–260.
  - [38] —, “Theoretical results on sparse representations of multiple-measurement vectors,” *IEEE Trans. Signal Process.*, vol. 54, no. 12, pp. 4634–4643, Dec. 2006.
  - [39] J. A. Tropp, A. C. Gilbert, and M. J. Strauss, “Algorithms for simultaneous sparse approximation. Part I: Greedy pursuit,” *Signal Processing*, vol. 86, no. 3, pp. 572–588, 2006.
  - [40] J. A. Tropp, “Algorithms for simultaneous sparse approximation. Part II: Convex relaxation,” *Signal Processing*, vol. 86, no. 3, pp. 589 – 602, 2006.
  - [41] H. L. Van Trees, *Detection, Estimation, and Modulation Theory, Part I*. New York: Wiley-Interscience, 2001.
  - [42] A. K. Gupta and D. K. Nagar, *Matrix Variate Distributions*. Boca Raton, FL: Chapman & Hall/CRC, 1999.
  - [43] Y. S. Chow and H. Teicher, *Probability Theory: Independence, Interchangeability, Martingales*. New York: Springer-Verlag, 1988.
  - [44] R. Baraniuk, V. Cevher, M. Duarte, and C. Hegde, “Model-based compressive sensing,” submitted for publication.
  - [45] M. Wainwright, “Information-theoretic bounds on sparsity recovery in the high-dimensional and noisy setting,” in *IEEE Int. Symp. Information Theory (ISIT 2007)*, Nice, France, Jun. 2007, pp. 961–965.
  - [46] D. S. Bernstein, *Matrix Mathematics: Theory, Facts, and Formulas with Application to Linear Systems Theory*. Princeton, NJ: Princeton University Press, 2005.
  - [47] T. M. Cover and J. A. Thomas, *Elements of Information Theory*, 2nd ed. Hoboken, NJ: Wiley-Interscience, 2006.
  - [48] S. Kullback, *Information Theory and Statistics*. New York: Dover Publications, 1997.
  - [49] Y. G. Yatracos, “A lower bound on the error in nonparametric regression type problems,” *Ann. Statist.*, vol. 16, no. 3, pp. 1180–1187, 1988.
  - [50] L. Birge, “Approximation dans les espaces métriques et théorie de l’estimation,” *Probability Theory and Related Fields*, vol. 65, no. 2, pp. 181–237, Jun. 1983.
  - [51] A. Dogandzic and A. Nehorai, “Space-time fading channel estimation and symbol detection in unknown spatially correlated noise,” *IEEE Trans. Signal Process.*, vol. 50, no. 3, pp. 457–474, Mar. 2002.
  - [52] H. L. V. Trees, *Optimum Array Processing*. New York: Wiley-Interscience, 2002.
  - [53] T. W. Anderson, *An Introduction to Multivariate Statistical Analysis*, 3rd ed. Hoboken, NJ: Wiley-Interscience, 2003.
  - [54] R. A. DeVore, “Deterministic constructions of compressed sensing matrices,” *J. Complexity*, vol. 23, no. 4-6, pp. 918–925, 2007.
  - [55] R. A. Horn and C. R. Johnson, *Matrix Analysis*. New York: Cambridge University Press, 1990.
  - [56] I. Gohberg, P. Lancaster, and L. Rodman, *Indefinite Linear Algebra and Applications*. Basel, Switzerland: Birkhäuser, 2005.
  - [57] G. H. Golub and C. F. V. Loan, *Matrix Computations*, 3rd ed. Baltimore, MD: Johns Hopkins University Press, 1996.

PLACE  
PHOTO  
HERE

**Gongguo Tang** earned his B.Sc. degree in Mathematics from the Shandong University, China in 2003, and the M.Sc. degree in System Science from Chinese Academy of Sciences, China, in 2006.

Currently, he is a Ph.D. candidate of the Department of Electrical and Systems Engineering, Washington University, under the guidance of Dr. Arye Nehorai. His research interests are in the area of compressive sensing, statistical signal processing, detection and estimation, and their applications.

PLACE  
PHOTO  
HERE

**Arye Nehorai** (S'80-M'83-SM'90-F'94) earned his B.Sc. and M.Sc. degrees in electrical engineering from the Technion-Israel Institute of Technology, Haifa, Israel, and the Ph.D. degree in electrical engineering from Stanford University, Stanford, CA.

From 1985 to 1995, he was a Faculty Member with the Department of Electrical Engineering at Yale University. In 1995, he became a Full Professor in the Department of Electrical Engineering and Computer Science at The University of Illinois at Chicago (UIC). From 2000 to 2001, he was Chair of the Electrical and Computer Engineering (ECE) Division, which then became a new department. In 2001, he was named University Scholar of the University of Illinois. In 2006, he became Chairman of the Department of Electrical and Systems Engineering at Washington University in St. Louis. He is the inaugural holder of the Eugene and Martha Lohman Professorship and the Director of the Center for Sensor Signal and Information Processing (CSSIP) at WUSTL since 2006.

Dr. Nehorai was Editor-in-Chief of the IEEE TRANSACTIONS ON SIGNAL PROCESSING from 2000 to 2002. From 2003 to 2005, he was Vice President (Publications) of the IEEE Signal Processing Society (SPS), Chair of the Publications Board, member of the Board of Governors, and member of the Executive Committee of this Society. From 2003 to 2006, he was the founding editor of the special columns on Leadership Reflections in the IEEE Signal Processing Magazine. He was co-recipient of the IEEE SPS 1989 Senior Award for Best Paper with P. Stoica, coauthor of the 2003 Young Author Best Paper Award, and co-recipient of the 2004 Magazine Paper Award with A. Dogandzic. He was elected Distinguished Lecturer of the IEEE SPS for the term 2004 to 2005 and received the 2006 IEEE SPS Technical Achievement Award. He is the Principal Investigator of the new multidisciplinary university research initiative (MURI) project entitled Adaptive Waveform Diversity for Full Spectral Dominance. He has been a Fellow of the Royal Statistical Society since 1996.