

PERFORMANCE BOUNDS FOR SPARSE PARAMETRIC COVARIANCE ESTIMATION IN GAUSSIAN MODELS

Alexander Jung¹, Sebastian Schmutzhard², Franz Hlawatsch¹, and Alfred O. Hero III³

¹Institute of Telecommunications, Vienna University of Technology, Austria; {ajung, fhlawats}@nt.tuwien.ac.at

²NuHAG, Faculty of Mathematics, University of Vienna, Austria; sebastian.schmutzhard@univie.ac.at

³Department of Electrical Engineering and Computer Science, University of Michigan, Ann Arbor, MI, USA; hero@eecs.umich.edu

ABSTRACT

We consider estimation of a sparse parameter vector that determines the covariance matrix of a Gaussian random vector via a sparse expansion into known “basis matrices.” Using the theory of reproducing kernel Hilbert spaces, we derive lower bounds on the variance of estimators with a given mean function. This includes unbiased estimation as a special case. We also present a numerical comparison of our lower bounds with the variance of two standard estimators (hard-thresholding estimator and maximum likelihood estimator).

Index Terms—Sparsity, sparse covariance estimation, variance bound, reproducing kernel Hilbert space, RKHS.

1. INTRODUCTION

We consider a Gaussian signal vector $\mathbf{s} \in \mathbb{R}^M$, $\mathbf{s} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{C})$ embedded in white Gaussian noise $\mathbf{n} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$. The observed vector is

$$\mathbf{y} = \mathbf{s} + \mathbf{n}, \quad (1)$$

where \mathbf{s} and \mathbf{n} are independent and the signal mean $\boldsymbol{\mu}$ and noise variance σ^2 are known. In what follows, we assume $\boldsymbol{\mu} = \mathbf{0}$ since a nonzero $\boldsymbol{\mu}$ can always be subtracted from \mathbf{s} . The signal covariance matrix \mathbf{C} is unknown; we will parameterize it according to

$$\mathbf{C} = \mathbf{C}(\mathbf{x}) \triangleq \sum_{k \in [N]} x_k \mathbf{C}_k, \quad (2)$$

with unknown nonrandom coefficients $x_k \geq 0$, known positive semidefinite “basis matrices” $\mathbf{C}_k \in \mathbb{R}^{M \times M}$, and $[N] \triangleq \{1, \dots, N\}$. Thus, estimation of the signal covariance matrix \mathbf{C} reduces to estimation of the coefficient vector $\mathbf{x} \triangleq (x_1, \dots, x_N)^T \in \mathbb{R}_+^N$.

Our central assumption is that \mathbf{x} is *S-sparse*, i.e., at most S coefficients x_k are nonzero. We can formulate this as

$$\mathbf{x} \in \mathcal{X}_{S,+} \triangleq \{\mathbf{x}' \in \mathbb{R}_+^N \mid \|\mathbf{x}'\|_0 \leq S\}. \quad (3)$$

The sparsity degree S is supposed known; however, the set of positions of the nonzero entries of \mathbf{x} (denoted by $\text{supp}(\mathbf{x})$; note that $|\text{supp}(\mathbf{x})| = \|\mathbf{x}\|_0 \leq S$) is unknown. Typically, $S \ll N$. We will refer to (1)–(3) as the *sparse parametric covariance model* (SPCM). The SPCM and estimation of \mathbf{x} are relevant, e.g., in time-frequency (TF) analysis [1, 2], where the basis matrices \mathbf{C}_k correspond to disjoint TF regions and x_k represents the mean signal power in the k th TF region. An application is cognitive radio scene analysis [3].

The problem we will study is estimation of $\mathbf{z} \triangleq \mathbf{g}(\mathbf{x}) \in \mathbb{R}^K$ from \mathbf{y} , where $\mathbf{g}(\cdot)$ is a known function. This includes estimation

This work was supported by the FWF under Grants S10602 and S10603 within the National Research Network SISE, and by the WWTF under Grant MA 07-004 (SPORTS).

of \mathbf{x} and, less trivially, of a linear combination of the x_k . In the TF application mentioned above, the latter case corresponds to a linear combination of the mean signal powers in the various TF regions.

In this paper, building on [4, 5], we use the theory of *reproducing kernel Hilbert spaces* (RKHS) to derive lower bounds on the variance of estimators of \mathbf{z} for an important special case of the SPCM that we term the *sparse diagonalizable parametric covariance model* (SDPCM). The estimators are required to have a prescribed differentiable mean function; this includes the case of unbiased estimation. They are allowed to exploit the known sparsity of \mathbf{x} . The RKHS framework has been previously proposed for a fundamentally different problem of sparsity-exploiting estimation in [6].

Sparsity-exploiting estimation of \mathbf{C} and of \mathbf{C}^{-1} was considered recently in [7] and in [8], respectively. In both cases, the sparsity assumption was placed on \mathbf{C}^{-1} , which corresponds to a sparse graphical model for \mathbf{s} . Our SPCM approach (2), (3) is clearly different: while the coefficient vector \mathbf{x} is assumed sparse, the matrices \mathbf{C} or \mathbf{C}^{-1} need not be sparse.

This paper is organized as follows. In Section 2, we review minimum-variance estimation and the RKHS framework. In Section 3, we use RKHS theory to derive lower variance bounds for the SDPCM. The special case of unbiased estimation is considered in Section 4. Finally, Section 5 presents a numerical comparison of our bounds with the variance of two established estimation schemes.

2. RKHS FORMULATION OF MINIMUM-VARIANCE ESTIMATION

2.1. Minimum-Variance Estimation

The performance of an estimator $\hat{\mathbf{z}}(\mathbf{y})$ of $\mathbf{z} = \mathbf{g}(\mathbf{x})$ can be quantified by the mean squared error (MSE) $\varepsilon(\hat{\mathbf{z}}(\cdot); \mathbf{x}) \triangleq \mathbb{E}_{\mathbf{x}} \{\|\hat{\mathbf{z}}(\mathbf{y}) - \mathbf{z}\|_2^2\}$, where the notation $\mathbb{E}_{\mathbf{x}}\{\cdot\}$ indicates that the expectation is taken with respect to the probability density function $f(\mathbf{y}; \mathbf{x})$ parameterized by \mathbf{x} . According to our assumptions in Section 1,

$$f(\mathbf{y}; \mathbf{x}) = \frac{\exp(-\frac{1}{2} \mathbf{y}^T \tilde{\mathbf{C}}^{-1}(\mathbf{x}) \mathbf{y})}{[(2\pi)^M \det\{\tilde{\mathbf{C}}(\mathbf{x})\}]^{1/2}}, \quad \text{with } \tilde{\mathbf{C}}(\mathbf{x}) \triangleq \mathbf{C}(\mathbf{x}) + \sigma^2 \mathbf{I}. \quad (4)$$

Let z_k and $\hat{z}_k(\mathbf{y})$ denote the k th entries of \mathbf{z} and $\hat{\mathbf{z}}(\mathbf{y})$, respectively. We have $\varepsilon(\hat{\mathbf{z}}(\cdot); \mathbf{x}) = \sum_{k \in [K]} \varepsilon(\hat{z}_k(\cdot); \mathbf{x})$, where $\varepsilon(\hat{z}_k(\cdot); \mathbf{x}) \triangleq \mathbb{E}_{\mathbf{x}} \{[\hat{z}_k(\mathbf{y}) - z_k]^2\}$ denotes the k th component MSE. For our scope, minimization of $\varepsilon(\hat{\mathbf{z}}(\cdot); \mathbf{x})$ with respect to $\hat{\mathbf{z}}(\cdot)$ is equivalent to separate minimization of each component MSE $\varepsilon(\hat{z}_k(\cdot); \mathbf{x})$ with respect to $\hat{z}_k(\cdot)$, for $k \in [K]$. We furthermore have

$$\varepsilon(\hat{\mathbf{z}}(\cdot); \mathbf{x}) = b^2(\hat{\mathbf{z}}(\cdot); \mathbf{x}) + v(\hat{\mathbf{z}}(\cdot); \mathbf{x}), \quad (5)$$

with the component bias $b(\hat{z}_k(\cdot); \mathbf{x}) \triangleq \mathbf{E}_{\mathbf{x}}\{\hat{z}_k(\mathbf{y})\} - z_k$ and the component variance $v(\hat{z}_k(\cdot); \mathbf{x}) \triangleq \mathbf{E}_{\mathbf{x}}\{[\hat{z}_k(\mathbf{y}) - \mathbf{E}_{\mathbf{x}}\{\hat{z}_k(\mathbf{y})\}]^2\}$. A common approach to defining a ‘‘locally optimal’’ estimator $\hat{z}_k(\cdot)$ is to require $b(\hat{z}_k(\cdot); \mathbf{x}) = c_k(\mathbf{x})$ for all $\mathbf{x} \in \mathcal{X}_{S,+}$, with a given bias function $c_k(\mathbf{x})$, and look for estimators that minimize the variance $v(\hat{z}_k(\cdot); \mathbf{x})$ at a given parameter vector $\mathbf{x} = \mathbf{x}_0 \in \mathcal{X}_{S,+}$. It follows from (5) that once the bias is fixed, minimizing $\varepsilon(\hat{z}_k(\cdot); \mathbf{x}_0)$ is equivalent to minimizing $v(\hat{z}_k(\cdot); \mathbf{x}_0)$. Furthermore, fixing the bias is equivalent to fixing the mean, i.e., requiring that $\mathbf{E}_{\mathbf{x}}\{\hat{z}_k(\mathbf{y})\} = \gamma_k(\mathbf{x})$ for all $\mathbf{x} \in \mathcal{X}_{S,+}$, where $\gamma_k(\mathbf{x}) \triangleq c_k(\mathbf{x}) + g_k(\mathbf{x})$.

In what follows, we consider a fixed component k and drop the subscript k for better readability. Furthermore, we consider a given mean function $\gamma(\mathbf{x})$ (short for $\gamma_k(\mathbf{x})$) and a given nominal parameter vector \mathbf{x}_0 . We are interested in the minimum variance at \mathbf{x}_0 achievable by estimators $\hat{z}(\cdot)$ (short for $\hat{z}_k(\cdot)$) that have mean function $\gamma(\mathbf{x})$ for all $\mathbf{x} \in \mathcal{X}_{S,+}$. In order to derive a lower bound on this achievable variance, let us consider some subset $\mathcal{D} \subseteq \mathcal{X}_{S,+}$. We denote by $\mathcal{B}_{\gamma}^{\mathcal{D}}(\mathbf{x}_0)$ the set of all scalar estimators $\hat{z}(\cdot)$ whose mean equals $\gamma(\mathbf{x})$ for all $\mathbf{x} \in \mathcal{D}$ (however, not necessarily for all $\mathbf{x} \in \mathcal{X}_{S,+}$) and whose variance at \mathbf{x}_0 is finite, i.e.,

$$\mathcal{B}_{\gamma}^{\mathcal{D}}(\mathbf{x}_0) \triangleq \{\hat{z}(\cdot) \mid \mathbf{E}_{\mathbf{x}}\{\hat{z}(\mathbf{y})\} = \gamma(\mathbf{x}) \ \forall \mathbf{x} \in \mathcal{D}, v(\hat{z}(\cdot); \mathbf{x}_0) < \infty\}.$$

If $\mathcal{B}_{\gamma}^{\mathcal{D}}(\mathbf{x}_0)$ is nonempty, we consider the minimum variance achievable at the given parameter vector \mathbf{x}_0 by estimators $\hat{z}(\cdot) \in \mathcal{B}_{\gamma}^{\mathcal{D}}(\mathbf{x}_0)$:

$$L_{\gamma}^{\mathcal{D}}(\mathbf{x}_0) \triangleq \min_{\hat{z}(\cdot) \in \mathcal{B}_{\gamma}^{\mathcal{D}}(\mathbf{x}_0)} v(\hat{z}(\cdot); \mathbf{x}_0). \quad (6)$$

The use of min (rather than inf) in (6) is justified by the fact that the existence of a finite minimum can always be guaranteed by a proper choice of \mathcal{D} ; a sufficient condition will be provided in Section 2.2.

Because $\mathcal{D} \subseteq \mathcal{X}_{S,+}$, $L_{\gamma}^{\mathcal{D}}(\mathbf{x}_0)$ is a lower bound on the variance at \mathbf{x}_0 of any estimator $\hat{z}(\cdot)$ whose mean is $\gamma(\mathbf{x})$ for all $\mathbf{x} \in \mathcal{X}_{S,+}$ (and not just for all $\mathbf{x} \in \mathcal{D}$), i.e.,

$$L_{\gamma}^{\mathcal{D}}(\mathbf{x}_0) \leq v(\hat{z}(\cdot); \mathbf{x}_0), \quad \text{for any } \hat{z}(\cdot) \text{ such that} \\ \mathbf{E}_{\mathbf{x}}\{\hat{z}(\mathbf{y})\} = \gamma(\mathbf{x}) \ \forall \mathbf{x} \in \mathcal{X}_{S,+}. \quad (7)$$

2.2. RKHS Formulation

An inner product of two real random variables $a = a(\mathbf{y})$, $b = b(\mathbf{y})$ can be defined as $\langle a, b \rangle_{\text{RV}} \triangleq \mathbf{E}_{\mathbf{x}_0}\{a(\mathbf{y})b(\mathbf{y})\}$, with induced norm $\|a\|_{\text{RV}} = \sqrt{\langle a, a \rangle_{\text{RV}}} = \sqrt{\mathbf{E}_{\mathbf{x}_0}\{a^2(\mathbf{y})\}}$. Note the dependence on \mathbf{x}_0 . One can show that (6) can be rewritten formally as the following constrained norm-minimization problem:

$$L_{\gamma}^{\mathcal{D}}(\mathbf{x}_0) = \min_{\hat{z}(\cdot)} \|\hat{z}\|_{\text{RV}}^2 - \gamma^2(\mathbf{x}_0) \\ \text{subject to } \langle \hat{z}, \rho_{\mathbf{x}} \rangle_{\text{RV}} = \gamma(\mathbf{x}) \ \forall \mathbf{x} \in \mathcal{D}, \quad (8)$$

where

$$\rho_{\mathbf{x}}(\mathbf{y}) \triangleq \frac{f(\mathbf{y}; \mathbf{x})}{f(\mathbf{y}; \mathbf{x}_0)}. \quad (9)$$

Furthermore, if $\mathcal{B}_{\gamma}^{\mathcal{D}}(\mathbf{x}_0)$ is nonempty, the existence of a finite minimum in (6), (8) can be guaranteed by choosing \mathcal{D} such that [4, 5]

$$\|\rho_{\mathbf{x}}\|_{\text{RV}}^2 \equiv \mathbf{E}_{\mathbf{x}_0}\{\rho_{\mathbf{x}}^2(\mathbf{y})\} < \infty \quad \forall \mathbf{x} \in \mathcal{D}. \quad (10)$$

According to [4], the solutions of (8) can be described using an RKHS $\mathcal{H}(R)$ with kernel $R(\mathbf{x}_1, \mathbf{x}_2) : \mathcal{D} \times \mathcal{D} \rightarrow \mathbb{R}$ given by

$$R(\mathbf{x}_1, \mathbf{x}_2) = \langle \rho_{\mathbf{x}_1}, \rho_{\mathbf{x}_2} \rangle_{\text{RV}} = \mathbf{E}_{\mathbf{x}_0}\{\rho_{\mathbf{x}_1}(\mathbf{y})\rho_{\mathbf{x}_2}(\mathbf{y})\}. \quad (11)$$

Note that $R(\mathbf{x}_1, \mathbf{x}_2)$ and $\mathcal{H}(R)$ depend on \mathbf{x}_0 . Inserting (9) and (4) into (11) yields the expression

$$R(\mathbf{x}_1, \mathbf{x}_2) = [\det\{\tilde{\mathbf{C}}(\mathbf{x}_0)\}]^{1/2} [\det\{\tilde{\mathbf{C}}(\mathbf{x}_1)\tilde{\mathbf{C}}(\mathbf{x}_2) \\ \cdot (\tilde{\mathbf{C}}^{-1}(\mathbf{x}_1) + \tilde{\mathbf{C}}^{-1}(\mathbf{x}_2) - \tilde{\mathbf{C}}^{-1}(\mathbf{x}_0))\}]^{-1/2}, \quad (12)$$

where as before $\tilde{\mathbf{C}}(\mathbf{x}) = \mathbf{C}(\mathbf{x}) + \sigma^2 \mathbf{I}$. The RKHS $\mathcal{H}(R)$ is a Hilbert space of functions $f(\cdot) : \mathcal{D} \rightarrow \mathbb{R}$ that is defined as the closure of the linear span of the set of functions $\{f_{\mathbf{x}'}(\mathbf{x}) = R(\mathbf{x}, \mathbf{x}')\}_{\mathbf{x}' \in \mathcal{D}}$. This closure is taken with respect to the topology that is given by the inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}(R)}$ defined via the *reproducing property* [9]

$$\langle f(\cdot), R(\cdot, \mathbf{x}') \rangle_{\mathcal{H}(R)} = f(\mathbf{x}'), \quad f \in \mathcal{H}(R), \ \mathbf{x}' \in \mathcal{D}.$$

The induced norm is $\|f\|_{\mathcal{H}(R)} = \sqrt{\langle f, f \rangle_{\mathcal{H}(R)}}$.

It can be shown [4] that if \mathcal{D} satisfies (10), then $\gamma \in \mathcal{H}(R)$ is necessary and sufficient (i) for $\mathcal{B}_{\gamma}^{\mathcal{D}}(\mathbf{x}_0)$ to be nonempty and (ii) for the minimum value $L_{\gamma}^{\mathcal{D}}(\mathbf{x}_0)$ in (6), (8) to exist and be given by

$$L_{\gamma}^{\mathcal{D}}(\mathbf{x}_0) = \|\gamma\|_{\mathcal{H}(R)}^2 - \gamma^2(\mathbf{x}_0). \quad (13)$$

3. LOWER BOUNDS ON THE ESTIMATOR VARIANCE

According to (13), any lower bound on $\|\gamma\|_{\mathcal{H}(R)}^2$ entails a lower bound on $L_{\gamma}^{\mathcal{D}}(\mathbf{x}_0)$. For tractability, we hereafter assume that the basis matrices $\mathbf{C}_k \in \mathbb{R}^{M \times M}$ in (2) are projection matrices on orthogonal subspaces of \mathbb{R}^M . Thus, they can be written as

$$\mathbf{C}_k = \sum_{i \in [r_k]} \mathbf{u}_{m_{k,i}} \mathbf{u}_{m_{k,i}}^T, \quad k \in [N], \quad (14)$$

where $\{\mathbf{u}_m\}_{m \in [M]}$ is an orthonormal basis for \mathbb{R}^M and the sets $\mathcal{U}_k \triangleq \{\mathbf{u}_{m_{k,i}}\}_{i \in [r_k]}$ are disjoint, so that they span orthogonal subspaces of \mathbb{R}^M . We note that (2) and (14) correspond to a latent variable model $\mathbf{s} = \sum_{k \in [N]} \mathbf{s}_k$ with $\mathbf{s}_k = \sum_{i \in [r_k]} \xi_{m_{k,i}} \mathbf{u}_{m_{k,i}}$, where the $\xi_{m_{k,i}}$ are independent zero-mean Gaussian with variance x_k for all i , i.e., $\xi_{m_{k,i}} \sim \mathcal{N}(0, x_k)$. This is similar to the latent variable model used in probabilistic principal component analysis [10] except that our ‘‘factors’’ \mathbf{u}_m are fixed. With (14), the kernel expression in (12) simplifies to

$$R(\mathbf{x}_1, \mathbf{x}_2) = \frac{\prod_{k \in [N]} (x_{0,k} + \sigma^2)^{r_k}}{\prod_{k \in [N]} [(x_{0,k} + \sigma^2)^2 - (x_{1,k} - x_{0,k})(x_{2,k} - x_{0,k})]^{\frac{r_k}{2}}},$$

where, e.g., $x_{0,k}$ denotes the k th entry of \mathbf{x}_0 . We will refer to the SPCM with basis matrices \mathbf{C}_k of the form (14) as the SDPCM.¹ It can be shown that, within the SDPCM, a sufficient condition for (10)—and, thus, for the existence of a minimum in (6), (8)—is $x_k < 2x_{0,k} + \sigma^2$ for all $k \in [N]$. Therefore, we choose our domain as

$$\mathcal{D} = \{\mathbf{x} \in \mathcal{X}_{S,+} \mid x_k < 2x_{0,k} + \sigma^2 \ \forall k \in [N]\}.$$

Note that \mathcal{D} depends on \mathbf{x}_0 .

We will now derive a lower bound on $\|\gamma\|_{\mathcal{H}(R)}^2$ for the SDPCM. Let us assume for now that $\gamma \in \mathcal{H}(R)$. Consider L functions $w_l(\mathbf{x})$, $l \in [L]$, with $w_l(\cdot) : \mathcal{D} \rightarrow \mathbb{R}$ and $w_l \in \mathcal{H}(R)$, which are orthogonal, i.e., $\langle w_l, w_{l'} \rangle_{\mathcal{H}(R)} = 0$ if $l \neq l'$. Let \mathcal{W} denote the subspace of $\mathcal{H}(R)$

¹For the SDPCM, the covariance matrix $\mathbf{C}(\mathbf{x})$ can be *diagonalized* by a signal transformation $\mathbf{s}' = \mathbf{U}\mathbf{s}$, with a unitary matrix \mathbf{U} that does not depend on the true parameter vector \mathbf{x} .

spanned by the w_l , and $\mathbf{P}_{\mathcal{W}}$ the orthogonal projection operator on \mathcal{W} . Clearly, a lower bound on $\|\gamma\|_{\mathcal{H}(R)}^2$ is given by

$$\|\mathbf{P}_{\mathcal{W}}\gamma\|_{\mathcal{H}(R)}^2 \leq \|\gamma\|_{\mathcal{H}(R)}^2. \quad (15)$$

This lower bound can be expressed as

$$\|\mathbf{P}_{\mathcal{W}}\gamma\|_{\mathcal{H}(R)}^2 = \sum_{l \in [L]} \frac{|\langle \gamma, w_l \rangle_{\mathcal{H}(R)}|^2}{\|w_l\|_{\mathcal{H}(R)}^2}. \quad (16)$$

A convenient construction of functions $w_l(\mathbf{x})$ is via partial derivatives of $R(\mathbf{x}_1, \mathbf{x}_2)$ with respect to \mathbf{x}_2 [4]. Consider an index set $\mathcal{K} \subseteq [N]$ containing S indices, i.e., $|\mathcal{K}| = S$. Furthermore let $\mathbf{p}_l = (p_{l,1}, \dots, p_{l,N}) \in \mathbb{N}_0^N$ be L different multi-indices satisfying $\text{supp}(\mathbf{p}_l) \subseteq \mathcal{K}$. We then define

$$w_l(\mathbf{x}) \triangleq \left. \frac{\partial^{\mathbf{p}_l} R(\mathbf{x}, \mathbf{x}_2)}{\partial \mathbf{x}_2^{\mathbf{p}_l}} \right|_{\mathbf{x}_2 = \mathbf{x}_0^{\mathcal{K}}}, \quad l \in [L], \quad (17)$$

where $\frac{\partial^{\mathbf{p}_l} f(\mathbf{x})}{\partial \mathbf{x}^{\mathbf{p}_l}} \triangleq \left(\prod_{k \in [N]} \frac{\partial^{p_{l,k}}}{\partial x_k^{p_{l,k}}} \right) f(\mathbf{x})$ and $\mathbf{x}_0^{\mathcal{K}}$ is obtained from \mathbf{x}_0 by zeroing all entries except those whose indices are in \mathcal{K} . It can be verified that the functions w_l are orthogonal, i.e.,

$$\langle w_l, w_{l'} \rangle_{\mathcal{H}(R)} = q_l(\mathbf{x}_0) \delta_{l,l'}, \quad (18)$$

where $q_l(\mathbf{x}_0) = \left. \frac{\partial^{\mathbf{p}_l} \partial^{\mathbf{p}_l} R(\mathbf{x}_1, \mathbf{x}_2)}{\partial \mathbf{x}_1^{\mathbf{p}_l} \partial \mathbf{x}_2^{\mathbf{p}_l}} \right|_{\mathbf{x}_1 = \mathbf{x}_2 = \mathbf{x}_0^{\mathcal{K}}}$. Furthermore [4],

$$\langle f, w_l \rangle_{\mathcal{H}(R)} = \left. \frac{\partial^{\mathbf{p}_l} f(\mathbf{x})}{\partial \mathbf{x}^{\mathbf{p}_l}} \right|_{\mathbf{x} = \mathbf{x}_0^{\mathcal{K}}} \quad \text{for any } f \in \mathcal{H}(R). \quad (19)$$

Using (18) and (19) in (16), we obtain

$$\|\mathbf{P}_{\mathcal{W}}\gamma\|_{\mathcal{H}(R)}^2 = \sum_{l \in [L]} \frac{1}{q_l(\mathbf{x}_0)} \left| \left. \frac{\partial^{\mathbf{p}_l} \gamma(\mathbf{x})}{\partial \mathbf{x}^{\mathbf{p}_l}} \right|_{\mathbf{x} = \mathbf{x}_0^{\mathcal{K}}} \right|^2. \quad (20)$$

Finally, combining (7), (13), (15), and (20), we arrive at the following bound. (Hereafter, we again explicitly indicate the index k .)

Theorem 3.1. *For the SDPCM, let $\hat{z}_k(\cdot)$ be any estimator of $z_k = g_k(\mathbf{x})$ whose mean equals $\gamma_k(\mathbf{x})$ for all $\mathbf{x} \in \mathcal{X}_{S,+}$ and whose variance at a fixed $\mathbf{x}_0 \in \mathcal{X}_{S,+}$ is finite. Then, this variance satisfies*

$$v(\hat{z}_k(\cdot); \mathbf{x}_0) \geq \sum_{l \in [L]} \frac{1}{q_l(\mathbf{x}_0)} \left| \left. \frac{\partial^{\mathbf{p}_l} \gamma_k(\mathbf{x})}{\partial \mathbf{x}^{\mathbf{p}_l}} \right|_{\mathbf{x} = \mathbf{x}_0^{\mathcal{K}}} \right|^2 - \gamma_k^2(\mathbf{x}_0), \quad (21)$$

for any choice of L different $\mathbf{p}_l \in \mathbb{N}_0^N$ such that $\text{supp}(\mathbf{p}_l) \subseteq \mathcal{K}$, where $\mathcal{K} \subseteq [N]$ is an arbitrary set of S different indices. The lower bound (21) is achieved by an estimator $\hat{z}_k(\cdot)$ if and only if

$$\hat{z}_k(\mathbf{y}) = \sum_{l \in [L]} a_l \left. \frac{\partial^{\mathbf{p}_l} \rho_{\mathbf{x}}(\mathbf{y})}{\partial \mathbf{x}^{\mathbf{p}_l}} \right|_{\mathbf{x} = \mathbf{x}_0^{\mathcal{K}}}$$

with nonrandom coefficients $a_l \in \mathbb{R}$ and with $\rho_{\mathbf{x}}(\mathbf{y})$ defined in (9).

Note that the bound in (21) depends on $\gamma_k(\mathbf{x})$ only via a finite number of partial derivatives of $\gamma_k(\mathbf{x})$ at $\mathbf{x} = \mathbf{x}_0^{\mathcal{K}}$. Thus, it only depends on the local behavior of the prescribed mean or bias. We furthermore note that Theorem 3.1 does not mention the condition $\gamma_k \in \mathcal{H}(R)$ we used in its derivation. This is no problem because—as mentioned at the end of Section 2.2—if $\gamma_k \notin \mathcal{H}(R)$, there exists no estimator that has mean $\gamma_k(\mathbf{x})$ for all $\mathbf{x} \in \mathcal{D}$. Since $\mathcal{D} \subseteq \mathcal{X}_{S,+}$, it follows that no estimator has mean $\gamma_k(\mathbf{x})$ for all $\mathbf{x} \in \mathcal{X}_{S,+}$, either.

4. SPECIAL CASE: UNBIASED ESTIMATION

In this section, we evaluate the bound (21) for the important special case of unbiased estimation of \mathbf{x} , i.e., for $z_k = g_k(\mathbf{x}) = x_k$ and $c_k(\mathbf{x}) \equiv 0$ or equivalently $\gamma_k(\mathbf{x}) = x_k$. To obtain a simple expression, we use $L = 2$ and particular choices of \mathcal{K} and \mathbf{p}_l ($l = 1, 2$). Specifically, using $\mathcal{K} = \{k\} \cup \mathcal{L}$, where \mathcal{L} consists of the indices of the $S-1$ largest entries of the vector that is obtained from \mathbf{x}_0 by zeroing the k th entry, and $\mathbf{p}_1 = \mathbf{0}$ and $\mathbf{p}_2 = \mathbf{e}_k$, where \mathbf{e}_k denotes the k th column of the identity matrix, the following variance bound is obtained from Theorem 3.1.

Corollary 4.1. *For the SDPCM, let $\hat{x}_k(\cdot)$ be any estimator of x_k that is unbiased (i.e., $\gamma_k(\mathbf{x}) = x_k$) for all $\mathbf{x} \in \mathcal{X}_{S,+}$ and whose variance at a fixed $\mathbf{x}_0 \in \mathcal{X}_{S,+}$ is finite. Then, this variance satisfies*

$$v(\hat{x}_k(\cdot); \mathbf{x}_0) \geq \begin{cases} \frac{2}{r_k} (x_{0,k} + \sigma^2)^2, & k \in \text{supp}(\mathbf{x}_0) \\ \frac{2\sigma^4}{r_k} \left[1 - \frac{\xi_0^2}{(\xi_0 + \sigma^2)^2} \right]^{r_{j_0}/2}, & k \notin \text{supp}(\mathbf{x}_0), \end{cases} \quad (22)$$

where ξ_0 and j_0 denote, respectively, the value and index of the S -largest entry of \mathbf{x}_0 .

The lower bound (22) can be achieved at least in the following two cases: (i) if $k \in \text{supp}(\mathbf{x}_0)$, and (ii) for any $k \in [N]$ if $\|\mathbf{x}_0\|_0 < S$ (note that this latter condition implies $\xi_0 = 0$). In both cases, the estimator given by

$$\hat{x}_k(\mathbf{y}) = \beta_k(\mathbf{y}) - \sigma^2, \quad \text{with } \beta_k(\mathbf{y}) \triangleq \frac{1}{r_k} \sum_{i \in [r_k]} (\mathbf{u}_{m_{k,i}}^T \mathbf{y})^2, \quad (23)$$

is unbiased and its variance achieves the bound (22). This estimator does not use the sparsity information and does not depend on \mathbf{x}_0 .

Let us define a “signal-to-noise ratio” (SNR) quantity as $\text{SNR}(\mathbf{x}_0) \triangleq \xi_0 / \sigma^2$. For $\text{SNR}(\mathbf{x}_0) \ll 1$, the lower bound (22) is approximately $\frac{2}{r_k} (x_{0,k} + \sigma^2)^2$ for any k , which does not depend on S and moreover equals the variance of the unbiased estimator (23). Since that estimator does not exploit any sparsity information, Corollary 4.1 suggests that, in the low-SNR regime, *unbiased* estimators cannot exploit the prior information that \mathbf{x} is S -sparse. However, in the high-SNR regime ($\text{SNR}(\mathbf{x}_0) \rightarrow \infty$), (22) becomes $\frac{2}{r_k} (x_{0,k} + \sigma^2)^2$ for $k \in \text{supp}(\mathbf{x}_0)$ and 0 for $k \notin \text{supp}(\mathbf{x}_0)$, which can be shown to equal the variance of the oracle estimator that knows $\text{supp}(\mathbf{x}_0)$ (this oracle estimator yields $\hat{x}_k = x_{0,k} = 0$ for all $k \notin \text{supp}(\mathbf{x}_0)$). The transition of the lower bound (22) from the low-SNR regime to the high-SNR regime has a polynomial characteristic; it is thus much slower than the exponential transition of an analogous lower bound recently derived in [6] for the *sparse linear model*. This slow transition suggests that the optimal *unbiased* estimator for low SNR—which ignores the sparsity information—will also be a nearly optimal unbiased estimator over a relatively wide SNR range. This further suggests that, for *unbiased* covariance estimation based on the SDPCM, prior information of sparsity is not as helpful as for estimating the parameter vector of the sparse linear model [6].

In the special case where $S = 1$ and $\mathbf{x}_0 \neq \mathbf{0}$, ξ_0 and j_0 are, respectively, the value and index of the single nonzero entry of $\mathbf{x}_0 \in \mathcal{X}_{1,+}$. Consider the estimator $\hat{\mathbf{x}}^{(\mathbf{x}_0)}(\cdot)$ given componentwise by

$$\hat{x}_k^{(\mathbf{x}_0)}(\mathbf{y}) = \begin{cases} \beta_k(\mathbf{y}) - \sigma^2, & k = j_0 \\ \alpha(\mathbf{y}; \mathbf{x}_0) (\beta_k(\mathbf{y}) - \sigma^2), & k \neq j_0, \end{cases}$$

where $\alpha(\mathbf{y}; \mathbf{x}_0) \triangleq a(\mathbf{x}_0) \exp(-r_{j_0} b(\mathbf{x}_0) \beta_{j_0}(\mathbf{y}))$ with $a(\mathbf{x}_0) \triangleq \left[\frac{(\xi_0 + \sigma^2)^2 - \xi_0^2}{\sigma^2 (\xi_0 + \sigma^2)} \right]^{r_{j_0}/2}$ and $b(\mathbf{x}_0) \triangleq \frac{1}{2} \left(\frac{1}{\sigma^2} - \frac{1}{\xi_0 + \sigma^2} \right)$. One can verify that $\hat{\mathbf{x}}^{(\mathbf{x}_0)}(\cdot)$ is unbiased and has the minimum variance achievable by unbiased estimators at any $\mathbf{x}_0 \in \mathcal{X}_{1,+}$ with $\mathbf{x}_0 \neq \mathbf{0}$, since its variance achieves the bound (22). Note that this estimator depends explicitly on the assumed \mathbf{x}_0 , at which it achieves minimum variance; its performance may be poor when the true parameter vector \mathbf{x} is different from \mathbf{x}_0 .

5. NUMERICAL RESULTS

We compare the lower bound (21) for $\mathbf{g}(\mathbf{x}) = \mathbf{x}$ with the variance of two standard estimators. The first is an ad-hoc adaptation of the hard-thresholding (HT) estimator [11] to SDPCM-based covariance estimation. It is defined componentwise as (cf. (23))

$$\hat{\mathbf{x}}_{k,\text{HT}}(\mathbf{y}) \triangleq \frac{1}{r_k} \sum_{i \in [r_k]} \varphi_\tau^2(\mathbf{u}_{m_{k,i}}^T \mathbf{y}) - \sigma^2,$$

where $\varphi_\tau : \mathbb{R} \rightarrow \mathbb{R}$ denotes the hard-thresholding function with threshold $\tau \geq 0$, i.e., $\varphi_\tau(y)$ is y for $|y| \geq \tau$ and 0 else. The second standard method is the maximum likelihood (ML) estimator

$$\hat{\mathbf{x}}_{\text{ML}}(\mathbf{y}) \triangleq \arg \max_{\mathbf{x}' \in \mathcal{X}_{S,+}} f(\mathbf{y}; \mathbf{x}').$$

For the SDPCM, one can show that

$$\hat{\mathbf{x}}_{k,\text{ML}}(\mathbf{y}) = \begin{cases} \beta_k(\mathbf{y}) - \sigma^2, & k \in \mathcal{L}_1 \cap \mathcal{L}_2 \\ 0, & \text{else,} \end{cases}$$

where \mathcal{L}_1 consists of the S indices k for which $r_k [\beta_k(\mathbf{y})/\sigma^2 - \ln(\beta_k(\mathbf{y})/\sigma^2) - 1]$ (with $\ln = \log_e$) is largest, and \mathcal{L}_2 consists of all indices k for which $\beta_k(\mathbf{y}) \geq \sigma^2$.

For a numerical evaluation, we considered the SDPCM with $N = 5$, $S = 1$, $\sigma^2 = 1$, and $\mathbf{C}_k = \mathbf{e}_k \mathbf{e}_k^T$. We generated parameter vectors \mathbf{x}_0 with $j_0 = 1$ and different ξ_0 . In Fig. 1, we show the variance at \mathbf{x}_0 , $v(\hat{\mathbf{x}}(\cdot); \mathbf{x}_0) = \sum_{k \in [N]} v(\hat{x}_k(\cdot); \mathbf{x}_0)$ (computed by means of numerical integration), for the HT estimator using various choices of τ and for the ML estimator. The variance is plotted versus $\text{SNR}(\mathbf{x}_0) = \xi_0/\sigma^2$. Along with each variance curve, we display a corresponding lower bound that was calculated by evaluating (21) for each k , using for $\gamma_k(\mathbf{x})$ the mean function of the respective estimator (HT or ML), and summing over all k . (The mean functions of the HT and ML estimators were computed by means of numerical integration.) In evaluating (21), we used partial derivatives of order at most 1 in (17), and we chose for the evaluation of the lower bound $L = 2$, $\mathcal{K} = \{k\}$, $\mathbf{p}_1 = \mathbf{0}$, and $\mathbf{p}_2 = \mathbf{e}_k$. In Fig. 1, all variances and bounds are normalized by $2(\xi_0 + \sigma^2)^2$, which is the variance of the oracle estimator knowing j_0 .

It can be seen from Fig. 1 that in the high-SNR regime, for both estimators, the gap between the variance and the corresponding lower bound is quite small. This indicates that the variances of both estimators are nearly optimal (for the respective bias functions). However, in the low-SNR regime, the variances of the estimators tend to be significantly higher than the bounds. This means that there *may* be estimators with the same bias function as that of the HT or ML estimator but a lower variance. However, the actual existence of such estimators is not shown by our analysis.

6. CONCLUSION

We considered estimation of (a function of) a sparse vector \mathbf{x} that determines the covariance matrix of a Gaussian random vector via

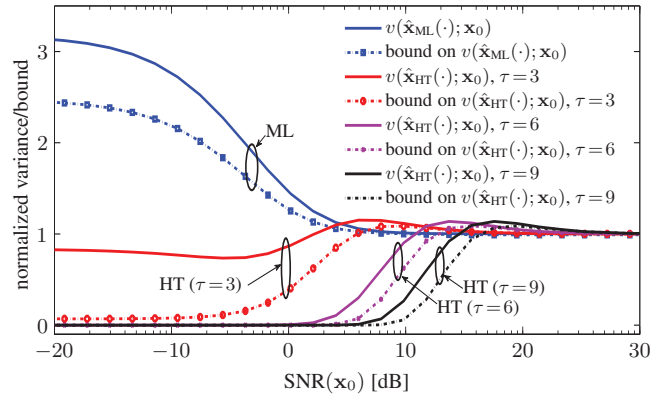


Fig. 1. Normalized variance of the HT and ML estimators and corresponding lower bounds versus $\text{SNR}(\mathbf{x}_0) = \xi_0/\sigma^2$, for the SDPCM with $N = 5$, $S = 1$, $\sigma^2 = 1$, and $\mathbf{C}_k = \mathbf{e}_k \mathbf{e}_k^T$.

a parametric covariance model. Using RKHS theory, we derived lower bounds on the estimator variance for a prescribed bias function. For the important special case of unbiased estimators of \mathbf{x} , we found that the transition of our bounds from low to high SNR is significantly slower than that of analogous bounds for the sparse linear model [6]. This suggests that for unbiased estimation, the prior information of sparsity is not as helpful as for the sparse linear model. Numerical results showed that for low SNR, the variance of two standard estimators (hard-thresholding estimator and maximum likelihood estimator) is significantly higher than our bounds. Hence, there might exist estimators that have the same bias function as these standard estimators but a smaller variance.

7. REFERENCES

- [1] P. Flandrin, *Time-Frequency/Time-Scale Analysis*. San Diego, CA: Academic Press, 1999.
- [2] F. Hlawatsch, *Time-Frequency Analysis and Synthesis of Linear Signal Spaces: Time-Frequency Filters, Signal Detection and Estimation, and Range-Doppler Estimation*. Boston, MA: Kluwer, 1998.
- [3] S. Haykin, "Cognitive radio: Brain-empowered wireless communication," *IEEE J. Sel. Areas Comm.*, vol. 23, pp. 201–220, Feb. 2005.
- [4] E. Parzen, "Statistical inference on time series by Hilbert space methods, I," Appl. Math. Stat. Lab., Stanford University, Stanford, CA, Tech. Rep. 23, Jan. 1959.
- [5] D. D. Duttweiler and T. Kailath, "RKHS approach to detection and estimation problems – Part V: Parameter estimation," *IEEE Trans. Inf. Theory*, vol. 19, no. 1, pp. 29–37, Jan. 1973.
- [6] S. Schmutzhard, A. Jung, F. Hlawatsch, Z. Ben-Haim, and Y. C. Eldar, "A lower bound on the estimator variance for the sparse linear model," in *Proc. 44th Asilomar Conf. Signals, Systems, Computers*, Pacific Grove, CA, Nov. 2010.
- [7] P. Rütimann and P. Bühlmann, "High dimensional sparse covariance estimation via directed acyclic graphs," *Electron. J. Statist.*, vol. 3, pp. 1133–1160, 2009.
- [8] Y. Chen, A. Wiesel, Y. C. Eldar, and A. O. Hero III, "Covariance estimation in decomposable Gaussian graphical models," *IEEE Trans. Signal Processing*, vol. 58, no. 3, pp. 1482–1492, March 2010.
- [9] N. Aronszajn, "Theory of reproducing kernels," *Trans. Am. Math. Soc.*, vol. 68, no. 3, pp. 337–404, May 1950.
- [10] M. E. Tipping and C. M. Bishop, "Probabilistic principal component analysis," *J. Roy. Stat. Soc. Ser. B*, vol. 61, pp. 611–622, 1999.
- [11] D. L. Donoho and I. M. Johnstone, "Ideal spatial adaptation by wavelet shrinkage," *Biometrika*, vol. 81, pp. 425–455, 1994.