# The Sparse Principal Component of a Constant-rank Matrix

Megasthenis Asteris, *Student Member, IEEE,* Dimitris S. Papailiopoulos, *Student Member, IEEE,* and George N. Karystinos, *Member, IEEE*

*Abstract*—The computation of the sparse principal component of a matrix is equivalent to the identification of its principal submatrix with the largest maximum eigenvalue. Finding this optimal submatrix is what renders the problem $\mathcal{NP}$-hard. In this work, we prove that, if the matrix is positive semidefinite and its rank is constant, then its sparse principal component is polynomially computable. Our proof utilizes the auxiliary unit vector technique that has been recently developed to identify problems that are polynomially solvable. Moreover, we use this technique to design an algorithm which, for any sparsity value, computes the sparse principal component with complexity $\mathcal{O}\left(N^{D+1}\right)$, where $N$ and $D$ are the matrix size and rank, respectively. Our algorithm is fully parallelizable and memory efficient.

*Index Terms*—Eigenvalues and eigenfunctions, feature extraction, information processing, machine learning algorithms, principal component analysis, signal processing algorithms.

## I. INTRODUCTION

Principal component analysis (PCA) is a well studied and broadly used dimensionality reduction tool. The principal components (PCs) of a set of observations on $N$ variables capture orthogonal directions of maximum variance and offer a distance-optimal, low-dimensional representation that -for many purposes- conveys sufficient amount of information. Without additional constraints, the PCs of a data set can be computed in polynomial time in $N$ using the eigenvalue decomposition.

A disadvantage of conventional PCA is that, in general, the extracted components are expected to have nonzero elements in all their entries. In many applications, sparse vectors that convey information are more favorable either due to sparsity of the actual signals [1], [2] or because sparsity implies interpretability [3] when each coordinate of a PC corresponds, for example, to a different word in text analysis applications

or the expression of a particular gene in bio data sets. Thus, provided that the application requires it, some of the maximum variance of the true PCs may be traded for sparsity. Recently, there has been an increasing interest in computing sparse components of data sets with applications that range from signal processing, communication networks, and machine learning, to bioinformatics, finance, and meteorology [4]-[13].

To enforce sparsity on the extracted components, a linearly constrained $l_0$-norm minimization problem is usually considered [1], [2], [14], [15]. This problem is equivalent to the sparse variance maximization, that is, the maximization of the Rayleigh quotient of a matrix under an $l_0$-norm constraint on the maximizing argument [4]-[6], [8], [11], [13], [16]-[23]. In both problems, due to the additional cardinality constraint that is enforced, the sparsity-aware flavor of PCA, termed sparse PCA, comes at a higher cost: sparse PCA is an $\mathcal{NP}$-hard problem [16].

To approximate sparse PCA, various methods have been introduced in the literature. A modified PCA technique based on the LASSO was introduced in [24]. In [3], a nonconvex regression-type optimization approach combined with LASSO penalty was used to approximately tackle the problem. A nonconvex technique, locally solving difference-of-convex-functions programs, was presented in [17]. Semidefinite programming (SDP) was used in [5], [22], while [18] augmented the SDP approach with an extra greedy step that offers favorable optimality guarantees under certain sufficient conditions. The authors of [4] considered greedy and branch-and-bound approaches. Generalized power method techniques using convex programs were also used to approximately solve sparse PCA [20]. A sparse-adjusted deflation procedure was introduced in [19] and in [6] optimality guarantees were shown for specific types of covariance matrices under thresholding and SDP relaxations. Iterative thresholding was also considered in [25] in conjunction with certain guarantees while a truncated power method was presented in [26].

In this present work, we prove that the sparse principal component of an $N \times N$ matrix $\mathbf{C}$ can be obtained in polynomial time under a new sufficient condition: when $\mathbf{C}$ can be written as a sum of a scaled identity matrix and a positive semidefinite update, i.e., $\mathbf{C} = \sigma\mathbf{I}_N + \mathbf{A}$, and the rank $D$ of the update $\mathbf{A}$ is not a function of the problem size.[1] Under this condition, we show that sparse PCA is solvable with complexity $\mathcal{O}\left(N^{D+1}\right)$. Our proof utilizes the

---

[1] If $\sigma = 0$, then we simply have a constant-rank matrix $\mathbf{C}$.

auxiliary unit vector technique that we developed in [27], [28]. This technique has been inspired by the work in [29], which reappeared in [30] and was used in [31]. It introduces an auxiliary unit vector that unlocks the constant-rank structure of a matrix (in this present work, matrix $\mathbf{A}$). The constant-rank property along with the auxiliary vector enable us to scan a constant-dimensional space and identify a polynomial number of candidate vectors (i.e., candidate solutions to the original problem). Interestingly, the optimal solution always lies among these candidates and a polynomial time search can always retrieve it. As a result, we have applied the auxiliary unit vector technique to identify the polynomial solvability of certain optimization problems and provide polynomial-time algorithms that are directly implementable, fully parallelizable, and memory efficient [32]-[35].

The rest of this paper is organized as follows. In Section II, we state the sparse PCA problem and indicate its $\mathcal{NP}$-hardness. Then, in Section III, we follow the principles of the auxiliary unit vector technique to present a proof of the polynomial solvability of the sparse PCA problem under a constant-rank condition. Moreover, we design a novel algorithm[2] which, for any sparsity value, computes the sparse principal component with complexity $\mathcal{O}\left(N^{D+1}\right)$. Our algorithm is simply implementable, fully parallelizable, and memory efficient. Especially for the case $D = 2$, an alternative nonparallelizable version of our algorithm with complexity $\mathcal{O}\left(N^2 \log N\right)$ is presented in Section IV.[3] A few conclusions are drawn in Section V.

## II. PROBLEM STATEMENT

We are interested in the computation of the real, unit-norm, and at most $K$-sparse principal component of a $N \times N$ matrix $\mathbf{C}$, i.e.,

$$\mathbf{x}_{\text{opt}} \triangleq \underset{\substack{\mathbf{x} \in \mathbb{R}^N \\ \|\mathbf{x}\|=1, \|\mathbf{x}\|_0 \leq K}}{\arg\max} \left\{ \mathbf{x}^T \mathbf{C} \mathbf{x} \right\}. \quad (1)$$

Interestingly, when $\mathbf{C}$ can be decomposed as a constant-rank positive semidefinite update of the identity matrix, i.e.,

$$\mathbf{C} = \sigma \mathbf{I}_N + \mathbf{A} \quad (2)$$

where $\sigma \in \mathbb{R}$, $\mathbf{I}_N$ is the $N \times N$ identity matrix, and $\mathbf{A}$ is a positive semidefinite matrix with rank $D$, then the optimization (1) can always be rewritten as

$$\mathbf{x}_{\text{opt}} = \underset{\substack{\mathbf{x} \in \mathbb{R}^N \\ \|\mathbf{x}\|=1, \|\mathbf{x}\|_0 \leq K}}{\arg\max} \left\{ \mathbf{x}^T \left(\sigma \mathbf{I}_N + \mathbf{A}\right) \mathbf{x} \right\} = \underset{\substack{\mathbf{x} \in \mathbb{R}^N \\ \|\mathbf{x}\|=1, \|\mathbf{x}\|_0 \leq K}}{\arg\max} \left\{ \mathbf{x}^T \mathbf{A} \mathbf{x} \right\}. \quad (3)$$

Since $\mathbf{A}$ is positive semidefinite and has rank $D$, it can be decomposed as

$$\mathbf{A} = \mathbf{V} \mathbf{V}^T, \quad (4)$$

where

$$\mathbf{V} \triangleq \left[\mathbf{v}_1 \; \mathbf{v}_2 \; \ldots \; \mathbf{v}_D\right] \quad (5)$$

---

[2] Early versions of our algorithm appeared in [36], [37].

[3] Alternative (non)parallelizable implementations of our algorithm for $D = 3$ have been presented in [36].

is an $N \times D$ matrix, and problem (1) can be written as

$$\mathbf{x}_{\text{opt}} = \underset{\substack{\mathbf{x} \in \mathbb{R}^N \\ \|\mathbf{x}\|=1, \|\mathbf{x}\|_0 \leq K}}{\arg\max} \left\{ \mathbf{x}^T \mathbf{V} \mathbf{V}^T \mathbf{x} \right\} = \underset{\substack{\mathbf{x} \in \mathbb{R}^N \\ \|\mathbf{x}\|=1, \|\mathbf{x}\|_0 \leq K}}{\arg\max} \left\| \mathbf{V}^T \mathbf{x} \right\|. \quad (6)$$

For the optimization problem in (6), we note that

$$\underset{\substack{\mathbf{x} \in \mathbb{R}^N \\ \|\mathbf{x}\|=1, \|\mathbf{x}\|_0 \leq K}}{\max} \left\| \mathbf{V}^T \mathbf{x} \right\| = \underset{\substack{\mathcal{I} \subseteq [N] \\ |\mathcal{I}|=K}}{\max} \underset{\substack{\mathbf{x} \in \mathbb{R}^K \\ \|\mathbf{x}\|=1}}{\max} \left\| \mathbf{V}_{\mathcal{I},:}^T \mathbf{x} \right\| \quad (7)$$

where $[N] \triangleq \{1, 2, \ldots, N\}$. In (7), set $\mathcal{I} \subseteq [N]$ (which we call the *support*) consists of the indices of the $K$ potentially nonzero elements of $\mathbf{x} \in \mathbb{R}^N$. For a given support $\mathcal{I}$, the inner maximization is a $K$-dimensional principal-component problem, where $\mathbf{V}_{\mathcal{I},:}$ is the corresponding $K \times D$ submatrix of $\mathbf{V}$. The solution to the inner maximization is denoted by

$$\mathbf{x}(\mathcal{I}) \triangleq \underset{\substack{\mathbf{x} \in \mathbb{R}^K \\ \|\mathbf{x}\|=1}}{\arg\max} \left\| \mathbf{V}_{\mathcal{I},:}^T \mathbf{x} \right\| \quad (8)$$

and equals the principal left singular vector of $\mathbf{V}_{\mathcal{I},:}$. Then, our optimization problem in (7) becomes

$$\mathcal{I}_{\text{opt}} \triangleq \underset{\substack{\mathcal{I} \subseteq [N] \\ |\mathcal{I}|=K}}{\arg\max} \left\{ \sigma_{\max}\left(\mathbf{V}_{\mathcal{I},:}\right) \right\} \quad (9)$$

where $\sigma_{\max}(\mathbf{V})$ denotes the principal singular value of matrix $\mathbf{V}$. That is, to solve our original problem in (1), according to (9), we need to find the $K$-row submatrix of $\mathbf{V}$ whose principal singular value is the maximum one among all submatrices. The indices that are contained in the optimal support $\mathcal{I}_{\text{opt}}$ that solves (9) correspond to the nonzero loadings of the solution $\mathbf{x}_{\text{opt}}$ to (1). Then, according to (8), the values of these nonzero loadings are directly computed by the left singular vector of $\mathbf{V}_{\mathcal{I}_{\text{opt}},:}$.

From the above discussion, it turns out that the hardness of the original problem in (1) comes from the identification of the optimal support $\mathcal{I}_{\text{opt}}$ in (9). To obtain the optimal support $\mathcal{I}_{\text{opt}}$, we could simply perform an exhaustive search among all $\binom{N}{K}$ possible supports $\mathcal{I}$ and compare them against the metric of interest in (9). However, if $K$ is not constant but grows with $N$, then such an approach has complexity that is exponential in $N$, indicating the $\mathcal{NP}$-hardness of (1), which was shown in [16]. In this present work, we show that, if the rank $D$ of $\mathbf{A}$ is constant, then (9) can be solved in time polynomial in $N$. In fact, we develop an algorithm that has complexity $\mathcal{O}\left(N^{D+1}\right)$ and returns $\mathcal{O}\left(N^D\right)$ candidate supports, one of which is guaranteed to be the solution to (9). Then, by an exhaustive search among only these candidate supports, we identify the optimal support in (9) and, hence, the sparse principal component of $\mathbf{A}$ and $\mathbf{C}$ with complexity polynomial in $N$, for any sparsity value $K$ between 1 and $N$ (that is, even if $K$ grows with $N$).

## III. COMPUTATION OF THE SPARSE PRINCIPAL COMPONENT IN TIME $\mathcal{O}\left(N^{D+1}\right)$

Prior to presenting the main result for the general rank-$D$ case, in the following subsection we provide insights as to why the sparse principal component of constant-rank matrices can be solved in polynomial time by first considering the trivial case $D = 1$.

### A. Rank-1: A motivating example

In this case, $\mathbf{A}$ has rank 1 and $\mathbf{V} = \mathbf{v} \in \mathbb{R}^N$. For a given support $\mathcal{I}$, we have $\mathbf{V}_{\mathcal{I},:} = \mathbf{v}_{\mathcal{I}}$. Then, our optimization problem in (7) becomes

$$\max_{\substack{\mathcal{I} \subseteq [N] \\ |\mathcal{I}|=K}} \max_{\substack{\mathbf{x} \in \mathbb{R}^K \\ \|\mathbf{x}\|=1}} \left| \mathbf{v}_{\mathcal{I}}^T \mathbf{x} \right| \qquad (10)$$

where, for any given support $\mathcal{I}$, the corresponding vector in (8) is

$$\mathbf{x}(\mathcal{I}) = \arg\max_{\substack{\mathbf{x} \in \mathbb{R}^K \\ \|\mathbf{x}\|=1}} \left| \mathbf{v}_{\mathcal{I}}^T \mathbf{x} \right| = \frac{\mathbf{v}_{\mathcal{I}}}{\|\mathbf{v}_{\mathcal{I}}\|}. \qquad (11)$$

Therefore, (7) becomes

$$\max_{\substack{\mathcal{I} \subseteq [N] \\ |\mathcal{I}|=K}} \left| \mathbf{v}_{\mathcal{I}}^T \frac{\mathbf{v}_{\mathcal{I}}}{\|\mathbf{v}_{\mathcal{I}}\|} \right| = \max_{\substack{\mathcal{I} \subseteq [N] \\ |\mathcal{I}|=K}} \|\mathbf{v}_{\mathcal{I}}\| \qquad (12)$$

and the optimal support is

$$\mathcal{I}_{\text{opt}} = \arg\max_{\substack{\mathcal{I} \subseteq [N] \\ |\mathcal{I}|=K}} \|\mathbf{v}_{\mathcal{I}}\| = \arg\max_{\substack{\mathcal{I} \subseteq [N] \\ |\mathcal{I}|=K}} \sum_{n \in \mathcal{I}} v_n^2. \qquad (13)$$

That is, to determine the solution to (9), we only need to compare the elements of $|\mathbf{v}|$ and select the $K$ largest ones. Then, their indices are the elements of $\mathcal{I}_{\text{opt}}$.

The above observation, although simple, turns out to be critical for the developments that follow. Hence, to simplify the presentation, we define function $top_k$ which is parameterized in an integer $k$, has as input a vector $\mathbf{u}$ of length $N \geq k$, and returns the indices of the $k$ largest values in $|\mathbf{u}|$:

$$\text{top}_k(\mathbf{u}) \triangleq \arg\max_{\substack{\mathcal{I} \subseteq [N] \\ |\mathcal{I}|=k}} \|\mathbf{u}_{\mathcal{I}}\|. \qquad (14)$$

Function $\text{top}_k(\mathbf{u})$ operates by selecting the indices of the $k$ largest values among $|u_1|, |u_2|, \ldots, |u_N|$. Its complexity is $\mathcal{O}(N)$ [38].

We conclude this subsection by noting that, if $D = 1$, then the optimal support in (9) is

$$\mathcal{I}_{\text{opt}} = \text{top}_K(\mathbf{v}) \qquad (15)$$

and is computed with linear complexity.

### B. Rank-D: Utilizing the auxiliary unit vector technique

We consider now the case where $\mathbf{A}$ has rank $D \geq 1$ and, hence, $\mathbf{V}$ is an $N \times D$ matrix. Without loss of generality (w.l.o.g.), we assume that each row of $\mathbf{V}$ has at least one nonzero element, i.e., $\mathbf{V}_{n,1:D} \neq \mathbf{0}_{1 \times D}, \forall\ n \in [N]$. Indeed, as explained in [28], if there exists an index $n \in [N]$ such that $\mathbf{V}_{n,:} = \mathbf{0}$, then, independently of the value of the corresponding element $x_n$ of $\mathbf{x}$, the contribution of this row to the value of $\|\mathbf{V}^T \mathbf{x}\|$ in (6) will be zero. Hence, there is no point in "spending" in $x_n$ a weight that could be distributed to other elements of $\mathbf{x}$; we can ignore the $n$th row of $\mathbf{V}$, replace $\mathbf{V}$ by $\mathbf{V}_{[N]-\{n\},:}$, and, hence, reduce the problem size from $N$ to $N-1$. In the final solution $\mathbf{x}_{\text{opt}}$, $x_n$ will be set to zero.

In our subsequent developments, we rely on the auxiliary unit vector technique that was introduced in [27] for matrices

of rank $D = 2$ and generalized in [28] for matrices of rank $D \geq 2$. This technique utilizes an auxiliary vector $\mathbf{c}$ to generate the subspace spanned by the columns of $\mathbf{V}$ and result in a rank-1 problem for each value of $\mathbf{c}$. Interestingly, for several problems, the number of different solutions that we obtain as $\mathbf{c}$ scans the unit-radius hypersphere has polynomial size. If the rank-1 problem for each value of $\mathbf{c}$ is polynomially solvable (as, for example, in the optimization problem of this work, as indicated in Subsection III-A), then the optimal solution is obtained with overall polynomial complexity. In a few words, the auxiliary unit vector technique of [27], [28] is a fully parallelizable and memory efficient technique that translates $D$-dimensional problems into a polynomial collection of rank-1 problems among which one results in the overall optimal solution.

For our sparse-principal-component problem, the auxiliary unit vector technique works as follows. Consider a unit vector $\mathbf{c} \in \mathbb{R}^D$. By Cauchy-Schwartz Inequality, for any $\mathbf{a} \in \mathbb{R}^D$,

$$\left| \mathbf{a}^T \mathbf{c} \right| \leq \|\mathbf{a}\| \|\mathbf{c}\| = \|\mathbf{a}\| \qquad (16)$$

with equality if and only if $\mathbf{c}$ is collinear to $\mathbf{a}$. Then,

$$\max_{\mathbf{c} \in \mathbb{R}^D, \|\mathbf{c}\|=1} \left| \mathbf{a}^T \mathbf{c} \right| = \|\mathbf{a}\|. \qquad (17)$$

Using (17), our optimization problem in (7) becomes

$$\max_{\substack{\mathcal{I} \subseteq [N] \\ |\mathcal{I}|=K}} \max_{\substack{\mathbf{x} \in \mathbb{R}^K \\ \|\mathbf{x}\|=1}} \left\| \mathbf{V}_{\mathcal{I},:}^T \mathbf{x} \right\| = \max_{\substack{\mathcal{I} \subseteq [N] \\ |\mathcal{I}|=K}} \max_{\substack{\mathbf{x} \in \mathbb{R}^K \\ \|\mathbf{x}\|=1}} \max_{\substack{\mathbf{c} \in \mathbb{R}^D \\ \|\mathbf{c}\|=1}} \left| \mathbf{x}^T \mathbf{V}_{\mathcal{I},:} \mathbf{c} \right|$$

$$= \max_{\substack{\mathbf{c} \in \mathbb{R}^D \\ \|\mathbf{c}\|=1}} \max_{\substack{\mathcal{I} \subseteq [N] \\ |\mathcal{I}|=K}} \max_{\substack{\mathbf{x} \in \mathbb{R}^K \\ \|\mathbf{x}\|=1}} \left| \mathbf{x}^T \mathbf{u}_{\mathcal{I}}(\mathbf{c}) \right| \quad (18)$$

where

$$\mathbf{u}(\mathbf{c}) \triangleq \mathbf{V}\mathbf{c}. \qquad (19)$$

The rightmost equality in (18) is obtained by interchanging the maximizations. This is a critical step of the auxiliary unit vector technique. It unlocks the constant-rank structure of $\mathbf{V}$ and allows us to consider a simple rank-1 problem for each value of $\mathbf{c}$. Indeed, for each $\mathbf{c} \in \mathbb{R}^D$, the inner double maximization problem

$$\max_{\substack{\mathcal{I} \subseteq [N] \\ |\mathcal{I}|=K}} \max_{\substack{\mathbf{x} \in \mathbb{R}^K \\ \|\mathbf{x}\|=1}} \left| \mathbf{u}_{\mathcal{I}}(\mathbf{c})^T \mathbf{x} \right| \qquad (20)$$

is equivalent to the rank-1 optimization problem in (10) that, according to (15), results in the optimal support (for fixed $\mathbf{c}$)

$$\mathcal{I}(\mathbf{c}) \triangleq \text{top}_K(\mathbf{u}(\mathbf{c})) \qquad (21)$$

which is obtained with complexity $\mathcal{O}(N)$. Then, according to (18), the solution to our original problem in (9) is met by collecting all possible supports $\mathcal{I}(\mathbf{c})$ as $\mathbf{c}$ scans the unit-radius $D$-dimensional hypersphere. That is, $\mathcal{I}_{\text{opt}}$ in (9) belongs to

$$\mathcal{S} \triangleq \bigcup_{\mathbf{c} \in \mathbb{R}^D, \|\mathbf{c}\|=1} \mathcal{I}(\mathbf{c}). \qquad (22)$$

Set $\mathcal{S}$ contains candidate supports $\mathcal{I} \subseteq [N]$ one of which is the solution to our original optimization problem. If $\mathcal{S}$ was available, then one would have to compare the elements of $\mathcal{S}$ against the metric of interest in (9) to obtain the optimal

This article has been accepted for publication in a future issue of IEEE Transactions on Information Theory, but has not been fully edited. Content may change prior to final publication.

4

support $\mathcal{I}_{\text{opt}}$. Therefore, the size of $\mathcal{S}$ and the complexity to build $\mathcal{S}$ determine the overall complexity to solve (9). Our major contribution in this work is that we show that the cardinality of $\mathcal{S}$ is upper bounded by

$$|\mathcal{S}| \le 2^{D-1} \binom{D}{\lfloor \frac{D}{2} \rfloor} \binom{N}{D} = \mathcal{O}\left(N^D\right) \qquad (23)$$

and develop an algorithm to build $\mathcal{S}$ with complexity $\mathcal{O}\left(N^{D+1}\right)$. After $\mathcal{S}$ is constructed, each element (support) $\mathcal{I}$ of it is mapped to the principal singular value of the $K \times D$ matrix $\mathbf{V}_{\mathcal{I},:}$ with complexity $\mathcal{O}\left(KD^2\right) = \mathcal{O}(K)$, since $D$ is constant. Finally, all computed singular values are compared with each other to obtain the optimal support $\mathcal{I}_{\text{opt}}$ in (9). Then, the solution to our original problem in (1) is the principal left singular vector of the $K \times D$ matrix $\mathbf{V}_{\mathcal{I}_{\text{opt}},:}$, computed with complexity $\mathcal{O}\left(KD^2\right) = \mathcal{O}(K)$. Therefore, we compute the optimal support $\mathcal{I}_{\text{opt}}$ and the sparse principal component of a rank-$D$ matrix with complexity $\mathcal{O}\left(N^{D+1} + N^D K\right) = \mathcal{O}\left(N^{D+1}\right)$.

A constructive proof is presented in detail in the next three subsections. To give some insight of the proof, we begin with the simple case $D = 2$. Then, we generalize our proof for the case of any arbitrary $D$.

*C. Rank-2: A simple instance of our proof*

If $D = 2$, then $\mathbf{V}$ has size $N \times 2$ and the auxiliary vector $\mathbf{c}$ is a length-2, unit vector that, as in [27], can be parameterized in an auxiliary angle $\phi \in \left(-\frac{\pi}{2}, \frac{\pi}{2}\right]$. That is,

$$\mathbf{c}(\phi) \triangleq \begin{bmatrix} \sin\phi \\ \cos\phi \end{bmatrix}, \quad \phi \in \Phi \triangleq \left(-\frac{\pi}{2}, \frac{\pi}{2}\right]. \qquad (24)$$

Hence, $\mathbf{c}(\phi)$ lies on the unit-radius semicircle.[4] Then, the candidate set in (22) is re-expressed as

$$\mathcal{S} = \bigcup_{\phi \in \Phi} \mathcal{I}(\phi) \qquad (25)$$

where, according to (21),

$$\mathcal{I}(\phi) \triangleq \text{top}_K(\mathbf{u}(\phi)) \qquad (26)$$

and, according to (19),

$$\mathbf{u}(\phi) \triangleq \mathbf{V}\mathbf{c}(\phi). \qquad (27)$$

That is, for any given $\phi \in \Phi$, the corresponding support $\mathcal{I}(\phi)$ is obtained with complexity $\mathcal{O}(N)$ by selecting the indices of the $K$ absolutely largest elements of $\mathbf{u}(\phi)$.

However, why should $\phi$ simplify the computation of a solution? The intuition behind the auxiliary unit vector technique is that every element of $\mathbf{u}(\phi)$ is actually a continuous function of $\phi$, i.e., a curve (or 1-manifold) in $\phi$, since

$$\mathbf{u}(\phi) = \mathbf{V}\mathbf{c}(\phi) = \begin{bmatrix} V_{1,1}\sin\phi + V_{1,2}\cos\phi \\ V_{2,1}\sin\phi + V_{2,2}\cos\phi \\ \vdots \\ V_{N,1}\sin\phi + V_{N,2}\cos\phi \end{bmatrix}. \qquad (28)$$

---

[4]We ignore the other semicircle because any pair of angles $\phi_1$ and $\phi_2$ with difference $\pi$ results in opposite vectors $\mathbf{c}(\phi_1) = -\mathbf{c}(\phi_2)$ which, however, are equivalent with respect to the optimization metric in (18) and produce the same support $\mathcal{I}(\mathbf{c}(\phi_1)) = \mathcal{I}(\mathbf{c}(\phi_2))$ in (21).
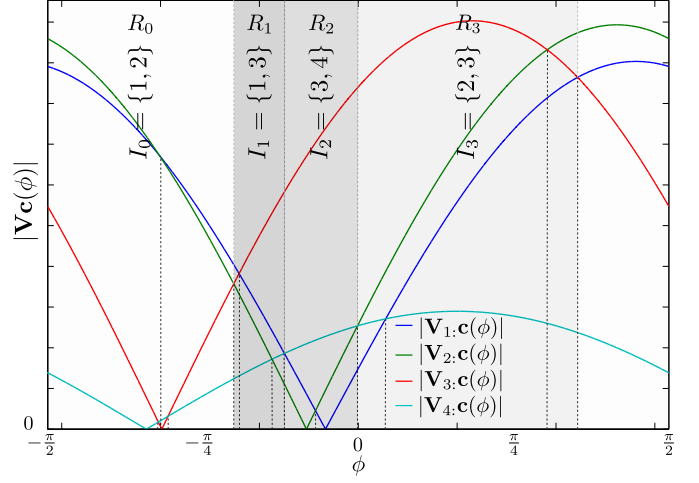


Fig. 1. Partition of $\Phi$ into $2\binom{4}{2} = 12$ intervals and 4 regions (for sparsity $K = 2$), generated by the rows of an arbitrary $4 \times 2$ matrix $\mathbf{V}$.

Hence, the support that corresponds to the $K$ absolutely largest elements of $\mathbf{u}(\phi)$ at a given point $\phi$ is a function of $\phi$. Due to the continuity of the curves and the discrete nature of the support, we expect that the support $\mathcal{I}(\phi)$ will retain the same elements in an area around $\phi$. Therefore, we expect the formation of intervals in $\Phi$, within which the indices of the $K$ absolutely largest elements of $\mathbf{u}(\phi)$ remain unaltered. A support $\mathcal{I}$ might change only if the sorting of the amplitudes of two elements in $\mathbf{u}(\phi)$, say $|u_i(\phi)|$ and $|u_j(\phi)|$, changes. This occurs at points $\phi$ where $|u_i(\phi)| = |u_j(\phi)|$, that is, points where two curves intersect. Finding all these *intersection points* is sufficient to determine intervals and construct all possible candidate supports $\mathcal{I}$. Among all candidate supports, lies the support that corresponds to the optimal $K$-sparse principal component. Exhaustively checking the supports $\mathcal{I}$ of all intervals suffices to retrieve the optimal one. The number of these intervals is exactly equal to number of possible intersections among the amplitudes of $|\mathbf{u}(\phi)|$, which is exactly equal to $2\binom{N}{2} = \mathcal{O}\left(N^2\right)$, counting all possible combinations of element pairs.

Before we proceed, in Fig. 1, we illustrate the interval partition of $\Phi$ for an arbitrary $4 \times 2$ matrix $\mathbf{V}$ (i.e., $N = 4$). We plot the 4 curves that originate from the 4 rows of $|\mathbf{u}(\phi)| = |\mathbf{V}\mathbf{c}(\phi)|$ and observe the intervals that are formed, within which the sorting of the curves does not change. The borders of the intervals are denoted by vertical dashed lines at points of curve intersections. Our approach creates $2\binom{4}{2} = 12$ intervals which exceeds the total number of possible supports, however this is not true for greater values of $N$. In addition, for sparsity $K = 2$, we observe that $\Phi$ is partitioned into *regions* (sets of adjacent intervals); within each region $R_i$, although the sorting changes, the set of $K$ largest curves does not change. For example, in Fig. 1, we identify the regions $R_0$, $R_1$, $R_2$, and $R_3$ where the candidate support remains fixed. These regions are an interesting feature that might further decrease the number of intervals we need to check. We exploit this feature in the serial implementation of our algorithm in Section IV.

## D. Rank-2: Algorithmic developments and complexity

Our goal is the construction of all possible candidate $K$-sparse vectors, determined by the support $\mathcal{I}$ of each interval in $\Phi$. This is a two-step process. First, we identify interval borders and, then, we determine the supports associated with these intervals.

*Algorithmic Steps:* We first determine all possible intersections of curve pairs in $|\mathbf{u}(\phi)|$. Any pair $\{i, j\}$ of distinct elements in $|\mathbf{u}(\phi)|$ is associated with two intersections: $u_i(\phi) = u_j(\phi)$ and $u_i(\phi) = -u_j(\phi)$. Solving these two equations with respect to $\phi$ determines two points

$$\hat{\phi} = \mp \tan^{-1}\left(\frac{V_{i,2} - V_{j,2}}{V_{i,1} - V_{j,1}}\right) \in \Phi \qquad (29)$$

where a new support of indices of the $K$-largest values of $|\mathbf{u}(\phi)|$ might occur. We note that, when $\phi$ varies in $\Phi$ (from $-\frac{\pi}{2}$ to $\frac{\pi}{2}$), changes of the support $\mathcal{I}$ may occur only over intersection points $\hat{\phi}$ given in (29), for any $i, j \in [N]$ with $i \neq j$.

Next, we observe that, at an intersection point $\hat{\phi}$, the intersecting curves are equal, i.e., $|u_i(\hat{\phi})| = |u_j(\hat{\phi})|$, while their relative order changes over $\hat{\phi}$. That is, if $|u_i(\phi)| > |u_j(\phi)|$ in the interval immediately preceding $\hat{\phi}$, then $|u_i(\phi)| < |u_j(\phi)|$ in the interval immediately following $\hat{\phi}$, and vice versa. Exactly at the point of intersection, we can determine the order of the curves in $|\mathbf{u}(\hat{\phi})|$.

- If both or none of the $i$th and $j$th curves are included in the set of $K$ largest curves at $\hat{\phi}$, then none of the two curves leaves or joins, respectively, that set at $\hat{\phi}$, despite the change in their relative order. Assuming that no other pair of curves intersect exactly at $\hat{\phi}$,[5] the two intervals incident to $\hat{\phi}$ are associated with the same support $\mathcal{I}$ which can be determined exactly at the intersection point $\hat{\phi}$.[6]
- If that is not the case, i.e., if the $i$th and $j$th curves occupy the $K$th and $(K+1)$th order at $\hat{\phi}$, then the change in their relative order affects the support $\mathcal{I}$: one of the two curves leaves and the other one joins the set of $K$ largest curves at $\hat{\phi}$. The support sets associated with the two adjacent intervals differ only in one element (one contains index $i$ and the other contains index $j$ instead), while the remaining $K-1$ common indices correspond to the $K-1$ largest curves at the intersection point $\hat{\phi}$.[7]

We have fully described a way to calculate the (at most) two support sets associated with the two intervals incident to one intersection point. Since all intervals are incident to at least one intersection point, it suffices to consider all intersection points, that is, all $\binom{N}{2}$ pairwise combinations of elements in $|\mathbf{u}(\phi)|$, to collect the corresponding support sets.

*Computational Complexity:* A single intersection point can be computed in time $\mathcal{O}(1)$. At an intersection point $\hat{\phi}$, the $K$th-order element of $|\mathbf{u}(\phi)|$ (and the $K-1$ elements larger than that) can be determined in time $\mathcal{O}(N)$ which equals the construction time of the (at most) two support sets associated with the intervals incident to $\hat{\phi}$. Collecting all candidate supports $\mathcal{I}$ requires examining all $2\binom{N}{2}$ intersection points, implying a total construction cost of $2\binom{N}{2} \times \mathcal{O}(N) = \mathcal{O}(N^3)$. Since we obtain (at most) two supports for any of the $2\binom{N}{2}$ intersection points, the size of the candidate support set $\mathcal{S}$ is $|\mathcal{S}| \leq 4\binom{N}{2}$.

### E. Rank-D: A generalized proof

In the general case, $\mathbf{V}$ is a $N \times D$ matrix. In this subsection, we present our main result where we prove that the problem of identifying the $K$-sparse principal component of a rank-$D$ matrix is solvable with complexity $\mathcal{O}(N^{D+1})$. The statement is true for any value of $K$ (that is, even if $K$ is a function of $N$). The rest of this subsection contains a constructive proof of this statement.

Since $\mathbf{V}$ has size $N \times D$, the auxiliary vector $\mathbf{c}$ is a length-$D$, unit vector. We begin our constructive proof by introducing the auxiliary-angle vector $\boldsymbol{\phi} \in \Phi^{D-1}$ and parameterizing $\mathbf{c}$, as in [28], according to

$$\mathbf{c}(\boldsymbol{\phi}) \triangleq \begin{bmatrix} \sin\phi_1 \\ \cos\phi_1 \sin\phi_2 \\ \cos\phi_1 \cos\phi_2 \sin\phi_3 \\ \vdots \\ \cos\phi_1 \cos\phi_2 \ldots \sin\phi_{D-1} \\ \cos\phi_1 \cos\phi_2 \ldots \cos\phi_{D-1} \end{bmatrix}. \qquad (30)$$

Hence, $\mathbf{c}(\boldsymbol{\phi})$ lies on the unit-radius semihypersphere.[8] Then, the candidate set in (22) is re-expressed as

$$\mathcal{S} = \bigcup_{\boldsymbol{\phi} \in \Phi^{D-1}} \mathcal{I}(\boldsymbol{\phi}) \qquad (31)$$

where, according to (21),

$$\mathcal{I}(\boldsymbol{\phi}) \triangleq \text{top}_K(\mathbf{u}(\boldsymbol{\phi})) \qquad (32)$$

and, according to (19),

$$\mathbf{u}(\boldsymbol{\phi}) \triangleq \mathbf{V}\mathbf{c}(\boldsymbol{\phi}). \qquad (33)$$

That is, for any given $\boldsymbol{\phi} \in \Phi^{D-1}$, the corresponding support $\mathcal{I}(\boldsymbol{\phi})$ is obtained with complexity $\mathcal{O}(N)$ by selecting the indices of the $K$ absolutely largest elements of $\mathbf{u}(\boldsymbol{\phi})$.

---

[5]Recall that the support $\mathcal{I}$ can change only at an intersection point. If a second pair of curves, say the $l$th and $k$th ones, intersect exactly at $\hat{\phi}$ causing the support $\mathcal{I}$ to change over $\hat{\phi}$, then the correct support sets for the two incident intervals will be determined when the intersections of the $l$th and $k$th curves are examined.

[6]In fact, since the support does not change at $\hat{\phi}$, there exists another intersection point where the same support will be collected. Hence, the support $\mathcal{I}$ at $\hat{\phi}$ may safely be ignored.

[7]Since any interval in $\Phi$ follows an intersection point, we actually need to evaluate the support $\mathcal{I}$ of only the interval that follows the intersection point $\hat{\phi}$. The other interval will be examined at the intersection point that precedes it. To identify which of the two indices $i$ and $j$ is contained in the support of the interval that follows the intersection point $\hat{\phi}$, we can visit the "rightmost" point of $\Phi$, that is, $\frac{\pi}{2}$. There, due to the continuity of $u_i(\phi)$ and $\pm u_j(\phi)$, the relative order of $|u_i(\phi)|$ and $|u_j(\phi)|$ within the interval that follows $\hat{\phi}$ will be the identical or opposite relative order of $u_i\left(\frac{\pi}{2}\right)$ and $\pm u_j\left(\frac{\pi}{2}\right)$, depending on whether $u_i(\hat{\phi})$ is positive or negative, respectively.

[8]As in the rank-2 case, we ignore the other semihypersphere because any pair of vectors $\boldsymbol{\phi}$ and $\tilde{\boldsymbol{\phi}}$ whose first elements $\phi_1$ and $\tilde{\phi}_1$, respectively, have difference $\pi$ results in opposite vectors $\mathbf{c}(\boldsymbol{\phi}) = -\mathbf{c}(\tilde{\boldsymbol{\phi}})$ which, however, are equivalent with respect to the optimization metric in (18) and produce the same support $\mathcal{I}(\mathbf{c}(\boldsymbol{\phi})) = \mathcal{I}(\mathbf{c}(\tilde{\boldsymbol{\phi}}))$ in (21).

This article has been accepted for publication in a future issue of IEEE Transactions on Information Theory, but has not been fully edited. Content may change prior to final publication.
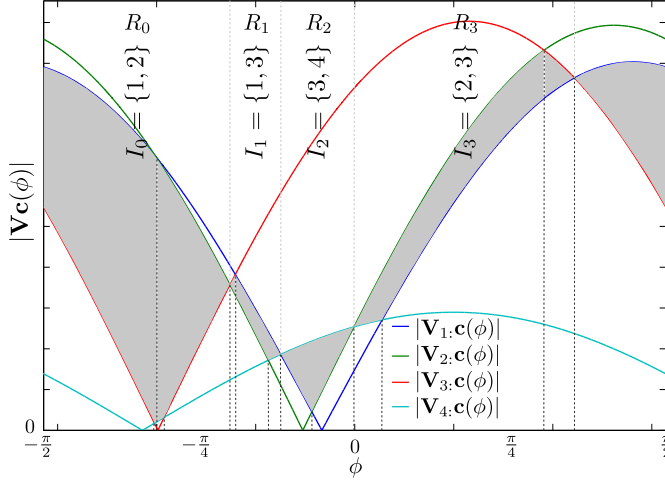
6



Fig. 2. Cells generated by the rows of an arbitrary $4 \times 2$ matrix $\mathbf{V}$. The shaded cells are the normal ones for sparsity $K = 2$.

To gain some intuition into the purpose of inserting the auxiliary-angle vector $\boldsymbol{\phi}$, notice that every element of $\mathbf{u}(\boldsymbol{\phi})$ in (34) is actually a continuous function of $\boldsymbol{\phi}$ and so are the elements of $|\mathbf{u}(\boldsymbol{\phi})|$. That is, each element of $|\mathbf{u}(\boldsymbol{\phi})|$ is a hypersurface (or $(D-1)$-manifold) in the $D$-dimensional space $\Phi^{D-1} \times [0, \infty)$. When we sort the $N$ elements of $|\mathbf{u}(\boldsymbol{\phi})|$ at a given point $\boldsymbol{\phi}$, we actually sort the $N$ hypersurfaces at point $\boldsymbol{\phi}$. The key observation in our algorithm is that, due to their continuity, the hypersurfaces will retain their sorting in an area "around" $\boldsymbol{\phi}$. This implies the partition of $\Phi^{D-1} \times [0, \infty)$ into cells $\mathcal{C}_1, \mathcal{C}_2, \ldots$, each of which (say, cell $\mathcal{C}$) is associated with a single set $\mathcal{I}^+(\mathcal{C}) \subseteq [N]$ of indices of hypersurfaces that lie above $\mathcal{C}$ and a single set $\mathcal{I}^-(\mathcal{C}) = [N] - \mathcal{I}^+(\mathcal{C})$ of indices of hypersurfaces that lie below it. Moreover, each cell $\mathcal{C}$ contains at least one vertex (that is, intersection of $D$ hypersurfaces). Finally, for any $\boldsymbol{\phi} \in \Phi$, there is a unique cell $\mathcal{C} \subset \Phi^{D-1} \times [0, \infty)$, called "normal," which contains uncountably many points in $\{\boldsymbol{\phi}\} \times [0, \infty)$ and is associated with a single index-set $\mathcal{I}^+(\mathcal{C})$ of cardinality $K$ (that is, exactly $K$ hypersurfaces lie above $\mathcal{C}$). In fact, the indices of these $K$ hypersurfaces (that is, the elements of $\mathcal{I}^+(\mathcal{C})$) are the elements of support $\mathcal{I}(\boldsymbol{\phi})$. Although our discussion refers to the general-$D$ case, for illustrative purposes we consider again the case $D = 2$ and, in Fig. 2, we revisit the example that we presented in Subsection III-C. The normal cells that are created by the $N = 4$ curves are the shaded ones. These cells carry the property that lie below exactly $K = 2$ curves. We observe that there is a one-to-one correspondence between normal cells and regions $R_i$.

According to (31) and the above observations, we need

to determine the index-set $\mathcal{I}^+(\mathcal{C})$ of every normal cell $\mathcal{C}$ in the partition. If we collect all such index-sets, then we have constructed $\mathcal{S}$ in (31). This will be achieved if, instead, we identify all cells in $\Phi^{D-1} \times [0, \infty)$ and, for each cell, determine the $K$ largest hypersurfaces that lie above an arbitrary point of it. The latter will return the desired index-set $\mathcal{I}^+(\mathcal{C})$ if the cell is normal. In Fig. 2, we observe that, for each normal cell, the indices of the $K = 2$ largest curves that lie above it can be computed at the leftmost vertex of it (we can ignore the leftmost normal cell because it produces the same indices with the rightmost one). In the following, we identify all cells in the partition and compute a size-$K$ support $\mathcal{I}$ for each such cell. This way, we obtain the index-set of any normal cell, among which one is the optimal support $\mathcal{I}$ in (9).

Since each cell contains at least one vertex, we only need to find all vertices in the partition and determine $\mathcal{I}^+$ for all neighboring cells. Recall that a vertex is an intersection of $D$ hypersurfaces. Consider $D$ arbitrary hypersurfaces; say, for example, $|u_1(\boldsymbol{\phi})|$, $|u_2(\boldsymbol{\phi})|$, $\ldots$, $|u_D(\boldsymbol{\phi})|$. Their intersection satisfies $|u_1(\boldsymbol{\phi})| = |u_2(\boldsymbol{\phi})| = \ldots = |u_D(\boldsymbol{\phi})|$ and is computed by solving the system of equations

$$\left\{ \begin{array}{c} u_1(\boldsymbol{\phi}) \pm u_2(\boldsymbol{\phi}) = 0 \\ u_1(\boldsymbol{\phi}) \pm u_3(\boldsymbol{\phi}) = 0 \\ \vdots \\ u_1(\boldsymbol{\phi}) \pm u_D(\boldsymbol{\phi}) = 0 \end{array} \right\} \tag{35}$$

or, equivalently,

$$\begin{bmatrix} \mathbf{V}_{1,:} \pm \mathbf{V}_{2,:} \\ \mathbf{V}_{1,:} \pm \mathbf{V}_{3,:} \\ \vdots \\ \mathbf{V}_{1,:} \pm \mathbf{V}_{D,:} \end{bmatrix} \mathbf{c}(\boldsymbol{\phi}) = \mathbf{0}. \tag{36}$$

For any sign combination, the solution to the latter consists of the spherical coordinates of the unit vector in the null space of the $(D-1) \times D$ leftmost matrix.[9] Then, the index-set $\mathcal{I}^+$ that corresponds to a neighboring cell is computed by

$$\text{top}_K(\mathbf{Vc}(\boldsymbol{\phi})). \tag{37}$$

Note that the $D$ intersecting hypersurfaces have the same value at $\boldsymbol{\phi}$. Hence, (37) returns ambiguity regarding the sorting of these particular $D$ hypersurfaces. If $d < D$ hypersurfaces of these belong to the $K$ largest ones, then, due to this ambiguity,

[9]If the $(D-1) \times D$ matrix is full-rank, then its null space has rank 1 and $\mathbf{c}(\boldsymbol{\phi})$ is uniquely determined (within a sign ambiguity which, however, does not affect the final decision on the index-set). If, instead, the $(D-1) \times D$ matrix is rank-deficient, then the intersection of the $D$ hypersurfaces (i.e., the solution of (36)) is a $p$-manifold (with $p \geq 1$) on the $D$-dimensional space and does not generate a new cell. Hence, combinations of $D$ rows of $\mathbf{V}$ that result in linearly dependent rows of the $(D-1) \times D$ matrix in (36) can be simply ignored.

$$\mathbf{u}(\boldsymbol{\phi}) = \mathbf{Vc}(\boldsymbol{\phi}) = \begin{bmatrix} \mathbf{V}_{1,:} \mathbf{c}(\boldsymbol{\phi}) \\ \mathbf{V}_{2,:} \mathbf{c}(\boldsymbol{\phi}) \\ \vdots \\ \mathbf{V}_{N,:} \mathbf{c}(\boldsymbol{\phi}) \end{bmatrix} = \begin{bmatrix} V_{1,1} \sin \phi_1 + \sum_{d=2}^{D-1} V_{1,d} \prod_{i=1}^{d-1} \cos \phi_i \sin \phi_d + V_{1,D} \prod_{i=1}^{D-1} \cos \phi_i \\ V_{2,1} \sin \phi_1 + \sum_{d=2}^{D-1} V_{2,d} \prod_{i=1}^{d-1} \cos \phi_i \sin \phi_d + V_{2,D} \prod_{i=1}^{D-1} \cos \phi_i \\ \vdots \\ V_{N,1} \sin \phi_1 + \sum_{d=2}^{D-1} V_{N,d} \prod_{i=1}^{d-1} \cos \phi_i \sin \phi_d + V_{N,D} \prod_{i=1}^{D-1} \cos \phi_i \end{bmatrix} \tag{34}$$

---

**Algorithm 1** Computation of the sparse principal component of a rank-$D$ matrix with complexity $\mathcal{O}\left(N^{D+1}\right)$

---

**Input:** $\mathbf{V} \in \mathbb{R}^{N \times D}$, $K \in [N]$

$\quad \mathcal{S} \leftarrow \{\}$ *(set of candidate supports)*
$\quad$ **for** all $\binom{N}{D}$ sets $\{i_1, \ldots, i_D\} \subseteq [N]$ **do**
$\quad\quad$ **for** all $2^{D-1}$ sequences $(b_1, \ldots, b_{D-1}) \in \{\pm 1\}^{D-1}$ **do**
$$\mathbf{c} \leftarrow \text{nullspace}\left( \begin{bmatrix} \mathbf{V}_{i_1,:} - b_1 \mathbf{V}_{i_2,:} \\ \vdots \\ \mathbf{V}_{i_1,:} - b_{D-1} \mathbf{V}_{i_D,:} \end{bmatrix} \right)$$
$\quad\quad\quad \mathcal{I} \leftarrow \text{top}_K(\mathbf{Vc})$
$\quad\quad\quad$ **if** $|\mathcal{I}| = K$ **then**
$\quad\quad\quad\quad \mathcal{S} \leftarrow \mathcal{S} \cup \{\mathcal{I}\}$
$\quad\quad\quad$ **else**
$\quad\quad\quad\quad \mathcal{T} \leftarrow \mathcal{I} - \{i_1, \ldots, i_D\}$
$\quad\quad\quad\quad r \leftarrow K - |\mathcal{T}|$
$\quad\quad\quad\quad$ **for** all $\binom{D}{r}$ $r$-subsets $\mathcal{M} \subseteq \{i_1, \ldots, i_D\}$ **do**
$\quad\quad\quad\quad\quad \hat{\mathcal{I}} \leftarrow \mathcal{T} \cup \mathcal{M}$
$\quad\quad\quad\quad\quad \mathcal{S} \leftarrow \mathcal{S} \cup \{\hat{\mathcal{I}}\}$
$\quad\quad\quad\quad$ **end for**
$\quad\quad\quad$ **end if**
$\quad\quad$ **end for**
$\quad$ **end for**
$\quad \mathcal{I}_{\text{opt}} = \arg\max_{\mathcal{I} \in \mathcal{S}} \sigma_{\max}(\mathbf{V}_{\mathcal{I},:})$
$\quad$ **return** $\mathcal{I}_{\text{opt}}$ & principal left singular vector of $\mathbf{V}_{\mathcal{I}_{\text{opt}},:}$

---

we have to consider all $\binom{D}{d}$ combinations of $d$ hypersurfaces among the $D$ intersecting ones, where $\binom{D}{d} < \binom{D}{\lfloor \frac{D}{2} \rfloor}$. Finally, we have to repeat the above procedure for all $2^{D-1}$ sign combinations in (36) and any combination of $D$ intersecting hypersurfaces among the $N$ ones. The total number of combinations is $\binom{N}{D}$, hence the cardinality of $\mathcal{S}$ is upper bounded by $2^{D-1} \binom{D}{\lfloor \frac{D}{2} \rfloor} \binom{N}{D} = \mathcal{O}\left(N^D\right)$.

A pseudocode that includes all the above steps is presented in Algorithm 1. Its complexity to build $\mathcal{S}$ is determined by the complexity to build each element of it (i.e., each index-set $\mathcal{I}^+$) for each examined intersection through (37). Note that function $\text{top}_K$ has complexity $\mathcal{O}(N)$ and the cardinality of $\mathcal{S}$ is $\mathcal{O}\left(N^D\right)$. Hence, the overall complexity to build $\mathcal{S}$ is $\mathcal{O}\left(N^{D+1}\right)$. Finally, we mention that the computation of each element of $\mathcal{S}$ is performed independently of each other. Therefore, the proposed Algorithm 1 that builds $\mathcal{S}$ and solves (1) or, equivalently, (9) with complexity $\mathcal{O}\left(N^{D+1}\right)$ is fully parallelizable and memory efficient.

## IV. AN ALGORITHM OF COMPLEXITY $\mathcal{O}\left(N^2 \log N\right)$ FOR RANK-2 MATRICES

In the special case of a rank-2 matrix $\mathbf{A}$ (i.e., $D = 2$), Algorithm 1 computes the sparse principal component of $\mathbf{A}$ with complexity $\mathcal{O}\left(N^3\right)$ for any sparsity value. In this section, we develop an algorithm that computes the sparse principal component of $\mathbf{A}$ with complexity $\mathcal{O}\left(N^2 \log N\right)$.

### A. A serial algorithm for rank-2 matrices

For the rank-2 case, Algorithm 1, discussed in Subsection III-D, relied on identifying a polynomial number of non-

overlapping intervals on $\Phi$, each associated with a candidate support set $\mathcal{I}$. These intervals are induced by the pairwise intersections of curves in $|\mathbf{u}(\phi)|$. The candidate support set associated with an interval is determined at the leftmost point of the interval. The two-step algorithm first computes all $2\binom{N}{2}$ pairwise intersection points. In the second step, it determines $\mathcal{I}(\hat{\phi})$, the set of indices of the $K$ largest elements of $|\mathbf{u}(\hat{\phi})|$, at each intersection point $\hat{\phi}$ individually, in linear time. In this section, we present an algorithmic enhancement that reduces the overall computational complexity, exploiting the correlation between candidate support sets of neighboring intervals.

The relative order of $|u_i(\phi)|$ and $|u_i(\phi)|$ changes only at the points where the $i$th and $j$th curves intersect. Conversely, if those two are the only curves intersecting at a point $\hat{\phi}$, then the ordering of the remaining elements of $|\mathbf{u}(\phi)|$ in the intervals immediately preceding and succeeding $\hat{\phi}$ is identical. The limited change on ordering of the curves in $|\mathbf{u}(\phi)|$ across an intersection point $\hat{\phi}$ implies that the differences between the candidate support sets of the adjacent intervals cannot be arbitrary. The support sets associated with two neighboring intervals differ in at most one element. More formally, let $\hat{\phi}$ be the intersection point of the $i$th and $j$th curves of $|\mathbf{u}(\phi)|$. There exists an $\epsilon > 0$ such that $\hat{\phi}$ is the only intersection point lying in $[\hat{\phi} - \epsilon, \hat{\phi} + \epsilon]$. If $|u_i(\phi)| < |u_j(\phi)|$ in $[\hat{\phi} - \epsilon, \hat{\phi})$, then $|u_i(\phi)| > |u_j(\phi)|$ in $(\hat{\phi}, \hat{\phi} + \epsilon]$, and vice versa. The ordering of all other curves remains unaltered over $\hat{\phi}$. If the $i$th and $j$th curves are both members of the candidate support in the interval preceding $\hat{\phi}$, then $\mathcal{I}(\hat{\phi} + \epsilon) = \mathcal{I}(\hat{\phi} - \epsilon)$. The same holds if neither is included in $\mathcal{I}(\hat{\phi} - \epsilon)$. On the contrary, if exactly one of $i$ and $j$ belongs to $\mathcal{I}(\hat{\phi} - \epsilon)$, then the candidate sets associated with the two neighboring intervals differ in exactly one element. If, w.l.o.g., the $i$th curve is the one belonging to $\mathcal{I}(\hat{\phi} - \epsilon)$, then $\mathcal{I}(\hat{\phi} + \epsilon) = \left( \mathcal{I}(\hat{\phi} - \epsilon) - \{i\} \right) \cup \{j\}$.

The key observation is that, if the candidate support set associated with a particular interval is known, then the set associated with a neighboring interval can be determined with a constant complexity. The above observations readily suggest the following procedure for determining the candidate support sets:

1) Compute all $2\binom{N}{2}$ pairwise intersection points of the curves in $|\mathbf{u}(\phi)|$.
2) Sort the intersection points in increasing order. Let $\phi_1, \phi_2, \ldots, \phi_{2\binom{N}{2}}$ be the sorted sequence.
3) Determine $\mathcal{I}(\phi_1)$, the set of indices of the $K$ largest curves in $|\mathbf{u}(\phi)|$ at the first intersection point $\phi_1$.
4) Successively determine $\mathcal{I}(\phi_t)$ at consecutive intersection points. At $\phi_t$, the candidate support set $\mathcal{I}(\phi_t)$ is determined by appropriately updating $\mathcal{I}(\phi_{t-1})$, the set associated with the previous interval.

Once all candidate support sets have been collected, the optimal solution is determined as in Algorithm 1.

An illustrative figure that explains the above steps is presented in Fig. 3, using the same $4 \times 2$ matrix $\mathbf{V}$ as the one examined in Section III (Figs. 2 and 3). We seek the optimal 2-sparse solution, i.e., $K = 2$. The algorithm starts at $\phi = -\pi/2$ and scans $\Phi$, keeping track of the $K$th order
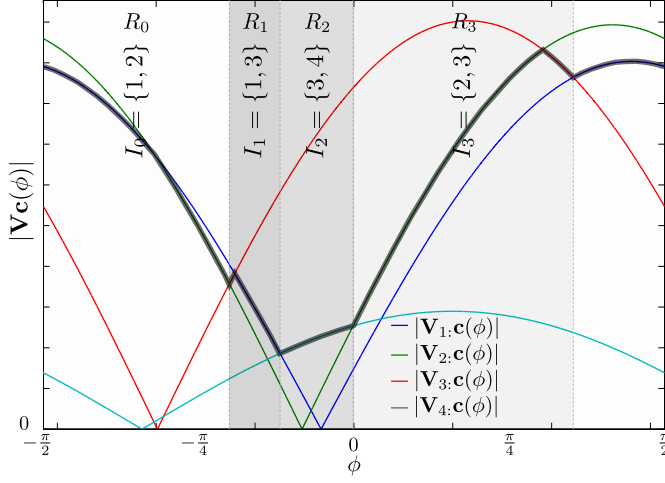
This article has been accepted for publication in a future issue of IEEE Transactions on Information Theory, but has not been fully edited. Content may change prior to final publication.

8



Fig. 3. Execution of Algorithm 2 on an arbitrary $4 \times 2$ matrix $\mathbf{V}$ for sparsity $K = 2$.
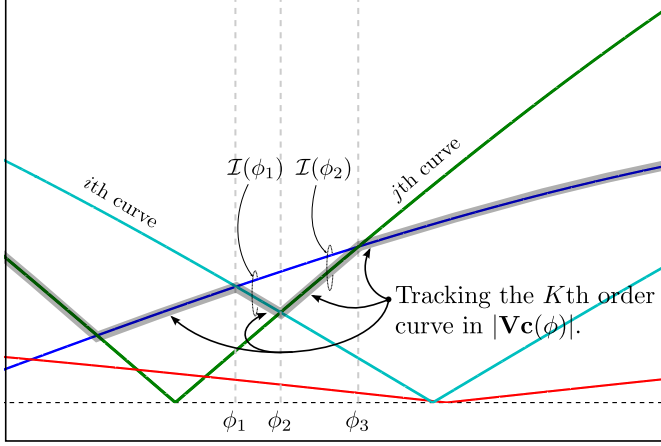


Fig. 4. Curves in $|\mathbf{u}(\phi)|$ in a constrained region of $\Phi$ for an arbitrary $4 \times 2$ matrix $\mathbf{V}$. We highlight the evolution of the $K$th order curve, for $K = 2$. Angles $\phi_1$, $\phi_2$, and $\phi_3$ are three consecutive intersection points. The candidate support sets in $[\phi_1, \phi_2)$ and $[\phi_2, \phi_3)$ differ by exactly one element. The $K$th order curve changes over the intersection point $\phi_2$.

curve, highlighted with a black solid line. Vertical dashed lines indicate the intersections points at which the support of the optimal $\mathbf{x}$ changes. Regions of $\Phi$ corresponding to different support have a different background color. The optimal support in each region is depicted in the top of the figure. Note that 2 of the $\binom{N}{2} = 6$ support sets need not be considered: sets $\{1, 4\}$ and $\{2, 4\}$ do not appear for any $\phi \in \Phi$.

A pseudocode that implements the above serial algorithmic steps is presented in Algorithm 2. It improves upon the computational complexity of its parallel counterpart (Algorithm 1) circumventing the $O(N)$ construction of the candidate set at each individual intersection point. The computation and sorting of all intersection points (steps 1 and 2) is performed in $O(N^2 \log N)$ operations. The construction of $\mathcal{I}(\phi_1)$ is performed in linear time. Finally, each successive update in the last step requires $O(1)$ operations. Therefore, the overall complexity of the serial Algorithm 2 is $O(N^2 \log N)$ as opposed to the complexity $O(N^3)$ of the parallel Algorithm 1 for $D = 2$. The disadvantage of Algorithm 2 is that it is not

---

**Algorithm 2** Serial computation of the sparse principal component of a rank-2 matrix with complexity $\mathcal{O}(N^2 \log N)$

---

**Input:** $\mathbf{V} \in \mathbb{R}^{N \times 2}$, $K$

  $\mathcal{S} \leftarrow \{\}$ *(set of candidate supports)*

  $\phi \leftarrow$ sort $\left(\text{all } 2\binom{N}{2} \text{ intersection points}\right)$

  *(remember the pair $\{i, j\}$ of curves intersecting at each $\phi_t$)*

  $\mathcal{I}_0 \leftarrow \text{top}_K(\mathbf{Vc}(\phi_0))$ at $\phi_0 = -\frac{\pi}{2}$

  **for** $t = 1 : \text{length}(\phi)$ **do**

    $\{i, j\} \leftarrow$ curves intersecting at $\phi_t$

    Compute $\mathcal{I}_t$ by modifying $\mathcal{I}_{t-1}$ in $\mathcal{O}(1)$:

    **if** both or none of $i$ and $j$ are in $\mathcal{I}_{t-1}$, **then**

      $\mathcal{I}_t \leftarrow \mathcal{I}_{t-1}$

    **else if** only $i$ (or $j$) is in $\mathcal{I}_{t-1}$, **then**

      $\mathcal{I}_t \leftarrow (\mathcal{I}_{t-1} - \{i \text{ (or } j)\}) \cup \{j \text{ (or } i)\}$

    **end if**

    $\mathcal{S} \leftarrow \mathcal{S} \cup \{\mathcal{I}_t\}$

  **end for**

  $\mathcal{I}_{\text{opt}} = \underset{\mathcal{I} \in \mathcal{S}}{\arg \max}\, \sigma_{\max}(\mathbf{V}_{\mathcal{I}, :})$

  **return** $\mathcal{I}_{\text{opt}}$ & principal left singular vector of $\mathbf{V}_{\mathcal{I}_{\text{opt}}, :}$

---

parallelizable.

### B. A modified serial algorithm for rank-2 matrices

Algorithms 1 and 2 described so far collect the candidate support sets by examining all pairwise intersection points. The number of distinct candidate support sets, however, can be significantly less: multiple intervals on $\Phi$ may be associated with the same support set. In the following, we describe a simple modification of Algorithm 2 that aims at reducing the total number of intersection points computed and examined. Although the worst case complexity remains the same, the modified serial algorithm may significantly speed up execution in certain cases.

Consider three consecutive intersection points $\phi_1$, $\phi_2$, and $\phi_3$, such that the candidate support sets on the two sides of $\phi_2$ are different (see, for example, Fig. 4). Let $\mathcal{I}(\phi_1)$ denote the set of indices of the $K$ largest curves in $|\mathbf{u}(\phi)|$ in the interval $[\phi_1, \phi_2)$. Similarly, $\mathcal{I}(\phi_2)$ is the set associated with $[\phi_2, \phi_3)$. The two sets differ by one element: if $\mathcal{C} = \mathcal{I}(\phi_1) \cap \mathcal{I}(\phi_2)$ is the set of the $K - 1$ common elements, then $\mathcal{I}(\phi_1) = \mathcal{C} \cup \{i\}$ and $\mathcal{I}(\phi_2) = \mathcal{C} \cup \{j\}$, for some curves $i$ and $j$.

At $\phi_2$, over which the candidate support set changes, the two curves intersect. The $j$th curve joins the set of $K$ largest curves of $|\mathbf{u}(\phi)|$, displacing the $i$th curve which was a member of $\mathcal{I}(\phi_1)$. In particular, the $i$th curve must be the smallest element of $|\mathbf{u}(\phi)|$ in $[\phi_1, \phi_2)$. To see that, assume for the sake of contradiction that among the curves in $\mathcal{I}(\phi_1)$, the $l$th curve was the smallest one in $[\phi_1, \phi_2)$, where $l \neq i$. Then, $u_i(\phi) > u_l(\phi) > u_j(\phi)$ in $[\phi_1, \phi_2)$. By assumption, $u_i(\phi) = u_j(\phi)$ at $\phi_2$. Due to the continuity of the curves, there must exist a point in $[\phi_1, \phi_2)$ where either $u_i(\phi) = u_l(\phi)$ or $u_l(\phi) = u_j(\phi)$. However, no intersection point lies in $[\phi_1, \phi_2)$. Following a similar argument, the $j$th curve must be the largest curve in $[\phi_1, \phi_2)$ among those not included in $\mathcal{I}(\phi_1)$. Moreover, the $j$th curve becomes the $K$th-order element of $|\mathbf{u}(\phi)|$ in $[\phi_2, \phi_3)$, i.e., it is the smallest curve in $\mathcal{I}(\phi_2)$.

---

**Algorithm 3** Modified serial computation of the sparse principal component of a rank-2 matrix with complexity $\mathcal{O}\left(N^2 \log N\right)$

---

**Input:** $\mathbf{V} \in \mathbb{R}^{N \times 2}$, $K$.

$\quad \phi_0 \leftarrow -\frac{\pi}{2}$, $t \leftarrow 0$

$\quad \mathcal{S} \leftarrow \{\}$ *(set of candidate supports)*

$\quad \mathcal{I}_0 \leftarrow \text{top}_K\left(\mathbf{V}\mathbf{c}(\phi_0)\right)$,

$\quad i \leftarrow$ index of the smallest magnitude curve in $\mathcal{I}_0$.

$\quad$**loop**

$\qquad t \leftarrow t + 1$.

$\qquad \mathcal{H} \leftarrow$ points of intersection of curve $i$ with all other

$\qquad\qquad$ curves in $\left(\phi_{t-1}, \frac{\pi}{2}\right)$, if not already computed

$\qquad$**if** $\mathcal{H} = \emptyset$, **then**

$\qquad\qquad$ exit loop

$\qquad$**end if**

$\qquad \phi_t \leftarrow \min\{\mathcal{H}\}$

$\qquad j \leftarrow$ index of the curve intersecting curve $i$ at $\phi_t$

$\qquad$**if** $j \in \mathcal{I}_{t-1}$, **then**

$\qquad\qquad \mathcal{I}_t \leftarrow \mathcal{I}_{t-1}$

$\qquad$**else**

$\qquad\qquad \mathcal{I}_t \leftarrow (\mathcal{I}_{t-1} - \{i\}) \cup \{j\}$

$\qquad$**end if**

$\qquad \mathcal{S} \leftarrow \mathcal{S} \cup \{\mathcal{I}_t\}$

$\qquad i \leftarrow j$ *(the new $K$th order curve)*

$\quad$**end loop**

$\quad \mathcal{I}_{\text{opt}} = \arg\max_{\mathcal{I} \in \mathcal{S}} \sigma_{\max}(\mathbf{V}_{\mathcal{I},:})$

$\quad$**return** $\mathcal{I}_{\text{opt}}$ & principal left singular vector of $\mathbf{V}_{\mathcal{I}_{\text{opt}},:}$

---

The key observation is that, along $\Phi$, the candidate support set $\mathcal{I}(\phi)$ changes only at the intersection points of the $K$th-order curve in $|\mathbf{u}(\phi)|$. Assume that the candidate support set $\mathcal{I}(\phi_0)$ is known at $\phi_0 = -\frac{\pi}{2}$ and the $i$th curve is the $K$th-order curve in $|\mathbf{u}(\phi)|$, i.e., the smallest curve in $\mathcal{I}(\phi_0)$. Moving along $\Phi$, the first point where the candidate support set can potentially change is the closest intersection point of the $i$th curve. Let $\phi_1 \in [\phi_0, \frac{\pi}{2})$ be that point and $j$ be the intersecting curve. If $j \notin \mathcal{I}(\phi_0)$, then the $j$th curve joins the set $\mathcal{I}(\phi_1)$ at $\phi_1$, displacing the $i$th curve. If $j \in \mathcal{I}(\phi_0)$, then $\phi_1$ is only a point of internal reordering for $\mathcal{I}(\phi_0)$, hence $\mathcal{I}(\phi_1) = \mathcal{I}(\phi_0)$. In either case, however, the $j$th curve becomes the $K$th-order curve immediately after $\phi_1$. Proceeding in a similar fashion, the next point where the candidate support set can potentially change is a point $\phi_2 \in [\phi_1, \frac{\pi}{2})$ closest to $\phi_1$, where the $j$th curve intersects one of the other $N$ curves.

A pseudocode that implements the above steps and modifications is presented in Algorithm 3. In summary, instead of computing all pairwise intersections at the first step, Algorithm 3 postpones the computation of the intersection points of the $i$th curve until the latter becomes the $K$th-order curve in $|\mathbf{u}(\phi)|$. The motivation behind this enhancement lies on the fact that multiple curves might never become $K$th in order.

To justify this, in Fig. 5, we present the average number of intersection points computed[10] by Algorithm 3 as a function of $N$, for sparsity $K = \sqrt{N}$, $5 \log N$, and 20. The average

---

[10]Although even fewer points will be eventually visited, the complexity of Algorithm 3 is dominated by the number of intersections computed.
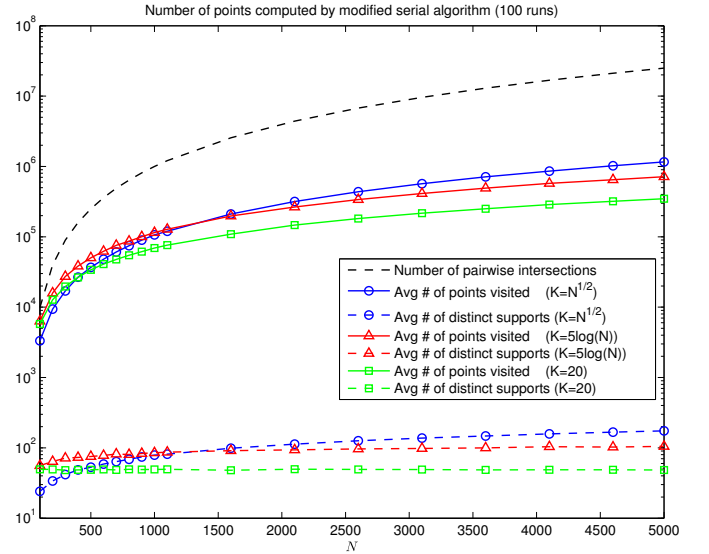


Fig. 5. Average number of intersection points and distinct candidate supports computed by Algorithm 3 for sparsity $K = \sqrt{N}$, $5 \log N$, and 20 and total number of intersection points $2\binom{N}{2}$, as a function of $N$.

is estimated over 100 independent instances of the $N \times 2$ matrix $\mathbf{V}$. We also present the average number of distinct candidate supports for the three different values of $K$ and the total number of intersection points $2\binom{N}{2}$. We observe that Algorithm 3 computes noticeably fewer points than the total number of $2\binom{N}{2}$ pairwise intersection points, indicating that a significant fraction of the curves never rises to the $K$th order. Finally, we note that the average number of distinct candidate supports, i.e., the cardinality of $\mathcal{S}$, is significantly smaller than the number of intersection points computed by Algorithm 3, revealing potential for further reduction of the computational complexity of the algorithms developed in this present work.

## V. CONCLUSIONS

We proved that the sparse principal component of an $N \times N$ matrix is computable with complexity $\mathcal{O}\left(N^{D+1}\right)$ if the matrix is positive semidefinite and its rank $D$ is constant. This holds true for any sparsity value $K$ (that is, even if $K$ grows with $N$). Our constructive proof was accompanied by a fully parallelizable and memory efficient algorithm which computes the sparse principal component with complexity $\mathcal{O}\left(N^{D+1}\right)$. For the special case of rank-2 matrices, we developed alternative serial algorithms of complexity $\mathcal{O}\left(N^2 \log N\right)$. The construction steps and properties of the algorithms that we presented in this work indicate that implementations of even lower complexity may be possible.

## REFERENCES

[1] D. L. Donoho, M. Elad, and V. N. Temlyakov, "Stable recovery of sparse overcomplete representations in the presence of noise," *IEEE Trans. Inf. Theory*, vol. 52, pp. 6-18, Jan. 2006.

[2] J. A. Tropp, "Just relax: Convex programming methods for identifying sparse signals in noise," *IEEE Trans. Inf. Theory*, vol. 52, pp. 1030-1051, Mar. 2006.

[3] H. Zou, T. Hastie, and R. Tibshirani, "Sparse principal component analysis," *J. Computational and Graphical Statistics*, vol. 15, pp. 265-286, 2006.

[4] B. Moghaddam, Y. Weiss, and S. Avidan, "Spectral bounds for sparse PCA: Exact and greedy algorithms," *Advances in Neural Information Processing Systems*, vol. 18, pp. 915-922, 2006.

[5] A. d' Aspremont, L. El Ghaoui, M. I. Jordan, and G. R. G. Lanckriet, "A direct formulation for sparse PCA using semidefinite programming," *SIAM Review*, vol. 49, pp. 434-448, 2007.

[6] A. A. Amini and M. J. Wainwright, "High-dimensional analysis of semidefinite relaxations for sparse principal components," in *Proc. IEEE ISIT 2008*, Toronto, ON, July 2008, pp. 2454-2458.

[7] E. Diederichs, A. Juditsky, V. Spokoiny, and C. Schütte, "Sparse non-Gaussian component analysis," *IEEE Trans. Inf. Theory*, vol. 56, pp. 3033-3047, June 2010.

[8] C. Shen, S. Paisitkriangkrai, and J. Zhang, "Efficiently learning a detection cascade with sparse eigenvectors," *IEEE Trans. Image Process.*, vol. 20, pp. 22-35, Jan. 2011.

[9] M. O. Ulfarsson and V. Solo, "Vector $l_0$ sparse variable PCA," *IEEE Trans. Signal. Process.*, vol. 59, pp. 1949-1958, May 2011.

[10] C. Boutsidis, P. Drineas, and M. Magdon-Ismail, "Sparse features for PCA-like linear regression," *Advances in Neural Information Processing Systems*, vol. 24, pp. 2285-2293, 2011.

[11] N. Singh, B. A. Miller, N. T. Bliss, and P. J. Wolfe, "Anomalous subgraph detection via sparse principal component analysis," in *Proc. IEEE SSP 2011*, Nice, France, June 2011, pp. 485-488.

[12] I. D. Schizas and G. B. Giannakis, "Covariance eigenvector sparsity for compression and denoising," *IEEE Trans. Signal. Process.*, vol. 60, pp. 2408-2421, May 2012.

[13] D. Wei, C. K. Sestok, and A. V. Oppenheim, "Sparse filter design under a quadratic constraint: Low-complexity algorithms," *IEEE Trans. Signal Process.*, vol. 61, pp. 857-870, Feb. 2013.

[14] M. Elad and I. Yavneh, "A plurality of sparse representations is better than the sparsest one alone," *IEEE Trans. Inf. Theory*, vol. 55, pp. 4701-4714, Oct. 2009.

[15] C. G. Tsinos, A. S. Lalos, and K. Berberidis, "Sparse subspace tracking techniques for adaptive blind channel identification in OFDM systems," in *Proc. IEEE ICASSP 2012*, Kyoto, Japan, Mar. 2012, pp. 3185-3188.

[16] B. Moghaddam, Y. Weiss, and S. Avidan, "Generalized spectral bounds for sparse LDA," in *Proc. 23rd Intern. Conf. Machine Learning*, Pittsburgh, PA, June 2006, pp. 641-648.

[17] B. K. Sriperumbudur, D. A. Torres, and G. R. G. Lanckriet, "Sparse eigen methods by DC programming," in *Proc. 24th Intern. Conf. Machine Learning*, Corvallis, OR, June 2007, pp. 831-838.

[18] A. d' Aspremont, F. Bach, and L. El Ghaoui, "Optimal solutions for sparse principal component analysis," *J. Machine Learning Research*, vol. 9, pp. 1269-1294, July 2008.

[19] L. Mackey, "Deflation methods for sparse PCA," *Advances in Neural Information Processing Systems*, vol. 21, pp. 1017-1024, 2009.

[20] M. Journée, Y. Nesterov, P. Richtárik, and R. Sepulchre, "Generalized power method for sparse principal component analysis," *J. Machine Learning Research*, vol. 11, pp. 517-553, Feb. 2010.

[21] Y. Zhang and L. El Ghaoui, "Large-scale sparse principal component analysis with application to text data," *Advances in Neural Information Processing Systems*, vol. 24, pp. 532-539, 2011.

[22] Y. Zhang, A. d' Aspremont, and L. El Ghaoui, "Sparse PCA: Convex relaxations, algorithms and applications," in *Handbook on Semidefinite, Cone and Polynomial Optimization*, vol. 166, pp. 915-940, 2012.

[23] M. Grbovic, C. R. Dance, and S. Vucetic, "Sparse principal component analysis with constraints," in *Proc. 26th AAAI Conf. Artificial Intelligence*, Toronto, ON, July 2012, pp. 935-941.

[24] T. Jolliffe, N. T. Trendafilov, and M. Uddin, "A modified principal component technique based on the LASSO," *J. Computational and Graphical Statistics*, vol. 12, pp. 531-547, 2003.

[25] Z. Ma, "Sparse principal component analysis and iterative thresholding," *Ann. Statist.*, vol. 41, pp. 772-801, Apr. 2013.

[26] X.-T. Yuan and T. Zhang, "Truncated power method for sparse eigenvalue problems," *J. Machine Learning Research*, vol. 14, pp. 899-925, Apr. 2013.

[27] G. N. Karystinos and D. A. Pados, "Rank-2-optimal adaptive design of binary spreading codes," *IEEE Trans. Inf. Theory*, vol. 53, pp. 3075-3080, Sept. 2007.

[28] G. N. Karystinos and A. P. Liavas, "Efficient computation of the binary vector that maximizes a rank-deficient quadratic form," *IEEE Trans. Inf. Theory*, vol. 56, pp. 3581-3593, July 2010.

[29] K. M. Mackenthun, Jr., "A fast algorithm for multiple-symbol differential detection of MPSK," *IEEE Trans. Commun.*, vol. 42, pp. 1471-1474, Feb./Mar./Apr. 1994.

[30] W. Sweldens, "Fast block noncoherent decoding," *IEEE Commun. Lett.*, vol. 5, pp. 132-134, Apr. 2001.

[31] I. Motedayen, A. Krishnamoorthy, and A. Anastasopoulos, "Optimal joint detection/estimation in fading channels with polynomial complexity," *IEEE Trans. Inf. Theory*, vol. 53, pp. 209-223, Jan. 2007.

[32] D. S. Papailiopoulos and G. N. Karystinos, "Maximum-likelihood noncoherent OSTBC detection with polynomial complexity," *IEEE Trans. Wireless Commun.*, vol. 9, pp. 1935-1945, June 2010.

[33] D. S. Papailiopoulos, G. Abou Elkheir, and G. N. Karystinos, "Maximum-likelihood noncoherent PAM detection," *IEEE Trans. Commun.*, vol. 61, pp. 1152-1159, Mar. 2013.

[34] A. Kyrillidis and G. N. Karystinos, "Fixed-rank Rayleigh quotient maximization by an $M$PSK sequence," submitted to *IEEE Trans. Commun.*, June 2013.

[35] P. P. Markopoulos, G. N. Karystinos, and D. A. Pados, "Some options for $L_1$-subspace signal processing," in *Proc. IEEE ISWCS 2013*, Ilmenau, Germany, Aug. 2013, pp. 622-626.

[36] M. Asteris, "Sparse rank-deficient variance maximization," Diploma Thesis, Department of ECE, Technical Univ. Crete, July 2010.

[37] M. Asteris, D. S. Papailiopoulos, and G. N. Karystinos, "Sparse principal component of a rank-deficient matrix," in *Proc. IEEE ISIT 2011*, Saint Petersburg, Russia, Aug. 2011, pp. 673-677.

[38] T. H Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein, *Introduction to Algorithms*, 2nd ed. Cambridge, MA: MIT Press, 2001.

**Megasthenis Asteris** (S'12) received the Diploma degree (five-year program) in electronic and computer engineering from the Technical University of Crete, Chania, Greece, in 2010 and the M.Sc. degree in electrical engineering from the University of Southern California, Los Angeles, in 2012. He is currently working toward the Ph.D. degree at the Department of Electrical and Computer Engineering, University of Texas at Austin. His research interests are in coding theory with an emphasis on distributed storage and large-scale data processing.

**Dimitris S. Papailiopoulos** (S'10) received the Diploma (five-year program) and M.Sc. degrees in electronic and computer engineering from the Technical University of Crete, Greece, in 2007 and 2009, respectively. Between 2009 and 2012, he was a graduate student with the Department of Electrical Engineering, University of Southern California, Los Angeles. Since January 2013, he has been with the Department of Electrical and Computer Engineering, University of Texas at Austin, where he is currently working towards the Ph.D. degree. His research interests are in the areas of coding theory, information theory, and large-scale data processing with an emphasis on codes for distributed storage, large-scale optimization for data and graph analysis, and succinct data representations.

**George N. Karystinos** (S'98-M'03) was born in Athens, Greece, on April 12, 1974. He received the Diploma degree in computer science and engineering (five-year program) from the University of Patras, Patras, Greece, in 1997 and the Ph.D. degree in electrical engineering from the State University of New York at Buffalo, Amherst, NY, in 2003.

From 2003 to 2005, he held an Assistant Professor position in the Department of Electrical Engineering, Wright State University, Dayton, OH. Since 2005, he has been with the Department of Electronic and Computer Engineering, Technical University of Crete, Chania, Greece, where he is presently an Associate Professor. His research interests are in the general areas of $L_1$-norm principal component analysis of data and signals, communication theory, and adaptive signal processing with applications to signal waveform design, low-complexity sequence detection, and secure wireless communications.

Dr. Karystinos is a member of the IEEE Signal Processing, Communications, Information Theory, and Computational Intelligence Societies. For articles that he coauthored with students and colleagues, he received a 2001 IEEE International Conference on Telecommunications best paper award, the 2003 IEEE TRANSACTIONS ON NEURAL NETWORKS Outstanding Paper Award, the 2011 IEEE International Conference on RFID-Technologies and Applications Second Best Student Paper Award, and the 2013 International Symposium on Wireless Communication Systems Best Paper Award in Signal Processing and Physical Layer Communications.