
Nonnegative Sparse PCA with Provable Guarantees

Megasthenis Asteris
Dimitris S. Papailiopoulos
Alexandros G. Dimakis

MEGAS@UTEXAS.EDU
DIMITRIS@UTEXAS.EDU
DIMAKIS@AUSTIN.UTEXAS.EDU

Department of Electrical and Computer Engineering, The University of Texas at Austin, TX, USA

Abstract

We introduce a novel algorithm to compute nonnegative sparse principal components of positive semidefinite (PSD) matrices. Our algorithm comes with approximation guarantees contingent on the spectral profile of the input matrix \mathbf{A} : the sharper the eigenvalue decay, the better the quality of the approximation.

If the eigenvalues decay like any asymptotically vanishing function, we can approximate nonnegative sparse PCA within any accuracy ϵ in time polynomial in the matrix dimension n and desired sparsity k , but not in $1/\epsilon$. Further, we obtain a data dependent bound that is computed by executing an algorithm on a given data set. This bound is significantly tighter than *a-priori* bounds and can be used to show that for all tested datasets our algorithm is provably within 40% – 90% from the unknown optimum.

Our algorithm is combinatorial and explores a subspace defined by the leading eigenvectors of \mathbf{A} . We test our scheme on several data sets, showing that it matches or outperforms the previous state of the art.

1. Introduction

Given a data matrix $\mathbf{S} \in \mathbb{R}^{n \times m}$ comprising m zero-mean vectors on n features, the first principal component (PC) is

$$\arg \max_{\|\mathbf{x}\|_2=1} \mathbf{x}^T \mathbf{A} \mathbf{x}, \quad (1)$$

where $\mathbf{A} = 1/m \cdot \mathbf{S} \mathbf{S}^T$ is the $n \times n$ positive semidefinite (PSD) empirical covariance matrix. Subsequent PCs can be computed after \mathbf{A} has been appropriately deflated to remove the first eigenvector. PCA is arguably the workhorse

of high dimensional data analysis and achieves dimensionality reduction by computing the directions of maximum variance. Typically, all n features affect positively or negatively these directions resulting in dense PCs, which explain the largest possible data variance, but are often not interpretable.

It has been shown that enforcing nonnegativity on the computed principal components can aid interpretability. This is particularly true in applications where features interact only in an additive manner. For instance, in bioinformatics, chemical concentrations are nonnegative (Kim & Park, 2007), or the expression level of genes is typically attributed to positive or negative influences of those genes, but not both (Badea & Tilivea, 2005). Here, enforcing nonnegativity, in conjunction with sparsity on the computed components can assist the discovery of local patterns in the data. In computer vision, where features may coincide with non negatively valued image pixels, nonnegative sparse PCA pertains to the extraction of the most informative image parts (Lee & Seung, 1999). In other applications, nonnegative weights admit a meaningful probabilistic interpretation.

Sparsity emerges as an additional desirable trait of the computed components because it further helps interpretability (Zou et al., 2006; d’Aspremont et al., 2007b), even independently of nonnegativity. From a machine learning perspective, enforcing sparsity serves as an unsupervised feature selection method: the active coordinates in an optimal l_0 -norm constrained PC should correspond to the most informative subset of features. Although nonnegativity inherently promotes sparsity, an explicit sparsity constraint enables precise control on the number of selected features.

Nonnegative Sparse PC. Nonnegativity and sparsity can be directly enforced on the principal component optimization by adding constraints to (1). The k -sparse nonnegative principal component of \mathbf{A} is

$$\mathbf{x}_* = \arg \max_{\mathbf{x} \in \mathbb{S}_k^n} \mathbf{x}^T \mathbf{A} \mathbf{x}, \quad (2)$$

where $\mathbb{S}_k^n = \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\|_2 = 1, \|\mathbf{x}\|_0 \leq k, \mathbf{x} \geq 0\}$, for a desired sparsity parameter $k \in [n]$.

The problem of computing the first eigenvector (1) is easily solvable, but with the additional sparsity and nonnegativity constraints problem (2) becomes computationally intractable. The cardinality constraint alone renders sparse PCA NP-hard (Moghaddam et al., 2006b). Even if the l_0 -norm constraint is dropped, we show that problem (2) remains computationally intractable by reducing it to checking matrix copositivity, a well known co-NP complete decision problem (Murty & Kabadi, 1987; Parrilo, 2000). Therefore, each of the constraints $\mathbf{x} \geq \mathbf{0}$ and $\|\mathbf{x}\|_0 \leq k$ individually makes the problem intractable.

Our Contribution: We introduce a novel algorithm for approximating the nonnegative k -sparse principal component with provable approximation guarantees.

Given any PSD matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$, sparsity parameter k , and accuracy parameter $d \in [n]$, our algorithm outputs a nonnegative, k -sparse, unit norm vector \mathbf{x}_d that achieves at least ρ_d fraction of the maximum objective value in (2), *i.e.*,

$$\mathbf{x}_d^T \mathbf{A} \mathbf{x}_d \geq \rho_d \cdot \mathbf{x}_*^T \mathbf{A} \mathbf{x}_*, \quad (3)$$

where

$$\rho_d \geq \max \left\{ \frac{k}{2n}, \frac{1}{1 + 2\frac{n}{k} \lambda_{d+1} / \lambda_1} \right\}. \quad (4)$$

Here, λ_i is the i^{th} largest eigenvalue of \mathbf{A} , and the accuracy parameter d specifies the rank of the approximation used and controls the running time. Specifically, our algorithm runs in time $O(n^d k^d + n^{d+1})$. As can be seen our result depends on the spectral profile of \mathbf{A} : the faster the eigenvalue decay, the tighter the approximation.

Near-Linear time approximation. Our algorithm has a running time $O(n^d k^d + n^{d+1})$, which in the linear sparsity regime can be as high as $O(n^{2d})$. This can be non-practical for large data sets, even if we set the rank parameter d to be two or three. We present a modification of our algorithm that can provably approximate the result of the first in near-linear time. Specifically, for any desired accuracy $\epsilon \in (0, 1]$ it computes a nonnegative, k -sparse, unit norm vector $\hat{\mathbf{x}}_d$ such that

$$\hat{\mathbf{x}}_d^T \mathbf{A} \hat{\mathbf{x}}_d \geq (1 - \epsilon) \cdot \rho_d \cdot \mathbf{x}_*^T \mathbf{A} \mathbf{x}_*, \quad (5)$$

where ρ_d is as described in (4). We show that the running time of our approximate algorithm is $O(\epsilon^{-d} \cdot n \log n)$, which is near-linear in n for any fixed accuracy parameters d and ϵ .

Our approximation theorem has several implications.

Exact solution for low-rank matrices. Observe that if the matrix \mathbf{A} has rank d , our algorithm returns the optimal k -sparse PC for any target sparsity k . The same holds in

the case of the rank- d update matrix $\mathbf{A} = \sigma \mathbf{I} + \mathbf{C}$, with $\text{rank}(\mathbf{C}) = d$ and arbitrary constant σ , since the algorithm can be equivalently applied on \mathbf{C} .

PTAS for any spectral decay. Consider the linear sparsity regime $k = c \cdot n$ and assume that the eigenvalues follow a decay law $\lambda_i \leq \lambda_1 \cdot f(i)$ for any decay function $f(i)$ which vanishes: $f(i) \rightarrow 0$ as $i \rightarrow \infty$. Special cases include power law decay $f(i) = 1/i^\alpha$ or even very slow decay functions like $f(i) = 1/\log \log i$. For all these cases, we can solve nonnegative sparse PCA for any desired accuracy ϵ in time polynomial in n and k , but not in $1/\epsilon$. Therefore, we obtain a polynomial-time approximation scheme (PTAS) for any spectral decay behavior.

Computable upper bounds. In addition to these theoretical guarantees, our method yields a data dependent upper bound on the maximum value of (2), that can be computed by running our algorithm. As it can be seen in Fig. 4-6, the obtained upper bound, combined with our achievable point, sandwiches the unknown optimum within a narrow region. Using this upper bound we are able to show that our solutions are within 40 – 90% from the optimal in all the datasets that we examine. To the best of our knowledge, this framework of data dependent bounds has not been considered in the previous literature.

1.1. Related Work

There is a substantial volume of work on sparse PCA, spanning a rich variety of approaches: from early heuristics in (Jolliffe, 1995), to the LASSO based techniques in (Jolliffe et al., 2003), the elastic net l_1 -regression in (Zou et al., 2006), a greedy branch-and-bound technique in (Moghaddam et al., 2006a), or semidefinite programming approaches (d’Aspremont et al., 2008; Zhang et al., 2012; d’Aspremont et al., 2007a). This line of work does not consider or enforce nonnegativity constraints.

When nonnegative components are desired, fundamentally different approaches have been used. Nonnegative matrix factorization (Lee & Seung, 1999) and its sparse variants (Hoyer, 2004; Kim & Park, 2007) fall within that scope: data is expressed as (sparse) nonnegative linear combinations of (sparse) nonnegative parts. These approaches are interested in finding a lower dimensionality representation of the data that reveals latent structure and minimizes a reconstruction error, but are not explicitly concerned with the statistical significance of individual output vectors.

Nonnegativity as an additional constraint on (sparse) PCA first appeared in (Zass & Shashua, 2007). The authors suggested a coordinate-descent scheme that jointly computes a set of nonnegative sparse principal components, maximizing the cumulative explained variance. An l_1 -penalty promotes sparsity of computed components on average,

but not on each component individually. A second convex penalty is incorporated to favor orthogonal components.

Similar convex optimization approaches for nonnegative PCA have been subsequently proposed in the literature. In (Allen & Maletić-Savatić, 2011) for instance, the authors suggest an alternating maximization scheme for the computation of the first nonnegative PC, allowing the incorporation of known structural dependencies.

A competitive algorithm for nonnegative sparse PCA was established in (Sigg & Buhmann, 2008), with the development of a framework stemming from Expectation-Maximization (EM) for a probabilistic generative model of PCA. The proposed algorithm, which enforces hard sparsity, or nonnegativity, or both constraints simultaneously, computes the first approximate PC in $O(n^2)$, *i.e.*, time quadratic in the number of features.

To the best of our knowledge, no prior works provide provable approximation guarantees for the nonnegative sparse PCA optimization problem. Further, no data dependent upper bounds have been present in the previous literature.

Differences from SPCA work. Our work is closely related to (Karystinos & Liavas, 2010; Asteris et al., 2011; Papailiopoulos et al., 2013) that introduced the ideas of solving low-rank quadratic combinatorial optimization problems on low-rank PSD matrices using hyperspectral transformations. Such transformations are called spannograms and follow a similar architecture. In this paper, we extend the spannogram framework to nonnegative sparse PCA. The most important technical issue compared to (Asteris et al., 2011; Papailiopoulos et al., 2013) is introducing nonnegativity constraints in spannogram algorithms.

To understand how this changes the problem, notice that in the original sparse PCA problem without nonnegativity constraints, if the support is known, the optimal principal component supported on that set can be easily found. However, under nonnegativity constraints, the problem is hard even if the optimal support is known. This is the fundamental technical problem that we address in this paper. We show that if the involved subspace is low-dimensional, it is possible to solve this problem.

2. Algorithm Overview

Given an $n \times n$ PSD matrix \mathbf{A} , the desired sparsity k , and an accuracy parameter $d \in [n]$, our algorithm computes a *nonnegative, k -sparse, unit norm* vector \mathbf{x}_d approximating the nonnegative, k -sparse PC of \mathbf{A} . We begin with a high-level description of the main steps of the algorithm.

Step 1. Compute \mathbf{A}_d , the rank- d approximation of \mathbf{A} . We compute \mathbf{A}_d , the best rank- d approximation of \mathbf{A} , zeroing

Algorithm 1 Spannogram Nonnegative Sparse PCA

input \mathbf{A} ($n \times n$ PSD matrix), $k \in [n]$, $d \in [n]$.
 1: $\mathbf{U}, \mathbf{\Lambda} \leftarrow \text{svd}(\mathbf{A}, d)$
 2: $\mathbf{V} = \mathbf{U}\mathbf{\Lambda}^{1/2}$ { $\mathbf{A}_d = \mathbf{V}\mathbf{V}^T$ }
 3: $\mathcal{S}_d \leftarrow \text{Spannogram}(\mathbf{V}, k)$ {Algo. 2}
 4: $\mathcal{X}_d \leftarrow \{\}$ { $|\mathcal{S}_d| \leq O(n^d)$ }
 5: **for all** $\mathcal{I} \in \mathcal{S}_d$ **do**
 6: $\mathbf{c}^{(\mathcal{I})} \leftarrow \arg \max_{\substack{\|\mathbf{c}\|_2=1 \\ \mathbf{V}_{\mathcal{I}}\mathbf{c} \geq \mathbf{0}}} \|\mathbf{V}_{\mathcal{I}}\mathbf{c}\|_2^2$ {Sec. 5}
 7: $\mathbf{x}_{\mathcal{I}}^{(\mathcal{I})} \leftarrow |\mathbf{V}_{\mathcal{I}}\mathbf{c}| / \|\mathbf{V}_{\mathcal{I}}\mathbf{c}\|$, $\mathbf{x}_{\mathcal{I}^c}^{(\mathcal{I})} \leftarrow \mathbf{0}$
 8: $\mathcal{X}_d \leftarrow \mathcal{X}_d \cup \{\mathbf{x}^{(\mathcal{I})}\}$
 9: **end for** { $|\mathcal{X}_d| \leq |\mathcal{S}_d|$ }
output $\mathbf{x}_d \leftarrow \arg \max_{\mathbf{x} \in \mathcal{X}_d} \mathbf{x}^T \mathbf{A}_d \mathbf{x}$

out the $n - d$ trailing eigenvalues of \mathbf{A} , that is,

$$\mathbf{A}_d = \sum_{i=1}^d \lambda_i \mathbf{u}_i \mathbf{u}_i^T,$$

where λ_i is the i^{th} largest eigenvalue of \mathbf{A} and \mathbf{u}_i the corresponding eigenvector.

Step 2. Compute \mathcal{S}_d , a set of $O(n^d)$ candidate supports. Enumerating the $\binom{n}{k}$ possible supports for k -sparse vectors in \mathbb{R}^n is computationally intractable. Using our *Spannogram* technique described in Section 4, we efficiently determine a collection \mathcal{S}_d of support sets, with cardinality $|\mathcal{S}_d| \leq 2^d \binom{n+1}{d}$, that provably contains the support of the nonnegative, k -sparse PC of \mathbf{A}_d .

Step 3. Compute \mathcal{X}_d , a set of candidate solutions. For each candidate support set $\mathcal{I} \in \mathcal{S}_d$, we compute a candidate solution \mathbf{x} supported only in \mathcal{I} :

$$\arg \max_{\substack{\|\mathbf{x}\|_2=1, \mathbf{x} \geq \mathbf{0}, \\ \text{supp}(\mathbf{x}) \subseteq \mathcal{I}}} \mathbf{x}^T \mathbf{A}_d \mathbf{x}. \quad (6)$$

The constant rank of \mathbf{A}_d is essential in solving (6): the constrained quadratic maximization is in general NP-hard, even for a given support.

Step 4. Output the best candidate solution in \mathcal{X}_d , *i.e.*, the candidate that maximizes the quadratic form.

If multiple components are desired, the procedure is repeated after an appropriate deflation has been applied on \mathbf{A}_d (Mackey, 2008). The steps are formally presented in Algorithm 1. A detailed description is the subject of subsequent sections.

2.1. Approximation Guarantees

Instead of the nonnegative, k -sparse, principal component \mathbf{x}_* of \mathbf{A} , which attains the optimal value $\text{OPT} = \mathbf{x}_*^T \mathbf{A} \mathbf{x}_*$, our algorithm outputs a nonnegative, k -sparse, unit norm vector \mathbf{x}_d . We measure the quality of \mathbf{x}_d as a surrogate of \mathbf{x}_* by the approximation factor $\mathbf{x}_d^T \mathbf{A}_d \mathbf{x}_d / \text{OPT}$. Clearly,

the approximation factor takes values in $(0, 1]$, with higher values implying tighter approximation.

Theorem 1. *For any $n \times n$ PSD matrix \mathbf{A} , sparsity parameter k , and accuracy parameter $d \in [n]$, Alg. 1 outputs a nonnegative, k -sparse, unit norm vector \mathbf{x}_d such that*

$$\mathbf{x}_d^T \mathbf{A} \mathbf{x}_d \geq \rho_d \cdot \mathbf{x}_*^T \mathbf{A} \mathbf{x}_*,$$

where

$$\rho_d \geq \max \left\{ \frac{k}{2n}, \frac{1}{1 + 2 \frac{n}{k} \lambda_{d+1} / \lambda_1} \right\},$$

in time $O(n^{d+1} + n^d k^d)$.

The approximation guarantee of Theorem 1 relies on establishing connections among the eigenvalues of \mathbf{A} , and the quadratic forms $\mathbf{x}_d^T \mathbf{A} \mathbf{x}_d$ and $\mathbf{x}_d^T \mathbf{A}_d \mathbf{x}_d$. The proof can be found in the supplemental material. The complexity of Algorithm 1 follows upon its detailed description.

3. Proposed Scheme

Our algorithm approximates the nonnegative, k -sparse PC of a PSD matrix \mathbf{A} by computing the corresponding PC of \mathbf{A}_d , a rank- d surrogate of the input argument \mathbf{A} :

$$\mathbf{A}_d = \sum_{i=1}^d \mathbf{v}_i \mathbf{v}_i^T = \mathbf{V} \mathbf{V}^T, \quad (7)$$

where $\mathbf{v}_i = \sqrt{\lambda_i} \mathbf{u}_i$ is the scaled eigenvector corresponding to the i^{th} largest eigenvalue of \mathbf{A} , and $\mathbf{V} = [\mathbf{v}_1 \cdots \mathbf{v}_d] \in \mathbb{R}^{n \times d}$. In this section, we delve into the details of our algorithmic developments and describe how the low rank of \mathbf{A}_d unlocks the computation of the desired PC.

3.1. Rank-1: A simple case

We begin with the rank-1 case because, besides its motivational simplicity, it is a fundamental component of the algorithmic developments for the rank- d case.

In the rank-1 case, \mathbf{V} reduces to a single vector in \mathbb{R}^n and \mathbf{x}_1 , the nonnegative k -sparse PC of \mathbf{A}_1 , is the solution to

$$\max_{\mathbf{x} \in \mathbb{S}_k^n} \mathbf{x}^T \mathbf{A}_1 \mathbf{x} = \max_{\mathbf{x} \in \mathbb{S}_k^n} (\mathbf{v}^T \mathbf{x})^2. \quad (8)$$

That is, \mathbf{x}_1 is the nonnegative, k -sparse, unit length vector that maximizes $(\mathbf{v}^T \mathbf{x})^2$. Let $\mathcal{I} = \text{supp}(\mathbf{x}_1)$, $|\mathcal{I}| \leq k$, be the unknown support of \mathbf{x}_1 . Then, $(\mathbf{v}^T \mathbf{x})^2 = (\sum_{i \in \mathcal{I}} v_i \cdot x_i)^2$. Since $\mathbf{x}_1 \geq \mathbf{0}$, it should not be hard to see that the active entries of \mathbf{x}_1 must correspond to nonnegative or nonpositive entries of \mathbf{v} , but not a combination of both. In other words, $\mathbf{v}_{\mathcal{I}}$, the entries of \mathbf{v} indexed by \mathcal{I} , must satisfy $\mathbf{v}_{\mathcal{I}} \geq \mathbf{0}$ or $\mathbf{v}_{\mathcal{I}} \leq \mathbf{0}$. In either case, by the Cauchy-Schwarz inequality,

$$(\mathbf{v}^T \mathbf{x})^2 = (\mathbf{v}_{\mathcal{I}}^T \mathbf{x}_{\mathcal{I}})^2 \leq \|\mathbf{v}_{\mathcal{I}}\|_2^2 \|\mathbf{x}_{\mathcal{I}}\|_2^2 = \|\mathbf{v}_{\mathcal{I}}\|_2^2. \quad (9)$$

Equality in (9) can always be achieved by setting $\mathbf{x}_{\mathcal{I}} = \mathbf{v}_{\mathcal{I}} / \|\mathbf{v}_{\mathcal{I}}\|_2$ if $\mathbf{v}_{\mathcal{I}} \geq \mathbf{0}$, and $\mathbf{x}_{\mathcal{I}} = -\mathbf{v}_{\mathcal{I}} / \|\mathbf{v}_{\mathcal{I}}\|_2$ if $\mathbf{v}_{\mathcal{I}} \leq \mathbf{0}$. The support of the optimal solution \mathbf{x}_1 is the set \mathcal{I} for which $\|\mathbf{v}_{\mathcal{I}}\|_2^2$ in (9) is maximized under the restriction that the entries of $\mathbf{v}_{\mathcal{I}}$ do not have mixed signs.

Def. 1. *Let $\mathcal{I}_k^+(\mathbf{v})$, $1 \leq k \leq n$ denote the set of indices of the (at most) k largest nonnegative entries in $\mathbf{v} \in \mathbb{R}^n$.*

Proposition 3.1. *Let \mathbf{x}_1 be the solution to problem (8). Then, $\text{supp}(\mathbf{x}_1) \in \mathcal{S}_1 = \{\mathcal{I}_k^+(\mathbf{v}), \mathcal{I}_k^+(-\mathbf{v})\}$.*

The collection \mathcal{S}_1 and the associated candidate vectors via (9) are constructed in $O(n)$. The solution \mathbf{x}_1 is the candidate that maximizes the quadratic.

3.2. Rank- d case

In the rank- d case, \mathbf{x}_d , the nonnegative, k -sparse PC of \mathbf{A}_d is the solution to the following problem:

$$\max_{\mathbf{x} \in \mathbb{S}_k^n} \mathbf{x}^T \mathbf{A}_d \mathbf{x} = \max_{\mathbf{x} \in \mathbb{S}_k^n} \|\mathbf{V}^T \mathbf{x}\|_2^2. \quad (10)$$

Consider an auxiliary vector $\mathbf{c} \in \mathbb{R}^d$, with $\|\mathbf{c}\|_2 = 1$. From the Cauchy-Schwarz inequality,

$$\|\mathbf{V}^T \mathbf{x}\|_2^2 = \|\mathbf{c}\|_2^2 \|\mathbf{V}^T \mathbf{x}\|_2^2 \geq |\mathbf{c}^T (\mathbf{V}^T \mathbf{x})|^2. \quad (11)$$

Equality in (11) is achieved if and only if \mathbf{c} is colinear to $\mathbf{V}^T \mathbf{x}$. Since \mathbf{c} spans the entire unit sphere, such a \mathbf{c} exists for every \mathbf{x} , yielding an alternative description for the objective function in (10):

$$\|\mathbf{V}^T \mathbf{x}\|_2^2 = \max_{\mathbf{c} \in \mathbb{S}^d} |(\mathbf{V} \mathbf{c})^T \mathbf{x}|^2, \quad (12)$$

where $\mathbb{S}^d = \{\mathbf{c} \in \mathbb{R}^d : \|\mathbf{c}\|_2 = 1\}$ is the d -dimensional unit sphere. The maximization in (10) becomes

$$\begin{aligned} \max_{\mathbf{x} \in \mathbb{S}_k^n} \|\mathbf{V}^T \mathbf{x}\|_2^2 &= \max_{\mathbf{x} \in \mathbb{S}_k^n} \max_{\mathbf{c} \in \mathbb{S}^d} |(\mathbf{V} \mathbf{c})^T \mathbf{x}|^2 \\ &= \max_{\mathbf{c} \in \mathbb{S}^d} \max_{\mathbf{x} \in \mathbb{S}_k^n} |(\mathbf{V} \mathbf{c})^T \mathbf{x}|^2. \end{aligned} \quad (13)$$

The set of candidate supports. A first key observation is that for fixed \mathbf{c} , the product $(\mathbf{V} \mathbf{c})$ is a vector in \mathbb{R}^n . Maximizing $|(\mathbf{V} \mathbf{c})^T \mathbf{x}|^2$ over all vectors $\mathbf{x} \in \mathbb{S}_k^n$ is a rank-1 instance of the optimization problem, as in (8). Let $(\mathbf{c}_d, \mathbf{x}_d)$ be the optimal solution of (10). By Proposition 3.1, the support of \mathbf{x}_d coincides with either $\mathcal{I}_k^+(\mathbf{V} \mathbf{c}_d)$ or $\mathcal{I}_k^+(-\mathbf{V} \mathbf{c}_d)$. Hence, we can safely claim that $\text{supp}(\mathbf{x}_d)$ appears in

$$\mathcal{S}_d = \bigcup_{\mathbf{c} \in \mathbb{S}^d} \{\mathcal{I}_k^+(\mathbf{V} \mathbf{c})\}. \quad (14)$$

Naively, one might think that \mathcal{S}_d can contain as many as $\binom{n}{k}$ distinct support sets. In Section 4, we show that $|\mathcal{S}_d| \leq 2^d \binom{n+1}{d}$ and present our Spannogram technique (Alg. 2) for efficiently constructing \mathcal{S}_d in $O(n^{d+1})$. Each support in \mathcal{S}_d corresponds to a candidate principal component.

Solving for a given support. We seek a pair (\mathbf{x}, \mathbf{c}) that maximizes (13) under the additional constraint that \mathbf{x} is supported only on a given set \mathcal{I} . By the Cauchy-Schwarz inequality, the objective in (13) satisfies

$$|(\mathbf{V}\mathbf{c})^T \mathbf{x}|^2 = |(\mathbf{V}_{\mathcal{I}}\mathbf{c})^T \mathbf{x}_{\mathcal{I}}|^2 \leq \|(\mathbf{V}_{\mathcal{I}}\mathbf{c})\|_2^2, \quad (15)$$

where $\mathbf{V}_{\mathcal{I}}$ is the matrix formed by the rows of \mathbf{V} indexed by \mathcal{I} . Equality in (15) is achieved if and only if $\mathbf{x}_{\mathcal{I}}$ is colinear to $\mathbf{V}_{\mathcal{I}}\mathbf{c}$. However, it is not achievable for arbitrary \mathbf{c} , as $\mathbf{x}_{\mathcal{I}}$ must be nonnegative. From Proposition 3.1, we infer that \mathbf{x} being supported in \mathcal{I} implies that all entries of $\mathbf{V}_{\mathcal{I}}\mathbf{c}$ have the same sign. Further, whenever the last condition holds, a nonnegative $\mathbf{x}_{\mathcal{I}}$ colinear to $\mathbf{V}_{\mathcal{I}}\mathbf{c}$ exists and equality in (15) can be achieved. Under the additional constraint that $\text{supp}(\mathbf{x}) = \mathcal{I} \in \mathcal{S}_d$, the maximization in (13) becomes

$$\max_{\mathbf{c} \in \mathbb{S}^d} \max_{\substack{\mathbf{x} \in \mathbb{S}_k^n \\ \text{supp}(\mathbf{x}) \subseteq \mathcal{I}}} |(\mathbf{V}\mathbf{c})^T \mathbf{x}|^2 = \max_{\substack{\mathbf{c} \in \mathbb{S}^d \\ \mathbf{V}_{\mathcal{I}}\mathbf{c} \geq \mathbf{0}}} \|(\mathbf{V}_{\mathcal{I}}\mathbf{c})\|_2^2. \quad (16)$$

The constraint $\mathbf{V}_{\mathcal{I}}\mathbf{c} \geq \mathbf{0}$ in (16), is equivalent to requiring that all entries in $\mathbf{V}_{\mathcal{I}}\mathbf{c}$ have the same sign, since \mathbf{c} and $-\mathbf{c}$ achieve the same objective value.

The optimization problem in (16) is NP-hard. In fact, it encompasses the original nonnegative PCA problem as a special case. Here, however, the constant dimension $d = \Theta(1)$ of the unknown variable \mathbf{c} permits otherwise intractable operations. In Section 5, we outline an $O(k^d)$ algorithm for solving this constrained quadratic maximization.

The algorithm. The previous discussion suggests a two-step algorithm for solving the rank- d optimization problem in (10). First, run the Spannogram algorithm to construct \mathcal{S}_d , the collection of $O(n^d)$ candidate supports for \mathbf{x}_d , in $O(n^{d+1})$. For each $\mathcal{I} \in \mathcal{S}_d$, solve (16) in $O(k^d)$ to obtain a candidate solution $\mathbf{x}^{(\mathcal{I})}$ supported on \mathcal{I} . Output the candidate solution that maximizes the quadratic $\mathbf{x}^T \mathbf{A}_d \mathbf{x}$. Efficiently combining the previous steps yields an $O(n^{d+1} + n^d k^d)$ procedure for approximating the nonnegative sparse PC, outlined in Alg. 1.

4. The Nonnegative Spannogram

In this section, we describe how to construct \mathcal{S}_d , the collection of candidate supports, defined in (14) as

$$\mathcal{S}_d = \bigcup_{\mathbf{c} \in \mathbb{S}^d} \{\mathcal{I}_k^+(\mathbf{V}\mathbf{c})\},$$

for a given $\mathbf{V} \in \mathbb{R}^{n \times d}$. \mathcal{S}_d comprises all support sets induced by vectors in the range of \mathbf{V} . The *Spannogram* of \mathbf{V} is a *visualization* of its range, and a valuable tool in efficiently collecting those supports.

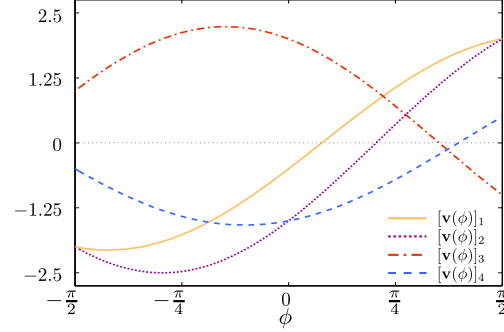


Figure 1. Spannogram of an arbitrary rank-2 matrix $\mathbf{V} \in \mathbb{R}^{4 \times 2}$. At a point ϕ , the values of the curves correspond to the entries of a vector $\mathbf{v}(\phi)$ in the range of \mathbf{V} and vice versa.

4.1. Constructing \mathcal{S}_2

We describe the $d = 2$ case, the simplest nontrivial case, to facilitate a gentle exposure to the Spannogram technique. The core ideas generalize to arbitrary d and a detailed description is provided in the supplemental material.

Spherical variables. Up to scaling, all vectors \mathbf{v} in the range of $\mathbf{V} \in \mathbb{R}^{n \times 2}$, $\mathcal{R}(\mathbf{V})$, can be written as $\mathbf{v} = \mathbf{V}\mathbf{c}$ for some $\mathbf{c} \in \mathbb{R}^2 : \|\mathbf{c}\| = 1$. We introduce a variable $\phi \in \Phi = (-\pi/2, \pi/2]$, and set \mathbf{c} to be the following function of ϕ :

$$\mathbf{c}(\phi) = [\sin(\phi) \quad \cos(\phi)]^T.$$

The range of \mathbf{V} , $\mathcal{R}(\mathbf{V}) = \{\pm \mathbf{v}(\phi) = \pm \mathbf{V}\mathbf{c}(\phi), \phi \in \Phi\}$, is also a function of ϕ , and in turn \mathcal{S}_2 can be expressed as

$$\mathcal{S}_2 = \bigcup_{\phi \in \Phi} \{\mathcal{I}_k^+(\mathbf{v}(\phi)), \mathcal{I}_k^+(-\mathbf{v}(\phi))\}.$$

Spannogram. The i^{th} entry of $\mathbf{v}(\phi)$ is a continuous function of ϕ generated by the i^{th} row of \mathbf{V} : $[\mathbf{v}(\phi)]_i = \mathbf{V}_{i,1} \sin(\phi) + \mathbf{V}_{i,2} \cos(\phi)$. Fig. 1 depicts the functions corresponding to the rows of an arbitrary matrix $\mathbf{V} \in \mathbb{R}^{4 \times 2}$. We call this a *spannogram*, because at each ϕ , the values of the curves coincide with the entries of a vector in the range of \mathbf{V} . A key observation is that the sorting of the curves at some ϕ is locally invariant for most points in Φ . In fact, due to the continuity of the curves, as we move along the ϕ -axis, the set $\mathcal{I}_k^+(\mathbf{v}(\phi))$ can only change at points where a curve intersects with (i) another curve, or (ii) the zero axis; a change in either the sign of a curve or the relative order of two curves is necessary, although not sufficient, for $\mathcal{I}_k^+(\mathbf{v}(\phi))$ to change.

Appending a zero $(n + 1)^{\text{th}}$ row to \mathbf{V} , the two aforementioned conditions can be merged into one: $\mathcal{I}_k^+(\mathbf{v}(\phi))$ can change only at the points where two of the $n + 1$ curves intersect. Finding the unique intersection point of two curves $[\mathbf{v}(\phi)]_i$ and $[\mathbf{v}(\phi)]_j$ for all pairs $\{i, j\}$ is the key to dis-

Back to the general (P_d) problem, if a linear inequality $\mathbf{R}_{i,:} \mathbf{c} \geq 0$ for some $i \in [k]$ is enforced with equality, the modified problem can be written as a quadratic maximization in the form of (P_d) , with dimension reduced to $d - 1$ and $k - 1$ linear constraints. This observation suggests a recursive algorithm for solving (P_d) : If $\pm \mathbf{u}_1$ is feasible, it is also the optimal solution. Otherwise, for $i = 1, \dots, k$, set the i^{th} inequality constraint active, solve recursively, and collect candidate solutions. Finally, output the candidate that maximizes the objective. The $O(k^d)$ recursive algorithm is formally presented in the supplemental material.

6. Near-Linear Time Nonnegative SPCA

Alg. 1 approximates the nonnegative, k -sparse PC of a PSD matrix \mathbf{A} by solving the nonnegative sparse PCA problem exactly on \mathbf{A}_d , the best rank- d approximation of \mathbf{A} . Albeit polynomial in n , the running time of Alg. 1 can be impractical even for moderate values of n .

Instead of pursuing the exact solution to the low-rank nonnegative sparse PCA problem $\max_{\mathbf{x} \in \mathbb{S}_k^+} \mathbf{x}^T \mathbf{A}_d \mathbf{x}$, we can compute an approximate solution in near-linear time, with performance arbitrarily close to optimal. The suggested procedure is outlined in Algorithm 3, and a detailed discussion is provided in the supplemental material. Alg. 3 relies on randomly sampling points from the range of \mathbf{A}_d and efficiently solving rank-1 instances of the nonnegative sparse PCA problem as described in Section 3.1.

Theorem 2. *For any $n \times n$ PSD matrix \mathbf{A} , sparsity parameter k , and accuracy parameters $d \in [n]$ and $\epsilon \in (0, 1]$, Alg. 3 outputs a nonnegative, k -sparse, unit norm vector $\hat{\mathbf{x}}_d$ such that*

$$\hat{\mathbf{x}}_d^T \mathbf{A} \hat{\mathbf{x}}_d \geq (1 - \epsilon) \cdot \rho_d \cdot \mathbf{x}_*^T \mathbf{A} \mathbf{x}_*,$$

with probability at least $1 - 1/n$, in time $O(\epsilon^{-d} \cdot n \log n)$ plus the time to compute the d leading eigenvectors of \mathbf{A} .

7. Experimental Evaluation

We empirically evaluate the performance of our algorithm on various datasets and compare it to the EM algorithm¹ for sparse and nonnegative PCA of (Sigg & Buhmann, 2008) which is known to outperform previous algorithms.

CBCL Face Dataset. The CBCL face image dataset (Sung, 1996), with 2429 gray scale images of size 19×19 pixels, has been used in the performance evaluation of both the NSPCA (Zass & Shashua, 2007) and EM (Sigg & Buhmann, 2008) algorithms.

Fig. 3 depicts samples from the dataset, as well as six orthogonal, nonnegative, k -sparse components ($k = 40$) successively computed by (i) Alg. 3 ($d = 3, \epsilon = 0.1$) and

¹ Matlab implementation available by the author.

Algorithm 3 Approximate Spannogram NSPCA (ϵ -net)

input \mathbf{A} ($n \times n$ PSD matrix), $k, d \in [n], \epsilon \in (0, 1]$
 1: $[\mathbf{U}, \mathbf{\Lambda}] = \text{svd}(\mathbf{A}, d)$
 2: $\mathbf{V} = \mathbf{U} \mathbf{\Lambda}^{1/2}$ { $\mathbf{A}_d = \mathbf{V} \mathbf{V}^T$ }
 3: $\mathcal{X}_d = \emptyset$
 4: **for** $i = 1 : O(\epsilon^{-d} \cdot \log n)$ **do**
 5: $\mathbf{c} = \text{randn}(d, 1)$
 6: $\mathbf{a} = \mathbf{V} \mathbf{c} / \|\mathbf{c}\|_2$
 7: $\mathbf{x} = \text{rank1solver}(\mathbf{a})$ { Section 3.1 }
 8: $\mathcal{X}_d = \mathcal{X}_d \cup \{\mathbf{x}\}$
 9: **end for**
output $\hat{\mathbf{x}}_d = \arg \max_{\mathbf{x} \in \mathcal{X}_d} \|\mathbf{V}^T \mathbf{x}\|_2^2$

(ii) the EM algorithm. Features active in one component are removed from the dataset prior to computing subsequent PCs to ensure orthogonality. Fig. 3 reveals the ability of nonnegative sparse PCA to extract significant parts.

In Fig. 4, we plot the variance explained by the computed approximate nonnegative, k -sparse PC (normalized by the leading eigenvalue) versus the sparsity parameter k . Alg. 3 for $d = 3$ and $\epsilon = 0.1$, and the EM algorithm exhibit nearly identical performance. For this dataset, we also compute the leading component using the NSPCA algorithm of (Zass & Shashua, 2007). Note that NSPCA does not allow for a precise control of the sparsity of its output; an appropriate sparsity penalty β was determined via binary search for each target sparsity k . We plot the explained variance only for those values of k for which a k -sparse component was successfully extracted. Finally, note that both the EM and NSPCA algorithms are randomly initialized. All depicted values are the best results over multiple random restarts.

Our theory allows us to obtain provable approximation guarantees: based on Theorem 2 and the output of Alg. 3, we compute a data dependent upper bound on the maximum variance, which provably lies in the shaded area. For instance, for $k = 180$, the extracted component explains at least 58% of the variance explained by the true nonnegative, k -sparse PC. The quality of the bound depends on the accuracy parameters d and ϵ , and the eigenvalue decay of the empirical covariance matrix of the data. There exist



Figure 3. We plot (a) six samples from the dataset, and the six leading orthogonal, nonnegative, k -sparse PCs for $k = 40$ extracted by (b) Alg. 3 ($d = 3, \epsilon = 0.1$), and (c) the EM algorithm.

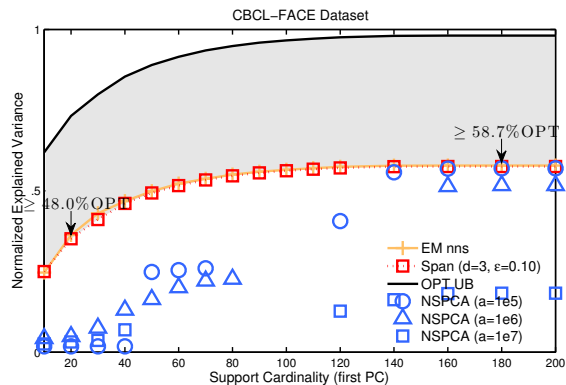


Figure 4. CBCL dataset (Sung, 1996). We plot the normalized variance explained by the approximate nonnegative, k -sparse PC versus the sparsity k . Our theory yields a provable data dependent approximation guarantee: the true unknown optimum provably lies in the shaded area.

datasets on which our algorithm provably achieves 70% or even 90% of the optimal.

Leukemia Dataset. The Leukemia dataset (Armstrong et al., 2001) contains 72 samples, each consisting of expression values for 12582 probe sets. The dataset was used in the evaluation of (Sigg & Buhmann, 2008). In Fig. 5, we plot the normalized variance explained by the computed nonnegative, k -sparse PC versus the sparsity parameter k . For low values of k , Alg. 3 outperforms the EM algorithm in terms of explained variance. For larger values, the two algorithms exhibit similar performance.

The approximation guarantees accompanying our algorithm allow us to upper bound the optimal performance. For k as small as 50, which roughly amounts to 0.4% of the features, the extracted component captures at least 44.6% of the variance corresponding to the true nonnegative k -sparse PC. The obtained upper bound is a significant improvement compared to the trivial bound given by λ_1 .

Low Resolution Spectrometer Dataset. The Low Resolution Spectrometer (LRS) dataset, available in (Bache & Lichman, 2013), originates from the Infra-Red Astronomy Satellite Project. It contains 531 high quality spectra (samples) measured in 93 bands. Fig. 6 depicts the normalized variance explained by the computed nonnegative, k -sparse PC versus the sparsity parameter k . The empirical covariance matrix of this dataset exhibits sharper decay in the spectrum than the previous examples, yielding tighter approximation guarantees according to our theory. For instance, for $k = 20$, the extracted nonnegative component captures at least 86% of the maximum variance. For values closer to $k = 90$, where the computed PC is nonnegative but no longer sparse, this value climbs to nearly 93%.

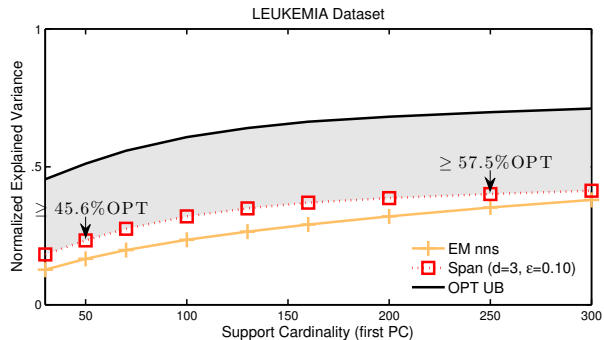


Figure 5. Leukemia dataset (Armstrong et al., 2001). We plot the normalized variance explained by the output of Alg. 3 ($d = 3$, $\epsilon = 0.1$) versus the sparsity k , and compare with the EM algorithm of (Sigg & Buhmann, 2008). By our approximation guarantees, the maximum variance provably lies in the shaded area.

8. Conclusions

We introduced a novel algorithm for nonnegative sparse PCA, expanding the spannogram theory to nonnegative quadratic optimization. We observe that the performance of our algorithm often matches and sometimes outperforms the previous state of the art (Sigg & Buhmann, 2008). Even though the theoretical running time of Alg. 3 scales better than EM, in practice we observed similar speed, both in the order of a few seconds. Our approach has the benefit of provable approximation, giving both theoretical a-priori guarantees and data dependent bounds that can be used to estimate the variance explained by nonnegative sparse PCs, as shown in our experiments.

9. Acknowledgements

The authors would like to acknowledge support from NSF grants CCF-1344364, CCF-1344179, DARPA XDATA and research gifts by Google and Docomo.

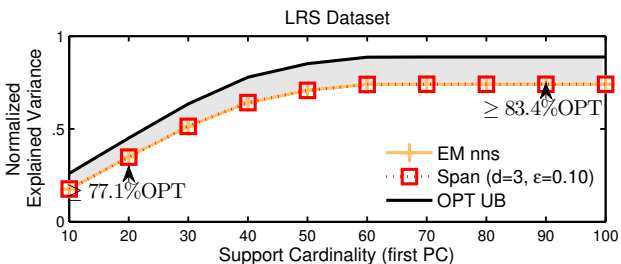


Figure 6. LRS dataset (Bache & Lichman, 2013). We plot the normalized explained variance versus the sparsity k . Alg. 3 ($d = 3$, $\epsilon = 0.1$) and the EM algorithm exhibit similar performance. The optimum value of the objective in (2) provably lies in the shaded area, which in this case is particularly tight.

References

- Allen, Genevera I. and Maletić-Savatić, Mirjana. Sparse non-negative generalized pca with applications to metabolomics. *Bioinformatics*, 2011.
- Armstrong, Scott A, Staunton, Jane E, Silverman, Lewis B, Pieters, Rob, den Boer, Monique L, Minden, Mark D, Sallan, Stephen E, Lander, Eric S, Golub, Todd R, and Korsmeyer, Stanley J. Mll translocations specify a distinct gene expression profile that distinguishes a unique leukemia. *Nature genetics*, 30(1):41–47, 2001.
- Asteris, M., Papailiopoulos, D.S., and Karystinos, G.N. Sparse principal component of a rank-deficient matrix. In *Information Theory Proceedings (ISIT), 2011 IEEE International Symposium on*, pp. 673–677, 2011.
- Bache, K. and Lichman, M. UCI machine learning repository, 2013. URL <http://archive.ics.uci.edu/ml>.
- Badea, Liviu and Tilivea, Doina. Sparse factorizations of gene expression guided by binding data. In *Pacific Symposium on Biocomputing*, 2005.
- d’Aspremont, A., El Ghaoui, L., Jordan, M.I., and Lanckriet, G.R.G. A direct formulation for sparse pca using semidefinite programming. *SIAM review*, 49(3):434–448, 2007a.
- d’Aspremont, Alexandre, Bach, Francis R., and Ghaoui, Laurent El. Full regularization path for sparse principal component analysis. In *Proceedings of the 24th international conference on Machine learning*, ICML ’07, pp. 177–184, New York, NY, USA, 2007b. ACM.
- d’Aspremont, Alexandre, Bach, Francis, and Ghaoui, Laurent El. Optimal solutions for sparse principal component analysis. *J. Mach. Learn. Res.*, 9:1269–1294, Jun 2008.
- Hoyer, Patrik O. Non-negative matrix factorization with sparseness constraints. *Journal of Machine Learning Research*, 5: 1457–1469, 2004.
- Jolliffe, I.T. Rotation of principal components: choice of normalization constraints. *Journal of Applied Statistics*, 22(1):29–35, 1995.
- Jolliffe, I.T., Trendafilov, N.T., and Uddin, M. A modified principal component technique based on the lasso. *Journal of Computational and Graphical Statistics*, 12(3):531–547, 2003.
- Karystinos, G.N. and Liavas, A.P. Efficient computation of the binary vector that maximizes a rank-deficient quadratic form. *Information Theory, IEEE Transactions on*, 56(7):3581–3593, 2010.
- Kim, Hyunsoo and Park, Haesun. Sparse non-negative matrix factorizations via alternating non-negativity-constrained least squares for microarray data analysis. *Bioinformatics*, 23(12): 1495–1502, 2007.
- Lee, Daniel D and Seung, H Sebastian. Learning the parts of objects by non-negative matrix factorization. *Nature*, 401(6755): 788–791, 1999.
- Mackey, Lester. Deflation methods for sparse pca. In *Advances in Neural Information Processing Systems 21*, NIPS ’08, pp. 1–8, Vancouver, Canada, Dec 2008.
- Moghaddam, B., Weiss, Y., and Avidan, S. Spectral bounds for sparse pca: Exact and greedy algorithms. *Advances in neural information processing systems*, 18:915, 2006a.
- Moghaddam, Baback, Weiss, Yair, and Avidan, Shai. Generalized spectral bounds for sparse l₁. In *Proceedings of the 23rd international conference on Machine learning*, ICML ’06, pp. 641–648, 2006b.
- Murty, Katta G. and Kabadi, Santosh N. Some np-complete problems in quadratic and nonlinear programming. *Mathematical Programming*, 39(2):117–129, 1987.
- Papailiopoulos, D. S., Dimakis, A. G., and Korokythakis, S. Sparse pca through low-rank approximations. In *Proceedings of the 30th International Conference on Machine Learning*, ICML ’13, pp. 767–774. ACM, 2013.
- Parrilo, Pablo A. *Structured Semidefinite Programs and Semialgebraic Geometry Methods in Robustness and Optimization*. PhD thesis, California Institute of Technology Pasadena, California, 2000.
- Sigg, Christian D. and Buhmann, Joachim M. Expectation-maximization for sparse and non-negative pca. In *Proceedings of the 25th International Conference on Machine Learning*, ICML ’08, pp. 960–967, New York, NY, USA, 2008. ACM.
- Sung, Kah-Kay. *Learning and example selection for object and pattern recognition*. PhD thesis, PhD thesis, MIT, Artificial Intelligence Laboratory and Center for Biological and Computational Learning, Cambridge, MA, 1996.
- Zass, Ron and Shashua, Amnon. Nonnegative sparse pca. In *Advances in Neural Information Processing Systems 19*, pp. 1561–1568, Cambridge, MA, 2007. MIT Press.
- Zhang, Y., d’Aspremont, A., and Ghaoui, L.E. Sparse pca: Convex relaxations, algorithms and applications. *Handbook on Semidefinite, Conic and Polynomial Optimization*, pp. 915–940, 2012.
- Zou, H., Hastie, T., and Tibshirani, R. Sparse principal component analysis. *Journal of computational and graphical statistics*, 15 (2):265–286, 2006.