

Row-sparse principal component analysis with guarantees

Santanu S. Dey · Marco Molinaro ·
Guanyi Wang

Received: date / Accepted: date

Abstract Row-sparse principal component analysis (rsPCA), also known as principal component analysis (PCA) with global support, is the problem of finding the top- r leading principal components such that all these principal components are linear combination of a subset of k variables. rsPCA is a popular dimension reduction tool in statistics that enhances interpretability compared to regular principal component analysis (PCA). Methods for solving rsPCA in the literature are either greedy heuristics (in the special case of $r = 1$) with guarantees under restrictive statistical-models, or algorithms with stationary point convergence for some regularized reformulation of rsPCA. Crucially, none of the existing computational methods can efficiently guarantee the quality of the solutions obtained by comparing against dual bounds.

In this work we first propose a convex relaxation based on operator norms that provably approximates the feasible region of rsPCA within a $O(\sqrt{\log r})$ factor. To prove this result we use a novel random sparsification procedure that uses the *Pietsch-Grothendieck factorization theorem* and may be of independent interest. We also propose a simpler relaxation that is second-order cone representable and gives a $(1 + \sqrt{r})$ -approximation for the feasible region.

Using these relaxations we then propose a convex integer program that provides a dual bound for the optimal value of rsPCA. Moreover, it also has worst-case guarantees: it is within a multiplicative/additive factor of the orig-

A preliminary version of this paper was published in [44]

Santanu S. Dey
ISyE, Georgia Institute of Technology
E-mail: santanu.dey@isye.gatech.edu

Marco Molinaro
Computer Science Department, Pontifical Catholic University of Rio de Janeiro
E-mail: mmolinaro@inf.puc-rio.br

Guanyi Wang
ISyE, Georgia Institute of Technology
E-mail: gwang93@gatech.edu

inal optimal value, the multiplicative factor being $O(\log r)$ or $O(r)$ depending on the relaxation used.

Finally, our experiments demonstrate both the viability of computing these dual bounds on instance with up to 2000 attributes and also their quality compared to baselines available.

Keywords Row sparse PCA · Randomized rounding · Convex hull

1 Introduction

Principal component analysis (PCA) is a popular tool for dimension reduction and data visualization. Given a *sample matrix* $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_M) \in \mathbb{R}^{d \times M}$ where each column denotes a d -dimensional zero-mean sample, the goal is to find the top- r leading eigenvectors $\mathbf{V} := (\mathbf{v}_1, \dots, \mathbf{v}_r) \in \mathbb{R}^{d \times r}$ (*principal components*), namely the matrix satisfying

$$\arg \max_{\mathbf{V}^\top \mathbf{V} = \mathbf{I}^r} \text{Tr}(\mathbf{V}^\top \mathbf{A} \mathbf{V}), \quad (\text{PCA})$$

where $\text{Tr}(\cdot)$ is the trace, $\mathbf{A} := \frac{1}{M} \mathbf{X} \mathbf{X}^\top$ is the *sample covariance matrix*, and \mathbf{I}^r denotes the $r \times r$ identity matrix.

Principal components usually tend to be dense, that is, the principal components usually involve almost all variables. This leads to a lack of interpretability of the results from PCA, especially in the high-dimensional setting, e.g. clinical analysis, biological gene analysis, computer vision [9, 47, 24]. Moreover, anecdotally the principal component analysis is also known to generate large generalization error, and therefore makes inaccurate prediction. To enhance the interpretability and reduce the generalization error, it is natural to consider alternatives to PCA where a sparsity constraint is incorporated. There are many different choices of sparsity constraint depending on the context and application.

In this paper, we consider the *row-sparse PCA* (rsPCA) problem (see, for example [42]) defined as follows: Given a sample covariance matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$, a *sparsity parameter* k ($\leq d$), the task is to find out the top- r k -sparsity principal components $\mathbf{V} \in \mathbb{R}^{d \times r}$ ($r \leq k$),

$$\arg \max_{\mathbf{V}^\top \mathbf{V} = \mathbf{I}^r, \|\mathbf{V}\|_0 \leq k} \text{Tr}(\mathbf{V}^\top \mathbf{A} \mathbf{V}), \quad (\text{rsPCA})$$

where the *row-sparsity constraint* $\|\mathbf{V}\|_0 \leq k$ denotes that there are at most k non-zero rows in matrix \mathbf{V} , i.e., the principal components share *global support*.

Let

$$\mathcal{F} := \{\mathbf{V} \mid \mathbf{V}^\top \mathbf{V} = \mathbf{I}^r, \|\mathbf{V}\|_0 \leq k\}$$

denote the feasible region of rsPCA and let $\text{opt}^{\mathcal{F}}(\mathbf{A})$ denote the optimal value of rsPCA for sample covariance matrix \mathbf{A} .

1.1 Literature review

There is an extensive literature on (approximately) solving the sparse PCA problem and existing approaches can be broadly classified into the following five categories.

In the first category, instead of dealing with the non-convex sparsity constraint directly, the papers [25, 51, 5, 31, 43, 7, 20, 13] incorporate additional regularizers to the objective function to enhance the sparsity of the solution. Similar to LASSO for sparse linear regression problem, these new formulations can be optimized via alternating-minimization type algorithms. We note here that the optimization problem presented in [25] is NP-hard to solve, and there is no convergence guarantee for the alternating-minimization method given in [51]. The papers [5, 31, 43, 7, 20, 13] propose their own formulations for sparse PCA problem, and show that the alternating-minimization algorithm converges to stationary (critical) points. However, the solutions obtained using the above methods cannot guarantee the row-sparsity constraint $\|\mathbf{V}\|_0 \leq k$. Moreover, none of these methods are able to provide worst-case guarantees.

The second category of methods work with the convex relaxations of sparsity constraint. A majority of this work is for solving rsPCA for the case where $r = 1$. The papers [16, 15, 50, 14, 28, 48] directly incorporate the sparsity constraint (for $r = 1$ case) and then relax the resulting optimization problem into some convex optimization problem – usually a semi-definite programming (SDP) relaxation. However, SDPs are often difficult to scale to large instances in practice. The paper [19] proposes a framework to find dual bounds of sparse PCA problem using convex quadratic integer program for the $r = 1$ case.

A third category of papers present fixed parameter tractable exact algorithms where the fixed parameter is usually the rank of the data matrix \mathbf{A} and r . The paper [36] proposes an exact algorithm to find the global optimal solution of rsPCA with $r = 1$ with running-time of $O(d^{\text{rank}(\mathbf{A})+1} \log d)$. Later the paper [3] gives a combinatorial method for sparse PCA problem with *disjoint* supports. They show that their algorithm outputs a feasible solution within $(1 - \epsilon)$ -multiplicative approximation ratio in time polynomial in data dimension d and reciprocal of ϵ , but exponential in the rank of sample covariance matrix \mathbf{A} and r . Recently [17] provides a general method for solving rsPCA exactly with computational complexity polynomial in d , but exponential in r and $\text{rank}(\mathbf{A})$. The paper [17] mentions that the results obtained are of theoretical nature, and these methods may not be practically implementable.

A fourth category of results is that of specialized iterative heuristic methods for finding good feasible solutions of rsPCA [38, 33, 26, 8, 4, 49, 36] for the $r = 1$ case. These methods do not come with worst-case guarantees. Moreover, to the best of our understanding, there is no natural way to generalize these methods for solving rsPCA when $r > 1$.

The final category of papers are those that present algorithms that perform well under the assumption of a statistical-model. Under the assumption of an underlying statistical-model, the paper [22] presents a family of estimators for rsPCA with so-called ‘oracle property’ via solving semidefinite relaxation

of sparse PCA. The paper [18] analyzes a covariance thresholding algorithm (first proposed by [29]) for the $r = 1$ case. They show that this algorithm correctly recovers the support with high probability for sparse parameter k within order \sqrt{M} , with M being the number of samples. This sample complexity, combining with the lower bounds results in [6,32], suggest that no polynomial time algorithm can do significantly better under their statistical assumptions. There are also a series of papers [42,11,45,10,30] that provide the minimax rate of estimation for sparse PCA. However, all these papers require underlying statistical models, thus do not have worst-case guarantees in the model-free case.

1.2 Our contributions

In this paper, we propose explicit convex relaxations for the sparse PCA with provable quality of their approximations.

Explicit convex relaxations of the feasible region \mathcal{F} . Note that the objective function of rsPCA is that of maximizing a convex function, and so at least one of the extreme points of the feasible region \mathcal{F} is an optimal solution. Hence, it is important to approximate the convex hull of the feasible region well.

Our first explicit convex relaxation for \mathcal{F} has constraints based on the operator norms of the matrix \mathbf{V} and is given by

$$\mathcal{CR1} := \left\{ \mathbf{V} \in \mathbb{R}^{d \times r} \mid \|\mathbf{V}\|_{\text{op}} \leq 1, \|\mathbf{V}\|_{2 \rightarrow 1} \leq \sqrt{k}, \sum_{i=1}^d \|\mathbf{V}_i\|_2 \leq \sqrt{rk} \right\}, \quad (\mathcal{CR1})$$

see Section 1.3 for the required definitions. Importantly, we prove that this relaxation is within a factor of $O(\log r)$ of the convex hull of \mathcal{F} .

Theorem 1 *For every positive integers d, r, k such that $1 \leq r \leq k \leq d$ the convex relaxation $\mathcal{CR1}$ satisfies*

$$\mathcal{F} \subseteq \mathcal{CR1} \subseteq \rho_{\mathcal{CR1}} \cdot \text{conv}(\mathcal{F})$$

for $\rho_{\mathcal{CR1}} = 2 + \max\{6\sqrt{2\pi}, 18\sqrt{\log 50r}\}$. In particular $\rho_{\mathcal{CR1}} = O(\sqrt{\log r})$.

Given the simplicity of the formulation of the relaxation and its provable approximation guarantee, $\mathcal{CR1}$ seems to be an important set not only for sparse PCA but for other problems with row-sparsity constraints. To prove this result we use a novel randomized matrix sparsification procedure that given a matrix \mathbf{V}^* in $\mathcal{CR1}$ produces a row-sparse matrix \mathbf{V} with controlled spectral norm (hence in \mathcal{F}) that is close to the starting point \mathbf{V}^* . The main difficulty is effectively leveraging the information provided by the simple formulation $\mathcal{CR1}$, mainly the control of the $\|\cdot\|_{2 \rightarrow 1}$ norm. For that, we employ in our sparsification a row sampling procedure where the weight of the rows are given

1 by the *Pietsch-Grothendieck factorization theorem* [37]. We believe this idea
 2 may also find uses in other problems with row-sparsity constraints.

3 The main computational drawback of this relaxation is that the norm $\|\cdot\|_{2 \rightarrow 1}$
 4 is known to be NP-hard to compute [39]. Therefore we also present a
 5 simpler convex relaxation of \mathcal{F} that is second-order cone representable, and
 6 hence can be efficiently optimized over using interior-point methods.
 7

8
 9 **Theorem 2** *For every d, r, k positive integers such that $1 \leq r \leq k \leq d$,
 10 there is a relaxation (CR2) that is second-order cone representable and has
 11 the guarantee*

$$12 \quad \text{conv}(\mathcal{F}) \subseteq \text{CR2} \subseteq \rho_{\text{CR2}} \cdot \text{conv}(\mathcal{F}),$$

13
 14 where $\rho_{\text{CR2}} \leq 1 + \sqrt{r}$.

15
 16 This generalizes the main theoretical result in [19] for the case $r = 1$.
 17

18
 19 *Convex integer programming formulation for obtaining dual bounds.* While
 20 the above relaxations allow us to convexify the feasible region, notice that
 21 the objective function of rsPCA is also non-convex (i.e. maximizing a convex
 22 function), and hence do not yet give a “full relaxation” that yield tractable
 23 upper (dual) bounds for the problem. Recall that dual bounds are typically
 24 crucial for effective computational procedures for non-convex problems, which
 25 are usually solved using branch-and-bound type algorithms, as dual bound
 26 allow pruning of the solution space.
 27

28 To handle the non-convex objective function, we consider the natural ap-
 29 proach of upper bounding the the objective function by piecewise linear func-
 30 tions, which can be modeled using binary variables and special ordered sets
 31 (SOS-2) [46]. Used together with a convex relaxation CR1 or CR2 this gives
 32 a convex integer programming relaxation for rsPCA.
 33

34 Interestingly, we show that this full relaxation has a provable approxima-
 35 tion guarantee, providing a dual bound that is within a multiplicative/additive
 36 factor of the optimal value of the original problem.
 37

38
 39 **Theorem 3** *Let $opt^{\mathcal{F}}$ be the optimal value of rsPCA. Then there is a convex
 40 integer program (CIP) using the relaxation CR1 or CR2 whose optimal value
 41 $ub^{\text{CR}i}$ satisfies the following:*

$$42 \quad opt^{\mathcal{F}} \leq ub^{\text{CR}i} \leq \rho_{\text{CR}i}^2 \cdot opt^{\mathcal{F}} + \text{additive-term},$$

43
 44 where the additive term depends on the input matrix \mathbf{A} and the parameters
 45 used in piecewise linear approximation of the objective function, and $\rho_{\text{CR}i}$ is
 46 the approximation guarantee from Theorems 1 and 2 for the relaxation used.
 47
 48
 49
 50
 51
 52
 53
 54
 55
 56
 57
 58
 59
 60
 61
 62
 63
 64
 65

1 *Computational experiments.* In order to evaluate both its practical viability
 2 and quality, we perform computational experiments with the full relaxation
 3 provided by our proposed convex integer program. For that we use both syn-
 4 thetic and real data with up to 2000 features, which yield input matrices of
 5 size 2000×2000 .
 6

7 Obtaining good feasible solutions, required to evaluate our dual bounds,
 8 already poses a challenge given the lack of heuristics for rsPCA when $r > 1$
 9 and the size of the instances. To mitigate this we use an optimized version of
 10 the natural greedy heuristic that looks for which rows of \mathbf{V} should be non-zero.
 11 This optimized version greatly saves on the number of eigenvalue computa-
 12 tions required, which is a bottleneck of the process. Solving our convex integer
 13 program for the larger instances also requires care, and we employ cutting
 14 planes and a submatrix splitting technique to speed up the computations.
 15

16 The numerical results show that our convex integer program provides
 17 within a reasonable time good upper bounds that are typically significantly
 18 better than an SDP relaxation and another baseline.
 19

20 We note that a preliminary version of this paper was published in [44].
 21 The current version has many new results, in particular $\mathcal{CR}1$ and results on
 22 its strength are completely new, and the numerical experiments have also been
 23 completely revamped.
 24

25 1.3 Notation

26 We use regular lower case letters, for example α , to denote scalars. For a
 27 positive integer n , let $[n] := \{1, \dots, n\}$. For a set $S \subseteq \mathbb{R}^n$ and a $\rho > 0$ denote
 28 $\rho \cdot S := \{\rho x \mid x \in S\}$.

29 We use bold lower case letters, for example \mathbf{a} , to be vectors. We denote the
 30 i -th component of a vector \mathbf{a} as \mathbf{a}_i . Given two vectors, $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$, we represent
 31 the inner product of \mathbf{u} and \mathbf{v} by $\langle \mathbf{u}, \mathbf{v} \rangle$. Sometimes it will be convenient to
 32 represent the outer product of vectors using \otimes , i.e., given two vectors $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$,
 33 $\mathbf{a} \otimes \mathbf{b}$ is the matrix where $[\mathbf{a} \otimes \mathbf{b}]_{i,j} = \mathbf{a}_i \mathbf{b}_j$. We denote the unit vector in the
 34 direction of the j th coordinate as \mathbf{e}^j .
 35

36 We use bold upper case letters, for example \mathbf{A} , to denote matrices. We
 37 denote the (i, j) -th component of a matrix \mathbf{A} as $\mathbf{A}_{i,j}$. We use $\text{supp}(\mathbf{A})$ to denote
 38 the support of non-zero rows of matrix \mathbf{A} . We use regular upper case letters,
 39 for example I , to denote the set of indices. Given any matrix $\mathbf{A} \in \mathbb{R}^{n \times m}$ and
 40 $I \subseteq [n], J \subseteq [m]$, we denote the sub-matrix of \mathbf{A} with rows in I and columns
 41 in J as $\mathbf{A}_{I,J}$. For $I \in [m]$, to simplify notation we denote the submatrix of
 42 $\mathbf{A} \in \mathbb{R}^{m \times n}$ corresponding to the rows with index in I as \mathbf{A}_I (instead of $\mathbf{A}_{I,[n]}$).
 43 Similarly for $i \in [m]$, we denote the i^{th} row of \mathbf{A} as \mathbf{A}_i . For $J \in [n]$ again to
 44 simplify the notation, we denote the submatrix of $\mathbf{A} \in \mathbb{R}^{m \times n}$ corresponding
 45 to the columns with index in J as $\mathbf{A}_{*,J}$ (instead of $\mathbf{A}_{[m],J}$), and for $j \in [n]$,
 46 we denote the j^{th} column of \mathbf{A} as $\mathbf{A}_{*,j}$.
 47
 48
 49
 50
 51
 52
 53
 54
 55
 56
 57
 58
 59
 60
 61
 62
 63
 64
 65

For a symmetric square matrix \mathbf{A} , we denote the largest eigen-value of \mathbf{A} as $\lambda_{\max}(\mathbf{A})$. Given $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$, two symmetric matrices, we say that $\mathbf{A} \preceq \mathbf{B}$ if $\mathbf{B} - \mathbf{A}$ is a positive semi-definite matrix. Given $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{m \times n}$, we let $\langle \mathbf{U}, \mathbf{V} \rangle = \sum_{i=1}^m \sum_{j=1}^n U_{ij} V_{ij}$ to be the inner product of matrices. We use $\mathbf{0}^{p,q}$ to denote the matrix of size $p \times q$ with all entries equal to zero. We use \oplus , as a sign of direct sum of matrices, i.e., given matrices $\mathbf{A} \in \mathbb{R}^{p \times q}, \mathbf{B} \in \mathbb{R}^{m \times n}$,

$$\mathbf{A} \oplus \mathbf{B} := \begin{bmatrix} \mathbf{A} & \mathbf{0}^{p,n} \\ \mathbf{0}^{m,q} & \mathbf{B} \end{bmatrix}.$$

The operator norm $\|\mathbf{A}\|_{p \rightarrow q}$ of a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ is defined as

$$\|\mathbf{A}\|_{p \rightarrow q} := \max_{\mathbf{x} \in \mathbb{R}^n, \|\mathbf{x}\|_p=1} \|\mathbf{A}\mathbf{x}\|_q.$$

We sometimes refer $\|\mathbf{A}\|_{2 \rightarrow 2}$ as $\|\mathbf{A}\|_{\text{op}}$. Note that $\|\mathbf{A}\|_{\text{op}}$ is the largest singular value of \mathbf{A} . The Frobenius norm of a matrix \mathbf{A} is denoted as $\|\mathbf{A}\|_F$.

2 Convex relaxations of \mathcal{F}

2.1 Operator-norms relaxation $\mathcal{CR1}$

In the vector case, i.e. $r = 1$ case, a natural convex relaxation for \mathcal{F} is to control the sparsity via the ℓ_2 and ℓ_1 norms, namely to consider the set $\{\mathbf{v} \in \mathbb{R}^d \mid \|\mathbf{v}\|_2 \leq 1, \|\mathbf{v}\|_1 \leq \sqrt{k}\}$ (see [19]). It is easy to see that this is indeed a relaxation in the case $r = 1$: if $\mathbf{v} \in \mathcal{F}$, then by definition $\langle \mathbf{v}, \mathbf{v} \rangle = 1$ and so $\|\mathbf{v}\|_2 = 1$, and since \mathbf{v} is a k -sparse vector we get, using the standard ℓ_1 - vs ℓ_2 -norm comparison in k -dimensional space, $\|\mathbf{v}\|_1 \leq \sqrt{k} \cdot \|\mathbf{v}\|_2 = \sqrt{k}$.

Here we consider the relaxation $\mathcal{CR1}$ that generalizes this idea for any r , which we now recall:

$$\mathcal{CR1} := \left\{ \mathbf{V} \in \mathbb{R}^{d \times r} \mid \|\mathbf{V}\|_{\text{op}} \leq 1, \|\mathbf{V}\|_{2 \rightarrow 1} \leq \sqrt{k}, \sum_{i=1}^d \|\mathbf{V}_i\|_2 \leq \sqrt{rk} \right\}.$$

Thus we now use both the $\ell_{2 \rightarrow 1}$ norm and the sum of the length of the rows of \mathbf{V} to take the role of the ℓ_1 -norm proxy for sparsity (by convexity of norms both constraints are convex). While it is not hard to see that this is a relaxation of \mathcal{F} , we further show that it has a provable approximation guarantee.

Remark 1 One can replace in $\mathcal{CR1}$ the constraint $\|\mathbf{V}\|_{\text{op}} \leq 1$ by the constraint $\begin{bmatrix} \mathbf{I}^r & -\mathbf{V} \\ -\mathbf{V} & \mathbf{I}^r \end{bmatrix} \succeq \mathbf{0}$, which is the convex hull of the *Stiefel manifold* $\{\mathbf{V} \mid \mathbf{V}^\top \mathbf{V} = \mathbf{I}^r\}$ [21].

For the remainder of the section we prove the approximation guarantee that $\mathcal{CR1}$ provides for the convex hull of \mathcal{F} , namely that

$$\mathcal{F} \subseteq \mathcal{CR1} \subseteq \rho_{\mathcal{CR1}} \cdot \text{conv}(\mathcal{F})$$

for $\rho_{\mathcal{CR1}} = 2 + \max\{6\sqrt{2\pi}, 18\sqrt{\log 50r}\}$, thus proving Theorem 1. We prove each of these inclusions separately.

2.1.1 Proof of first inclusion in Theorem 1: $\mathcal{F} \subseteq \mathcal{CR1}$

Consider a matrix \mathbf{V} in \mathcal{F} ; we show that it satisfies the 3 constraints of $\mathcal{CR1}$. First, observe that

$$\begin{aligned} \|\mathbf{V}\|_{\text{op}} &= \max_{\mathbf{x} \in \mathbb{R}^n, \|\mathbf{x}\|_2=1} \|\mathbf{V}\mathbf{x}\|_2 = \max_{\mathbf{x} \in \mathbb{R}^n, \|\mathbf{x}\|_2=1} \sqrt{\langle \mathbf{V}\mathbf{x}, \mathbf{V}\mathbf{x} \rangle} \\ &= \max_{\mathbf{x} \in \mathbb{R}^n, \|\mathbf{x}\|_2=1} \sqrt{\langle \mathbf{x}, \mathbf{V}^\top \mathbf{V}\mathbf{x} \rangle} = 1. \end{aligned}$$

Therefore, we obtain that for $\mathbf{V} \in \mathcal{F}$, we have $\|\mathbf{V}\|_{\text{op}} \leq 1$.

For the second constraint, by definition of $\|\cdot\|_{2 \rightarrow 1}$ it is equivalent to verify that $\|\mathbf{V}\mathbf{x}\|_1 \leq \sqrt{k}$ for all $\mathbf{x} \in \mathbb{R}^r$ such that $\|\mathbf{x}\|_2 \leq 1$. Since \mathbf{V} is k -row-sparse, $\mathbf{V}\mathbf{x}$ is a k -sparse vector and hence by ℓ_1 - vs ℓ_2 -norm comparison in k -dim space we get $\|\mathbf{V}\mathbf{x}\|_1 \leq \sqrt{k} \cdot \|\mathbf{V}\mathbf{x}\|_2 \leq \sqrt{k}$, where the last inequality follows $\|\mathbf{V}\mathbf{x}\|_2 \leq \|\mathbf{V}\|_{\text{op}}$ for all \mathbf{x} satisfying $\|\mathbf{x}\|_2 \leq 1$.

For the third constraint of $\mathcal{CR1}$, since $\|\mathbf{V}\|_{\text{op}} \leq 1$ each column of \mathbf{V} , i.e., $\mathbf{V}_{*,j}$ has a 2-norm of at most 1, and since there are r columns we have:

$$r \geq \|\mathbf{V}\|_F^2 = \sum_{i=1}^d \|\mathbf{V}_i\|_2^2.$$

Since V is k -row-sparse, at most k of the terms in the right-hand side is non-zero. Then again applying the ℓ_1 - vs ℓ_2 -norm comparison in k -dim space we get

$$\sum_{i=1}^d \|\mathbf{V}_i\|_2 \leq \sqrt{k} \cdot \sqrt{\sum_i \|\mathbf{V}_i\|_2^2}.$$

Combining the displayed inequalities gives $\sum_{i=1}^d \|\mathbf{V}_i\|_2 \leq \sqrt{rk}$, and so the third constraint of $\mathcal{CR1}$ is satisfied.

2.1.2 Proof of second inclusion in Theorem 1: $\mathcal{CR1} \subseteq \rho_{\mathcal{CR1}} \cdot \text{conv}(\mathcal{F})$

We assume that $k \geq 40$, otherwise $r \leq k < 40$ and the result follows from Theorem 2. We prove the desired inclusion by comparing the support function of these sets (Proposition C.3.3.1 of [23]), namely we show that for every matrix $\mathbf{C} \in \mathbb{R}^{d \times r}$

$$\max_{\mathbf{V} \in \mathcal{CR1}} \langle \mathbf{C}, \mathbf{V} \rangle \leq \rho_{\mathcal{CR1}} \cdot \max_{\mathbf{V} \in \text{conv}(\mathcal{F})} \langle \mathbf{C}, \mathbf{V} \rangle. \quad (1)$$

It will suffice to prove the following sparsification result for the optimum of the left-hand side.

Lemma 1 *Assume $k \geq 40$. Consider $\mathbf{C} \in \mathbb{R}^{d \times r}$ and let \mathbf{V}^* be a matrix attaining the maximum on the left-hand side of (1), namely $\mathbf{V}^* \in \arg \max_{\mathbf{V} \in \mathcal{CR1}} \langle \mathbf{C}, \mathbf{V} \rangle$. Then there is a matrix \mathbf{V} with the following properties:*

1. (Operator norm) $\|\mathbf{V}\|_{op} \leq 1 + \max\{\sqrt{18\pi}, 6\sqrt{\log 50r}\}$
2. (Sparsity) \mathbf{V} is k -row-sparse, namely $\|\mathbf{V}\|_0 \leq k$
3. (Value) $\langle \mathbf{C}, \mathbf{V} \rangle \geq \frac{1}{2} \langle \mathbf{C}, \mathbf{V}^* \rangle$.

Indeed, if we have such a matrix \mathbf{V} then $\frac{\mathbf{V}}{\|\mathbf{V}\|_{op}}$ belongs to the sparse set \mathcal{F} and has value $\langle \mathbf{C}, \frac{\mathbf{V}}{\|\mathbf{V}\|_{op}} \rangle \geq \frac{1}{2 \cdot (1 + \max\{\sqrt{18\pi}, 6\sqrt{\log 50r}\})} \cdot \langle \mathbf{C}, \mathbf{V}^* \rangle$, showing that (1) holds.

For the remainder of the section we prove Lemma 1. The idea is to randomly sparsify \mathbf{V}^* while controlling for operator norm and value. A standard procedure is to sample the rows of \mathbf{V}^* with probability proportional to their squared length (see [27] for this and other sampling methods). However these more standard methods do not seem effectively leverage the information that $\|\mathbf{V}^*\|_{2 \rightarrow 1} \leq \sqrt{k}$.

Instead, we use a novel sampling more adapted to the $\ell_{2 \rightarrow 1}$ -norm based on a weighting of the rows of \mathbf{V}^* given by the so-called *Pietsch-Grothendieck factorization* [37]. We state it in a convenient form that follows by applying Theorem 2.2 of [40] to the transpose.

Theorem 4 (Pietsch-Grothendieck factorization) *Any matrix $\mathbf{V} \in \mathbb{R}^{d \times r}$ can be factorized as $\mathbf{V} = \mathbf{W}\mathbf{T}$ of size $\mathbf{T} \in \mathbb{R}^{d \times r}$, $\mathbf{W} \in \mathbb{R}^{d \times d}$, where*

- \mathbf{W} is a nonnegative, diagonal matrix with $\sum_i \mathbf{W}_{ii}^2 = 1$
- $\|\mathbf{T}\|_{op} \leq \sqrt{\pi/2} \cdot \|\mathbf{V}\|_{2 \rightarrow 1}$.

So first apply this theorem to obtain a decomposition $\mathbf{V}^* = \mathbf{W}\mathbf{T}$. Notice that this means the i th row of \mathbf{V}^* is just the i th row of \mathbf{T} multiplied by the weight \mathbf{W}_{ii} . Define the “probability”

$$p_i := \frac{k}{6} \left(\mathbf{W}_{ii}^2 + \frac{\|\mathbf{V}_i^*\|_2}{\sum_{i'} \|\mathbf{V}_{i'}^*\|_2} \right),$$

and the truncation $\bar{p}_i = \min\{p_i, 1\}$ to make it a bonafide probability.¹ We then randomly sparsify \mathbf{V}^* by keeping each row i with probability \bar{p}_i and normalizing it: define the random matrix $\tilde{\mathbf{V}} := \tilde{\mathbf{W}}\mathbf{T}$, where $\tilde{\mathbf{W}}$ is the random diagonal matrix with

$$\tilde{\mathbf{W}}_{ii} := \varepsilon_i \frac{\mathbf{W}_{ii}}{\bar{p}_i},$$

and ε_i (the indicator that we keep row i) takes value 1 with probability \bar{p}_i and 0 with probability $1 - \bar{p}_i$ (and the ε_i 's are independent). Since $\mathbb{E}\tilde{\mathbf{W}} = \mathbf{W}$ notice this is procedure is unbiased: $\mathbb{E}\tilde{\mathbf{V}} = \mathbf{V}^*$.

We first show that $\tilde{\mathbf{V}}$ satisfies each of the desired items from Lemma 1 with good probability, and then use a union bound to exhibit a matrix that proves the lemma.

¹ For some intuition: The first term parenthesis in p_i controls the variance of $\tilde{\mathbf{W}}_{ii}$, which is $\text{Var}(\tilde{\mathbf{W}}_{ii}) \leq \frac{\mathbf{W}_{ii}^2}{p_i} \leq \frac{6}{k}$; the second term controls the largest size of a row of $\tilde{\mathbf{V}}$, which is $\|\tilde{\mathbf{V}}_i\|_2 \leq \frac{\|\mathbf{V}_i^*\|_2}{p_i} \leq \frac{6}{k} \sum_{i'} \|\mathbf{V}_{i'}^*\|_2$, which is at most 6 because $\mathbf{V}^* \in \mathcal{CR}1$.

Sparsity. The number of rows $\|\tilde{\mathbf{V}}\|_0$ of $\tilde{\mathbf{V}}$ is precisely $\sum_{i=1}^d \varepsilon_i$, whose expectation is at most

$$\sum_{i=1}^d p_i = \frac{k}{6} \left(\sum_i \mathbf{W}_{ii}^2 + 1 \right) = \frac{k}{3}.$$

Employing the multiplicative Chernoff bound (Lemma 3) we get

$$\Pr \left(\|\tilde{\mathbf{V}}\|_0 > k \right) \leq \left(\frac{2e}{6} \right)^k < \frac{1}{50}, \quad (2)$$

where the last inequality uses that $k \geq 40$.

Operator norm. Let I be the indices i where $p_i \leq 1$ (so $\bar{p}_i = p_i$), and $I^c = [d] \setminus I$ (so $\bar{p}_i = 1$ and hence $\tilde{\mathbf{V}}_i = \mathbf{V}_i^*$). From triangle inequality we can see that $\|\tilde{\mathbf{V}}\|_{\text{op}} \leq \|\tilde{\mathbf{V}}_I\|_{\text{op}} + \|\tilde{\mathbf{V}}_{I^c}\|_{\text{op}}$. Moreover,

$$\|\tilde{\mathbf{V}}_{I^c}\|_{\text{op}} = \|\mathbf{V}_{I^c}^*\|_{\text{op}} \leq \|\mathbf{V}^*\|_{\text{op}} \leq 1,$$

where the first equality is because the rows of $\tilde{\mathbf{V}}_{I^c}$ are exactly equal to the rows of $\mathbf{V}_{I^c}^*$ and the first inequality is because deleting rows cannot increase the operator norm, and the last inequality because $\mathbf{V}^* \in \mathcal{F}$. Combining these observations we get that $\|\tilde{\mathbf{V}}\|_{\text{op}} \leq \|\tilde{\mathbf{V}}_I\|_{\text{op}} + 1$, and so we focus on controlling the operator norm of $\tilde{\mathbf{V}}_I$. We do that by applying a concentration inequality to the largest eivengalue of the PSD matrix $(\tilde{\mathbf{V}}_I)^\top \tilde{\mathbf{V}}_I$; the following is Theorem 1.1 of [41] plus a simple estimate (see for example page 65 of [35]).

Theorem 5 *Let $\mathbf{X}_1, \dots, \mathbf{X}_n \in \mathbb{R}^{r \times r}$ be independent, random, symmetric matrices of dimension r . Assume with probability 1 each \mathbf{X}_i is PSD and has largest eigenvalue $\lambda_{\max}(\mathbf{X}_i) \leq R$. Then*

$$\Pr \left(\lambda_{\max} \left(\sum_i \mathbf{X}_i \right) \geq \alpha \right) < r \cdot 2^{-\alpha/R}$$

for every $\alpha \geq 6\lambda_{\max}(\mathbb{E} \sum_i \mathbf{X}_i)$.

First notice that indeed $(\tilde{\mathbf{V}}_I)^\top \tilde{\mathbf{V}}_I$ can be written as a sum of independent PSD matrices:

$$(\tilde{\mathbf{V}}_I)^\top \tilde{\mathbf{V}}_I = \sum_{i \in I} \tilde{\mathbf{V}}_i \otimes \tilde{\mathbf{V}}_i = \sum_{i \in I} \tilde{\mathbf{W}}_{ii}^2 (\mathbf{T}_i \otimes \mathbf{T}_i) = \sum_{i \in I} \varepsilon_i \frac{\mathbf{W}_{ii}^2}{p_i^2} (\mathbf{T}_i \otimes \mathbf{T}_i). \quad (3)$$

To estimate the max eigenvalue of the expected matrix, $\lambda_{\max}(\mathbb{E} (\tilde{\mathbf{V}}_I)^\top \tilde{\mathbf{V}}_I)$, by definition of p_i we have $\mathbb{E} \varepsilon_i \frac{\mathbf{W}_{ii}^2}{p_i^2} \leq \frac{6}{k}$ and hence

$$\mathbb{E} (\tilde{\mathbf{V}}_I)^\top \tilde{\mathbf{V}}_I \preceq \sum_{i \in I} \frac{6}{k} (\mathbf{T}_i \otimes \mathbf{T}_i) \preceq \frac{6}{k} \sum_i (\mathbf{T}_i \otimes \mathbf{T}_i) = \frac{6}{k} \mathbf{T}^\top \mathbf{T}.$$

By the guarantee of the Pietsch-Grothendieck factorization $\|\mathbf{T}\|_{\text{op}} \leq \sqrt{\pi/2} \|\mathbf{V}^*\|_{2 \rightarrow 1}$ and since $\mathbf{V}^* \in \mathcal{CR1}$ we have $\|\mathbf{V}^*\|_{2 \rightarrow 1} \leq \sqrt{k}$, so applying these bounds to the previous displayed inequality gives

$$\lambda_{\max}\left(\mathbb{E}(\tilde{\mathbf{V}}_I)^\top \tilde{\mathbf{V}}_I\right) \leq \frac{6}{k} \lambda_{\max}(\mathbf{T}^\top \mathbf{T}) = \frac{6}{k} \|\mathbf{T}\|_{\text{op}}^2 \leq 3\pi.$$

To control the R term in Theorem 5 we look at the first equation in (3) and notice that for $i \in I$

$$\begin{aligned} \lambda_{\max}\left(\tilde{\mathbf{V}}_i \otimes \tilde{\mathbf{V}}_i\right) &= \lambda_{\max}\left(\left(\frac{\varepsilon_i}{p_i} \mathbf{V}_i^*\right) \otimes \left(\frac{\varepsilon_i}{p_i} \mathbf{V}_i^*\right)\right) \\ &\leq \lambda_{\max}\left(\frac{1}{p_i^2} (\mathbf{V}_i^* \otimes \mathbf{V}_i^*)\right) = \frac{\|\mathbf{V}_i^*\|_2^2}{p_i^2} \leq \frac{36(\sum_{i'} \|\mathbf{V}_{i'}^*\|_2)^2}{k^2} \leq 36, \end{aligned}$$

where the last inequality uses the fact $\mathbf{V}^* \in \mathcal{CR1}$ and hence $\sum_{i'} \|\mathbf{V}_{i'}^*\|_2 \leq \sqrt{rk} \leq k$.

Then applying Theorem 5 with $\mathbf{X}_i = \tilde{\mathbf{V}}_i \otimes \tilde{\mathbf{V}}_i$, $R = 16$ and $\alpha = \max\{6 \cdot 3\pi, 36 \log 50r\}$ we get

$$\Pr\left(\|\tilde{\mathbf{V}}_I\|_{\text{op}} \geq \sqrt{\alpha}\right) = \Pr\left(\lambda_{\max}((\tilde{\mathbf{V}}_I)^\top \tilde{\mathbf{V}}_I) \geq \alpha\right) < \frac{1}{50}.$$

Recalling we have $\|\tilde{\mathbf{V}}\|_{\text{op}} \leq 1 + \|\tilde{\mathbf{V}}_I\|_{\text{op}}$, this gives that

$$\|\tilde{\mathbf{V}}\|_{\text{op}} > 1 + \max\{\sqrt{18\pi}, 6\sqrt{\log 50r}\} \quad \text{happens with probability at most } \frac{1}{50}. \quad (4)$$

Value. We want to show that with good probability $\langle \mathbf{C}, \tilde{\mathbf{V}} \rangle \geq \frac{1}{2} \langle \mathbf{C}, \mathbf{V}^* \rangle$. We use throughout the following observation: for each row i we have $\langle \mathbf{C}_i, \mathbf{V}_i^* \rangle \geq 0$, since the set $\mathcal{CR1}$ is symmetric with respect to flipping the sign of a row and \mathbf{V}^* maximizes $\langle \mathbf{C}, \mathbf{V}^* \rangle = \sum_i \langle \mathbf{C}_i, \mathbf{V}_i^* \rangle$.

Since $\mathbb{E}\tilde{\mathbf{V}} = \mathbf{V}^*$, we have $\mathbb{E}\langle \mathbf{C}_I, \tilde{\mathbf{V}}_I \rangle = \langle \mathbf{C}_I, \mathbf{V}_I^* \rangle$ and

$$\begin{aligned} \text{Var}(\langle \mathbf{C}_I, \tilde{\mathbf{V}}_I \rangle) &= \sum_{i \in I} \text{Var}(\langle \mathbf{C}_i, \tilde{\mathbf{V}}_i \rangle) = \sum_{i \in I} \text{Var}\left(\frac{\varepsilon_i}{p_i} \langle \mathbf{C}_i, \mathbf{V}_i^* \rangle\right) \leq \sum_{i \in I} \frac{\langle \mathbf{C}_i, \mathbf{V}_i^* \rangle^2}{p_i} \\ &\leq \frac{6 \sum_{i'} \|\mathbf{V}_{i'}^*\|_2}{k} \cdot \sum_{i \in I} \frac{\langle \mathbf{C}_i, \mathbf{V}_i^* \rangle^2}{\|\mathbf{V}_i^*\|_2} \leq 6 \cdot \left(\max_{i \in I} \left\langle \mathbf{C}_i, \frac{\mathbf{V}_i^*}{\|\mathbf{V}_i^*\|_2} \right\rangle\right) \cdot \langle \mathbf{C}_I, \mathbf{V}_I^* \rangle, \end{aligned}$$

where the second inequality uses the definition of p_i and the last inequality uses that $\sum_{i'} \|\mathbf{V}_{i'}^*\|_2 \leq \sqrt{rk} \leq k$ (since $\mathbf{V}^* \in \mathcal{CR1}$). Moreover, since $\frac{\mathbf{V}_i^*}{\|\mathbf{V}_i^*\|_2}$ also belongs to $\mathcal{CR1}$, the optimality of \mathbf{V}^* guarantees that $\langle \mathbf{C}_i, \frac{\mathbf{V}_i^*}{\|\mathbf{V}_i^*\|_2} \rangle \leq \langle \mathbf{C}, \mathbf{V}^* \rangle$, and so we have the variance upper bound

$$\text{Var}(\langle \mathbf{C}_I, \tilde{\mathbf{V}}_I \rangle) \leq 6 \cdot \langle \mathbf{C}, \mathbf{V}^* \rangle^2.$$

Using the fact that $\langle \mathbf{C}_{I^c}, \tilde{\mathbf{V}}_{I^c} \rangle = \langle \mathbf{C}_{I^c}, \mathbf{V}_{I^c}^* \rangle$ and the one-sided Chebychev inequality (Lemma 4) we get

$$\Pr \left(\langle \mathbf{C}, \tilde{\mathbf{V}} \rangle \leq \frac{1}{2} \langle \mathbf{C}, \mathbf{V}^* \rangle \right) = \Pr \left(\langle \mathbf{C}_I, \tilde{\mathbf{V}}_I \rangle \leq \langle \mathbf{C}_I, \mathbf{V}_I^* \rangle - \frac{1}{2} \langle \mathbf{C}, \mathbf{V}^* \rangle \right) \quad (5)$$

$$\leq \frac{6}{6 + \frac{1}{4}} = 1 - \frac{1}{25}. \quad (6)$$

Concluding the proof of Lemma 1. Taking a union bound over inequalities (2), (4), and (6), we see that with positive probability $\tilde{\mathbf{V}}$ satisfies all items from Lemma 1. This shows the existence of the desired matrix \mathbf{V} and concludes the proof.

2.2 Second order representable relaxation $\mathcal{CR2}$

Since an optimization problem involving the semi-definite constraint $\mathbf{V}^\top \mathbf{V} \preceq \mathbf{I}^r$ (equivalent to $\|\mathbf{V}\|_{op} \leq 1$) and the $\ell_{2 \rightarrow 1}$ -norm constraint $\|\mathbf{V}\|_{2 \rightarrow 1} \leq \sqrt{k}$ may be challenging to solve in practice we consider the following further relaxation involving second-order cone constraints:

$$\mathcal{CR2} := \left\{ \mathbf{V} \in \mathbb{R}^{d \times r} \left| \begin{array}{l} \|\mathbf{V}_{*,j}\|_2^2 \leq 1 \quad \forall j \in [r] \\ \|\mathbf{V}_{*,j_1} \pm \mathbf{V}_{*,j_2}\|_2^2 \leq 2 \quad \forall j_1 \neq j_2 \in [r] \\ \|\mathbf{V}_{*,j}\|_1 \leq \sqrt{k} \quad \forall j \in [r] \\ \sum_{i=1}^d \|\mathbf{V}_i\|_2 \leq \sqrt{rk} \end{array} \right. \right\}. \quad (\mathcal{CR2})$$

This set is a relaxation of $\mathcal{CR1}$ obtained by considering the constraint $\max_{\mathbf{x}: \|\mathbf{x}\|_2 \leq 1} \|\mathbf{V}\mathbf{x}\|_2 = \|\mathbf{V}\|_{op} \leq 1$ only for the vectors $\mathbf{x} = \mathbf{e}^j$ and $\mathbf{x} = \frac{1}{\sqrt{2}}(\mathbf{e}^{j_1} \pm \mathbf{e}^{j_2})$, and considering the constraint $\max_{\mathbf{x}: \|\mathbf{x}\|_2 \leq 1} \|\mathbf{V}\mathbf{x}\|_1 = \|\mathbf{V}\|_{2 \rightarrow 1} \leq \sqrt{k}$ only for the vectors $\mathbf{x} = \mathbf{e}^j$. In particular this shows that $\mathcal{CR2}$ is a relaxation of $\mathcal{CR1}$ and hence a relaxation of \mathcal{F} . Moreover, we now show that it still gives a guaranteed approximation to the convex hull of \mathcal{F} , namely that

$$\text{conv}(\mathcal{F}) \subseteq \mathcal{CR2} \subseteq (1 + \sqrt{r}) \text{conv}(\mathcal{F}),$$

thus proving Theorem 2.

Proof (Proof of Theorem 2) We have just argued that $\mathcal{CR2}$ is a relaxation of \mathcal{F} so it suffices to show the second inclusion $\mathcal{CR2} \subseteq (1 + \sqrt{r}) \text{conv}(\mathcal{F})$. So consider any $\mathbf{V} \in \mathcal{CR2}$, and we will show $\mathbf{V} \in (1 + \sqrt{r}) \text{conv}(\mathcal{F})$.

Since the sets \mathcal{F} and $\mathcal{CR2}$ are symmetric to row permutations, assume without loss of generality that the rows of \mathbf{V} are sorted in non-decreasing length, namely $\|\mathbf{V}_1\|_2 \geq \|\mathbf{V}_2\|_2 \geq \dots$. Decompose \mathbf{V} based on its top- k largest rows, second top- k largest rows, and so on, i.e., let $m = \lceil d/k \rceil$, $\mathbf{V} = \mathbf{V}^1 + \dots + \mathbf{V}^m$ with $\mathbf{V}^p \in \mathbb{R}^{d \times r}$ and

$$\text{supp}(\mathbf{V}^1) = \{1, \dots, k\} =: I^1, \quad \dots, \quad \text{supp}(\mathbf{V}^m) = \{d - (m-1)k, \dots, d\} =: I^m.$$

For each $p = 1, \dots, m$ we have $\|\mathbf{V}^p / \|\mathbf{V}^p\|_{\text{op}}\|_0 \leq k$ and $\|\mathbf{V}^p / \|\mathbf{V}^p\|_{\text{op}}\|_{\text{op}} = 1$, thus $\mathbf{V}^p / \|\mathbf{V}^p\|_{\text{op}} \in \mathcal{F}$. Observe that:

$$\begin{aligned} \mathbf{V} &= \mathbf{V}^1 + \dots + \mathbf{V}^m = \|\mathbf{V}^1\|_{\text{op}} \frac{\mathbf{V}^1}{\|\mathbf{V}^1\|_{\text{op}}} + \dots + \|\mathbf{V}^m\|_{\text{op}} \frac{\mathbf{V}^m}{\|\mathbf{V}^m\|_{\text{op}}} \quad (7) \\ \Rightarrow \frac{\mathbf{V}}{\sum_{p=1}^m \|\mathbf{V}^p\|_{\text{op}}} &= \left(\frac{\|\mathbf{V}^1\|_{\text{op}}}{\sum_{p=1}^m \|\mathbf{V}^p\|_{\text{op}}} \right) \frac{\mathbf{V}^1}{\|\mathbf{V}^1\|_{\text{op}}} + \dots + \left(\frac{\|\mathbf{V}^m\|_{\text{op}}}{\sum_{p=1}^m \|\mathbf{V}^p\|_{\text{op}}} \right) \frac{\mathbf{V}^m}{\|\mathbf{V}^m\|_{\text{op}}} \\ &\in \text{conv}(\mathcal{F}). \end{aligned}$$

Notice that $\|\mathbf{V}^1\|_{\text{op}} \leq 1$, since $\|\mathbf{V}\|_{\text{op}} \leq 1$ and zeroing out rows cannot increase the operator norm, and also by standard relationship between $\|\cdot\|_2$ and $\|\cdot\|_F$ we have:

$$\|\mathbf{V}^p\|_{\text{op}} \leq \sqrt{\sum_{i \in I^p} \|\mathbf{V}_i\|_2^2}.$$

Furthermore, we can bound the norm of each of these rows of \mathbf{V}^p by the average of the rows of \mathbf{V}^{p-1} , since the rows of \mathbf{V} are sorted in non-decreasing length. Employing these bounds we get

$$\begin{aligned} \sum_{p=1}^m \|\mathbf{V}^p\|_{\text{op}} &= \|\mathbf{V}^1\|_{\text{op}} + \sum_{p=2}^m \|\mathbf{V}^p\|_{\text{op}} \\ &\leq 1 + \sum_{p=2}^m \sqrt{\left(\frac{\sum_{i \in I^{p-1}} \|\mathbf{V}_i\|_2}{k} \right)^2 \cdot k} \\ &= 1 + \frac{1}{\sqrt{k}} \cdot \sum_{p=2}^m \sum_{i \in I^{p-1}} \|\mathbf{V}_i\|_2 \\ &\leq 1 + \frac{1}{\sqrt{k}} \sum_{i=1}^d \|\mathbf{V}_i\|_2 \leq 1 + \sqrt{r} \quad (8) \end{aligned}$$

where the final inequality holds since the constraint $\sum_{i=1}^d \|\mathbf{V}_i\|_2 \leq \sqrt{rk}$ is in the description of $\mathcal{CR}2$.

Combining inequalities (7) and (8) we have

$$\mathbf{V} \in \left(\sum_{p=1}^m \|\mathbf{V}^p\|_{\text{op}} \right) \cdot \text{conv}(\mathcal{F}) \subseteq (1 + \sqrt{r}) \cdot \text{conv}(\mathcal{F}).$$

concluding the proof of the theorem.

3 Convex IP formulation for obtaining dual bounds for rsPCA

Based on the results in the Section 2, we can set-up the following optimization problem:

$$\text{opt}^{\mathcal{CR}i} := \max_{\mathbf{V} \in \mathcal{CR}i} \text{Tr}(\mathbf{V}^\top \mathbf{A} \mathbf{V}). \quad (\text{CRi-Relax})$$

The following is a straightforward Corollary of Theorem 1 and Theorem 2:

Corollary 1 $\text{opt}^{\mathcal{F}} \leq \text{opt}^{\mathcal{CR}i} \leq \rho_{\mathcal{CR}i}^2 \text{opt}^{\mathcal{F}}$ for $i \in \{1, 2\}$.

The challenge of solving CRi-Relax is that the objective function is non-concave. Indeed, for the case $r = 1$, Corollary 1 provides constant multiplicative approximation ratios to rsPCA; thus, the inapproximability results for rsPCA with $r = 1$ from [12, 34] imply that solving CRi-Relax to optimality is NP-hard. Therefore we construct a further concave relaxation of the objective function.

3.1 Piecewise linear upper approximation of objective function

Let $\mathbf{A} = \sum_{j=1}^d \lambda_j \mathbf{a}_j \mathbf{a}_j^\top$ be the eigenvalue decomposition of sample covariance matrix \mathbf{A} with $\lambda_1 \geq \dots \geq \lambda_d \geq 0$. The objective function then can be represented as a summation

$$\text{Tr}(\mathbf{V}^\top \mathbf{A} \mathbf{V}) = \sum_{j=1}^d \lambda_j \sum_{i=1}^r (\mathbf{a}_j^\top \mathbf{v}_i)^2$$

where \mathbf{v}_i denotes the i th column of \mathbf{V} such that $\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_r)$. Set auxiliary variables $g_{ji} = \mathbf{a}_j^\top \mathbf{v}_i$ for $(j, i) \in [r] \times [d]$. Let $\mathbf{a}_j \in \mathbb{R}^d$ satisfy

$$|[\mathbf{a}_j]_{j_1}| \geq \dots \geq |[\mathbf{a}_j]_{j_k}| \geq \dots \geq |[\mathbf{a}_j]_{j_d}|,$$

and let

$$\theta_j = \sqrt{[\mathbf{a}_j]_{j_1}^2 + \dots + [\mathbf{a}_j]_{j_k}^2}$$

be the square root of sum of top- k largest absolute entries of \mathbf{a}_j . Since \mathbf{v}_i is supposed to be k -sparse, it is easy to observe that g_{ji} is within the interval $[-\theta_j, \theta_j]$.

Piecewise linear approximation: To relax the non-convex objective, we can upper approximate each quadratic term g_{ji}^2 by a piecewise linear function based on a new auxiliary variable ξ_{ji} via *special ordered sets type 2* (SOS-II) constraints (PLA) as follows,

$$\text{PLA}([d] \times [r]) := \left\{ (g, \xi, \eta) \left| \begin{array}{ll} g_{ji} = \mathbf{a}_j^\top \mathbf{v}_i & (j, i) \in [d] \times [r] \\ g_{ji} = \sum_{\ell=-N}^N \gamma_{ji}^\ell \eta_{ji}^\ell & (j, i) \in [d] \times [r] \\ \xi_{ji} = \sum_{\ell=-N}^N (\gamma_{ji}^\ell)^2 \eta_{ji}^\ell & (j, i) \in [d] \times [r] \\ (\eta_{ji}^\ell)_{\ell=-N}^N \in \text{SOS-II} & (j, i) \in [d] \times [r] \end{array} \right. \right\}$$

where for each $(j, i) \in [d] \times [r]$, $(\eta_{ji}^\ell)_{\ell=-N}^N$ is the set of SOS-II variables, and $(\gamma_{ji}^\ell)_{\ell=-N}^N$ is the corresponding set of splitting points that satisfy:

$$\underbrace{\gamma_{ji}^{-N}}_{=-\theta_j} \leq \dots \leq \underbrace{\gamma_{ji}^0}_{=0} \leq \dots \leq \underbrace{\gamma_{ji}^N}_{=\theta_j}$$

which split the region $[-\theta_j, \theta_j]$ into $2N$ equal intervals. See Figure 1 for an example. By using PLA, we arrive at the following *convex integer programming*

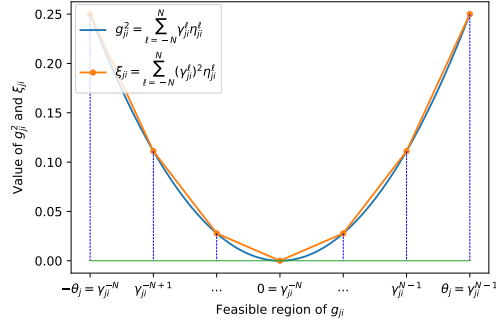


Fig. 1 The quadratic function g_{ji}^2 is upper approximated by a piecewise linear function ξ_{ji} by SOS-II constraints for all $(j, i) \in [d] \times [r]$.

problem,

$$\begin{aligned} \text{ub}^{\mathcal{CR}i} &:= \max \sum_{j=1}^d \lambda_j \sum_{i=1}^r \xi_{ji} \\ \text{s.t. } &\mathbf{V} \in \mathcal{CR}i \\ &(g, \xi, \eta) \in \text{PLA}([d] \times [r]) \end{aligned} \quad (\text{CIP})$$

where $\mathcal{CR}i$ is the convex set defined in Section 2.1 or Section 2.2 for $i \in \{1, 2\}$ respectively, and PLA is the set of constraints for piecewise-linear upper approximation of objective. Note that we say this is a convex integer program since SOS-II is modeled using binary variables.

3.2 Guarantees on the upper bounds from the convex integer program

Here we present the worst-case guarantee on the upper bound from solving convex integer program in the form of an affine function of $\text{opt}^{\mathcal{F}}$. This is a more precise restatement of Theorem 3 from the introduction.

Theorem 3 (restated) *Let $\text{opt}^{\mathcal{F}}$ be the optimal value of rsPCA. Let $\text{ub}^{\mathcal{CR}i}$ be the upper bound obtained from solving the convex integer program using $\mathcal{CR}i$*

convex relaxation of \mathcal{F} for $i \in \{1, 2\}$. Then:

$$\text{opt}^{\mathcal{F}} \leq \text{ub}^{\mathcal{CR}_i} \leq \rho_{\mathcal{CR}_i}^2 \cdot \text{opt}^{\mathcal{F}} + \underbrace{\sum_{j=1}^d \frac{r\lambda_j\theta_j^2}{4N^2}}_{\text{additive term}}, \quad \text{for } i \in \{1, 2\}.$$

Proof Based on the construction for CIP, the objective function $\text{Tr}(\mathbf{V}^\top \mathbf{A} \mathbf{V})$ satisfies

$$\sum_{j=1}^d \lambda_j \sum_{i=1}^r (\mathbf{a}_j^\top \mathbf{v}_i)^2 = \sum_{j=1}^d \lambda_j \sum_{i=1}^r g_{ji}^2.$$

By Corollary 1, we have

$$\max_{\mathbf{V} \in \mathcal{CR}_i} (\mathbf{V}^\top \mathbf{A} \mathbf{V}) = \max_{\mathbf{V} \in \mathcal{CR}_i} \sum_{j=1}^d \lambda_j \sum_{i=1}^r g_{ji}^2 \leq \rho_{\mathcal{CR}_i}^2 \cdot \text{opt}^{\mathcal{F}},$$

for $i \in \{1, 2\}$. Note that $g_{ji} \in [-\theta_j, \theta_j]$ and we have split the interval $[-\theta_j, \theta_j]$ evenly via splitting points $(\gamma_{ji}^\ell)_{\ell=-N}^N$ such that $\gamma_{ji}^\ell = \frac{\ell}{N} \cdot \theta_j$. For a given $j \in [d]$ and $i \in [r]$, by the definition of SOS-II sets, let $g_{ij} = \gamma_{ji}^{\ell^*} \eta_{j,i}^{\ell^*} + \gamma_{ji}^{\ell^*+1} \eta_{j,i}^{\ell^*+1}$, $\xi_{ji} = (\gamma_{ji}^{\ell^*})^2 \eta_{j,i}^{\ell^*} + (\gamma_{ji}^{\ell^*+1})^2 \eta_{j,i}^{\ell^*+1}$ and $\eta_{j,i}^{\ell^*} + \eta_{j,i}^{\ell^*+1} = 1$ for some $\ell^* \in \{-N, \dots, N-1\}$. Thus we have:

$$\begin{aligned} \xi_{ji} - g_{ji}^2 &= \left((\gamma_{ji}^{\ell^*})^2 \eta_{j,i}^{\ell^*} + (\gamma_{ji}^{\ell^*+1})^2 \eta_{j,i}^{\ell^*+1} \right) - \left(\gamma_{ji}^{\ell^*} \eta_{j,i}^{\ell^*} + \gamma_{ji}^{\ell^*+1} \eta_{j,i}^{\ell^*+1} \right)^2 \\ &= (\gamma_{ji}^{\ell^*})^2 \eta_{j,i}^{\ell^*} + (\gamma_{ji}^{\ell^*+1})^2 \eta_{j,i}^{\ell^*+1} - (\gamma_{ji}^{\ell^*})^2 (\eta_{j,i}^{\ell^*})^2 - (\gamma_{ji}^{\ell^*+1})^2 (\eta_{j,i}^{\ell^*+1})^2 \\ &\quad - 2\gamma_{ji}^{\ell^*} \eta_{j,i}^{\ell^*} \gamma_{ji}^{\ell^*+1} \eta_{j,i}^{\ell^*+1} \\ &= \left(\gamma_{ji}^{\ell^*+1} - \gamma_{ji}^{\ell^*} \right)^2 \eta_{j,i}^{\ell^*} \eta_{j,i}^{\ell^*+1} = \frac{\theta_j^2}{N^2} \eta_{j,i}^{\ell^*} \eta_{j,i}^{\ell^*+1} \leq \frac{\theta_j^2}{4N^2}. \end{aligned}$$

Therefore, the objective function in CIP satisfies

$$\sum_{j=1}^d \lambda_j \sum_{i=1}^r \xi_{ji} \leq \sum_{j=1}^d \lambda_j \sum_{i=1}^r g_{ji}^2 + \sum_{j=1}^d \frac{r\lambda_j\theta_j^2}{4N^2} \leq \rho_{\mathcal{CR}_i}^2 \cdot \text{opt}^{\mathcal{F}} + \sum_{j=1}^d \frac{r\lambda_j\theta_j^2}{4N^2},$$

which completes the proof.

4 Greedy heuristic for rsPCA

In order to evaluate the dual bounds produced by the convex integer program from the previous section we also need good feasible solutions for rsPCA. As mentioned in the introduction, we are not aware of any heuristics for the general case $r > 1$, so in this section we describe the optimized version of the natural greedy heuristic that we will use.

We can view rsPCA as the problem

$$\max_{S \subseteq [d], |S|=k} f(S) \text{ where,} \\ f(S) := \left(\max_{\mathbf{V} \in \mathbb{R}^{d \times r} \mid \mathbf{V}^\top \mathbf{V} = \mathbf{I}^r, \text{supp}(\mathbf{V})=S} \text{Tr}(\mathbf{V}^\top \mathbf{A} \mathbf{V}) \right), \quad (9)$$

and hence solving rsPCA reduces to selecting the correct support set S . Thus, a natural algorithm is the *1-neighborhood* local search that starts with a support set S and removes/adds one index to improve the value $f(S)$. The main issue with this strategy is that it requires an expensive eigendecomposition computation for each candidate pair i/j of indices to be removed/added in order to evaluate the function f . Here we propose a much more efficient strategy that solves a proxy version of this local search move that requires only 1 eigendecomposition per round.

For that we rewrite the problem as follows. Given a sample covariance matrix \mathbf{A} , let $\mathbf{A}^{1/2}$ be its positive semi-definite square root such that $\mathbf{A} = \mathbf{A}^{1/2} \mathbf{A}^{1/2}$. Observe that $\|\mathbf{A}^{1/2} - \mathbf{V} \mathbf{V}^\top \mathbf{A}^{1/2}\|_F^2 = \text{Tr}(\mathbf{A}) - \text{Tr}(\mathbf{V}^\top \mathbf{A} \mathbf{V})$, and therefore we may equivalently solve the following problem:

$$\min_{\mathbf{V} \in \mathbb{R}^{d \times r}} \|\mathbf{A}^{1/2} - \mathbf{V} \mathbf{V}^\top \mathbf{A}^{1/2}\|_F^2 \text{ s.t. } \mathbf{V}^\top \mathbf{V} = \mathbf{I}^r, \|\mathbf{V}\|_0 \leq k. \quad (\text{SPCA-alt})$$

Therefore, SPCA-alt can be reformulated into a *two-stage (inner & outer) optimization problem*:

$$\min_{S \subseteq [d], |S| \leq k} \min_{\mathbf{V}_S} \bar{f}(S, \mathbf{V}_S) \text{ s.t. } \mathbf{V}_S^\top \mathbf{V}_S = \mathbf{I}^r$$

where

$$\bar{f}(S, \mathbf{M}) := \|(\mathbf{A}^{1/2})_S - \mathbf{M} \mathbf{M}^\top (\mathbf{A}^{1/2})_S\|_F^2 + \|(\mathbf{A}^{1/2})_{S^C}\|_F^2 \quad (10)$$

and $S^C := [d] \setminus S$.

In order to find a solution with small $\bar{f}(S, \mathbf{V}_S)$ again we use a greedy swap heuristic that removes/adds one index to S . However, we avoid eigenvalue computations by keeping $\mathbf{M} = \mathbf{V}_S$ fixed and finding an improved set S' (i.e., with $\bar{f}(S', \mathbf{M}) \leq \bar{f}(S, \mathbf{M})$), and only then updating the term \mathbf{M} ; only the second only needs 1 eigendecomposition of \mathbf{A}_{S_t, S_t} . We describe this in more detail, letting S_t and $\mathbf{V}_{S_t}^t$ be the iterates at round t .

Leaving Candidate: In the t -th iteration, given the iterates S_{t-1} and $\mathbf{V}_{S_{t-1}}^{t-1}$ from the previous iteration, for each index $j \in S_{t-1}$, let Δ_j^{out} be

$$\Delta_j^{\text{out}} := \|\mathbf{A}_j^{1/2}\|_2^2 - \left\| \mathbf{A}_{S_{t-1}}^{1/2} - \mathbf{V}_{S_{t-1}} \mathbf{V}_{S_{t-1}}^\top \mathbf{A}_{S_{t-1}}^{1/2} \right\|_F^2.$$

Then let $j^{\text{out}} := \arg \min_{j \in S_{t-1}} \Delta_j^{\text{out}}$ be the candidate to leave the set S_{t-1} .

Entering Candidate: Similarly, for each $j \in S_{t-1}^C$ define Δ_j^{in} as

$$\Delta_j^{\text{in}} := \|\mathbf{A}_j^{1/2}\|_2^2 - \left\| (\mathbf{A}^{1/2})_{S_{t-1}^j} - \mathbf{V}_{S_{t-1}} \mathbf{V}_{S_{t-1}}^\top (\mathbf{A}^{1/2})_{S_{t-1}^j} \right\|_F^2,$$

where $S_{t-1}^j := S_{t-1} - \{j^{\text{out}}\} + \{j\}$. Then let $j^{\text{in}} := \arg \max_{j \in S_{t-1}^C} \Delta_j^{\text{in}}$.

1 *Update Rule:* If $\Delta_{j^{\text{out}}}^{\text{out}} < \Delta_{j^{\text{in}}}^{\text{in}}$ we perform the exchange with the candidates
 2 above, namely set $S_t = S_{t-1} - \{j^{\text{out}}\} + \{j^{\text{in}}\}$. In addition, we set $\mathbf{V}_{S_t}^t$ to be
 3 the minimizer of $\min\{f(S_t, \mathbf{M}) : \mathbf{M}^\top \mathbf{M} = \mathbf{I}^r\}$; for that we compute the
 4 eigendecomposition $\mathbf{A}_{S_t, S_t} = \mathbf{U}_{S_t} \mathbf{\Lambda}_{S_t} \mathbf{U}_{S_t}^\top$ of \mathbf{A}_{S_t, S_t} and set $\mathbf{V}_{S_t}^t = (\mathbf{U}_{S_t})_{*, [r]}$
 5 to be the eigenvectors corresponding to top r eigenvalues.
 6

7 If $\Delta_{j^{\text{out}}}^{\text{out}} \geq \Delta_{j^{\text{in}}}^{\text{in}}$ the algorithm stops and return the matrix \mathbf{V} where in
 8 rows S_{t-1} equals $\mathbf{V}_{S_{t-1}}^{t-1}$ (i.e., $\mathbf{V}_{S_{t-1}} = \mathbf{V}_{S_{t-1}}^{t-1}$) and in rows S_{t-1}^C equals zero.
 9 The complete pseudocode is presented in Appendix B.1.
 10

11 We observe that even though our procedure works only with a proxy of
 12 the original function f of the natural greedy heuristic, by construction it still
 13 finds support sets S that monotonically decrease this objective function (see
 14 Appendix B.2 for a proof).

15 **Lemma 2** *Algorithm 1 is a monotonically decreasing algorithm with respect*
 16 *to the objective function f , namely $f(S_t) < f(S_{t-1})$ for every iteration t .*
 17

18 5 Computational experiments

19 In this section we conduct computational experiments on fairly large instances
 20 to illustrate the efficiency of our proposed methods and to asses their qualities
 21 both in terms of finding good primal solutions and proving good dual bounds.
 22 We also compare our dual bound against that obtained from an SDP relaxation
 23 and from another baseline.
 24

25 5.1 Methods for comparison

26 5.1.1 Methods for dual bounds

27 In order to generate dual bounds we implemented a version of our convex
 28 integer programming formulation (CIP), adding several enhancements like re-
 29 duction of the number of SOS-II constraints and cutting planes in order to
 30 improve its efficiency (see [19] for related ideas for the case of $r = 1$). This im-
 31 plemented version is called CIP-impl, and is described in detail in Appendix C.
 32 For all experiments we use $N = 40$ as the level of discretization for the objec-
 33 tive function in CIP-impl. (For large instances we additionally use a dimension
 34 reduction technique, which we discuss later.)
 35

36 We compare our proposed dual bound with the following two baselines:

- 37 – **Baseline 1:** Sum of diagonal entries of sub-matrix:

$$38 \text{ Baseline1} := \mathbf{A}_{j_1, j_1} + \dots + \mathbf{A}_{j_k, j_k}, \text{ where } \mathbf{A}_{j_1, j_1} \geq \mathbf{A}_{j_2, j_2} \geq \dots \mathbf{A}_{j_d, j_d}.$$

39 Note the sum of $\mathbf{A}_{j_1, j_1}, \dots, \mathbf{A}_{j_k, j_k}$ is equal to sum of eigenvalues of sub-
 40 matrix indexed by $\{j_1, \dots, j_k\}$ in \mathbf{A} , then Baseline-1 can be viewed as an
 41 upper bound for the optimal value of rsPCA. Moreover, Baseline-1 is tight
 42 when we have $r = k$.
 43
 44
 45
 46
 47
 48
 49
 50
 51
 52
 53
 54
 55
 56
 57
 58
 59
 60
 61
 62
 63
 64
 65

– **Baseline 2:** The semi-definite programming relaxation:

$$\text{SDP} := \max_{\mathbf{P}} \text{Tr}(\mathbf{A}\mathbf{P}), \text{ s.t. } \mathbf{I}_d \succeq \mathbf{P} \succeq \mathbf{0}, \text{Tr}(\mathbf{P}) = r, \mathbf{1}^\top |\mathbf{P}| \mathbf{1} \leq rk.$$

Note that this is an SDP relaxation of rsPCA obtained by lifting the variables \mathbf{V} into the product space $\mathbf{P} = \mathbf{V}\mathbf{V}^\top$.

5.1.2 Parameter for primal algorithm (lower bounds)

To obtain good feasible solutions we implemented the modified greedy neighborhood search (Algorithm 1) proposed in Section 4. For each instance we run this algorithm 400 times, where each time we pick the initial support set S_0 as a uniformly random subset of $[d]$ of size k . We allow a maximum of d iterations. The objective function value corresponding to the best solution from the 400 runs is declared as the lower bound.

5.2 Instances for numerical experiments

We conducted numerical experiments on two types of instances.

5.2.1 Artificial instances

These instances were generated artificially using ideas similar to that of the *spiked covariance matrix* [18] that have been used often to test algorithms in the $r = 1$ case. An instance **Artificial- k^A** is generated as follows.

We first choose a sparsity parameter $k^A \leq \frac{d}{2}$ (which will be in the range [30]) and the orthonormal vectors \mathbf{u}_1 and \mathbf{u}_2 of dimension k^A given by

$$\mathbf{u}_1^\top = \left(\frac{1}{\sqrt{k^A}}, \dots, \frac{1}{\sqrt{k^A}} \right), \quad \mathbf{u}_2^\top = \left(\frac{1}{\sqrt{k^A}}, -\frac{1}{\sqrt{k^A}}, \dots, \frac{1}{\sqrt{k^A}}, -\frac{1}{\sqrt{k^A}} \right).$$

The *block spiked covariance matrix* $\Sigma \in \mathbb{R}^{d \times d}$ is then computed as

$$\Sigma := \Sigma_1 \oplus \Sigma_2 \oplus \mathbf{I}^{d-2k^A},$$

where $\Sigma_1 := 55\mathbf{u}_1\mathbf{u}_1^\top + 52\mathbf{u}_2\mathbf{u}_2^\top \in \mathbb{R}^{k^A \times k^A}$, $\Sigma_2 := 50\mathbf{I}_{k^A} \in \mathbb{R}^{k^A \times k^A}$. Finally, we sample M i.i.d. random vectors $\mathbf{x}_1, \dots, \mathbf{x}_M \sim N(\mathbf{0}_d, \Sigma)$ from the normal distribution with covariance matrix Σ and create the instance \mathbf{A} as the sample covariance matrix of these vectors:

$$\mathbf{A} := \frac{1}{M} (\mathbf{x}_1\mathbf{x}_1^\top + \dots + \mathbf{x}_M\mathbf{x}_M^\top).$$

In our experiments we use $d = 500$ (thus generating 500×500 matrices) and $M = 3000$ samples. Our experiments will focus on the cases $r = 2$ and $r = 3$ and we note that in these instances the optimal support set with cardinality k^A is different for both choices of r .

5.2.2 Real instances

The second type of instances are four real instances using the colon cancer dataset (CovColon) from [2], the lymphoma dataset (Lymph) from [1], and Reddit instances Reddit1500 and Reddit2000 from [19]. Table 1 presents the size of each instance.

name	CovColon	Lymph	Reddit1500	Reddit2000
size	500×500	500×500	1500×1500	2000×2000

Table 1 Real instances

5.3 Software & hardware

Software & Hardware: All numerical experiments are implemented on MacBookPro13 with 2GHz Intel Core i5 CPU and 8GB 1867MHz LPDDR3 Memory. The (CIP-impl) model was solved using Gurobi 7.0.2. The Baseline-2 model was solved using Mosek.

5.4 Performance measure

We measure the performances of CIP-impl and the baselines based on the primal-dual gap, defined as

$$\text{Gap} := \frac{\text{ub} - \text{lb}}{\text{lb}}.$$

Here $\text{ub} \in \{\text{ub}^{\text{impl}}, (\text{ub}^{\text{sub-mat}}$ in Section 5.6.1), Baseline-1, Baseline-2} denotes the dual bound obtained from CIP-impl or baselines. The term lb denotes the primal bound from the primal heuristic.

5.5 Numerical results for smaller instances

First we perform experiments on smaller instances of size 100×100 . These instances were constructed by picking the submatrix corresponding to the top 100 largest diagonal entries from each instance listed in Section 5.2. We append a “prime” in the name of the instances to denote these smaller instances, e.g., Artificial- k^A and CovColon’.

Time limits. We set the time limit for CIP-impl to 60 seconds and imposed no time limit on SDP. (We note that on these smaller instances SDP terminated within 600 seconds.) We also did not impose a time limit on the primal heuristic, and just note that it took less than 120 seconds on all smaller instances.

The gaps obtained by the dual bounds using CIP-impl, Baseline1, and SDP on these instances are presented in Tables 2 and 3.

name	param (r, k) :	(2, 10)	(2, 20)	(2, 30)	(3, 10)	(3, 20)	(3, 30)
Artificial-10' 100×100	CIP-impl	0.031	0.0004	0.0003	0.04	0.0005	0.0004
	Baseline1	3.523	4.309	4.403	2.108	2.625	2.689
	SDP	0.032	0.0004	0.0003	0.043	0.0005	0.0003
Artificial-20' 100×100	CIP-impl	0.027	0.011	0.007	0.026	0.011	0.006
	Baseline1	3.58	7.838	8.251	2.094	4.942	5.216
	SDP	0.02	0.014	0.008	0.027	0.014	0.006
Artificial-30' 100×100	CIP-impl	0.071	0.022	0.015	0.074	0.023	0.012
	Baseline1	3.503	7.614	11.68	2.066	4.814	7.508
	SDP	0.03	0.021	0.02	0.051	0.026	0.014

Table 2 Gap values for smaller artificial instances with size 100×100

name	param (r, k) :	(2, 10)	(2, 20)	(2, 30)	(3, 10)	(3, 20)	(3, 30)
CovColon' 100×100	CIP-impl	0.12	0.119	0.094	0.127	0.124	0.104
	Baseline1	0.063	0.117	0.132	0.052	0.086	0.098
	SDP	0.674	0.688	0.663	1.244	1.186	1.052
Lymp' 100×100	CIP-impl	0.329	0.272	0.269	0.225	0.296	0.32
	Baseline1	0.095	0.277	0.392	0.049	0.178	0.297
	SDP	0.529	0.449	0.362	0.943	0.695	0.567
Reddit1500' 100×100	CIP-impl	0.155	0.139	0.126	0.129	0.109	0.025
	Baseline1	0.695	0.396	0.99	1.197	0.811	1.294
	SDP	0.265	0.294	0.242	0.175	0.146	0.033
Reddit2000' 100×100	CIP-impl	0.029	0.014	0.011	0.092	0.054	0.011
	Baseline1	0.876	1.426	1.794	0.638	1.075	1.333
	SDP	0.106	0.062	0.036	0.160	0.084	0.034

Table 3 Gap values for smaller real instances with size 100×100

Observations:

- In Table 2 we see that for the relatively easy artificial instances both CIP-impl and SDP find quite tight upper bounds.
- In Table 3 we see that for real instances SDP is substantially dominated by both CIP-impl and Baseline1.

Overall, on the 42 instances, the dual bounds from CIP-impl are best for 28 instances, the dual bounds from Baseline-1 are best for 9 instances, and the dual bounds from SDP are best for 9 instances. Since the computation of Baseline-1 scales trivially in comparison to solving the SDP, and since SDP seems to produce dual bounds of poorer quality for the more difficult real instances — in the next section we discarded SDP from the comparison.

5.6 Larger instances

5.6.1 Sub-matrix technique for larger instances

In order to scale the convex integer program CIP-impl to handle the larger matrices, that are now up to 2000×2000 , we employ the following “sub-matrix technique” to reduce the dimension.

Given a *sub-matrix ratio parameter* $m \geq 1$ satisfying $\lceil mk \rceil \leq d$, let $S := \{j_1, \dots, j_{\lceil mk \rceil}\}$, where $\mathbf{A}_{j_1, j_1} \geq \dots \geq \mathbf{A}_{j_{\lceil mk \rceil}, j_{\lceil mk \rceil}}$, be the index set of the top- $\lceil mk \rceil$ largest diagonal entries of \mathbf{A} . Consider the blocked representation of the sample covariance matrix \mathbf{A} :

$$\mathbf{A} = \begin{pmatrix} \mathbf{A}_{S,S} & \mathbf{A}_{S,S^C} \\ \mathbf{A}_{S,S^C}^\top & \mathbf{A}_{S^C,S^C} \end{pmatrix},$$

where $S^C := [d] \setminus S$. Then the optimal value $\text{opt}^{\mathcal{F}}$ satisfies

$$\begin{aligned} \text{opt}^{\mathcal{F}} &= \max_{\mathbf{V} \in \mathcal{F}} \text{Tr}(\mathbf{V}^\top \mathbf{A} \mathbf{V}) \\ &= \max_{\mathbf{V} \in \mathcal{F}} \text{Tr}((\mathbf{V}_S)^\top \mathbf{A}_{S,S} \mathbf{V}_S) + 2 \text{Tr}((\mathbf{V}_S)^\top \mathbf{A}_{S,S^C} \mathbf{V}_{S^C}) \\ &\quad + \text{Tr}((\mathbf{V}_{S^C})^\top \mathbf{A}_{S^C,S^C} \mathbf{V}_{S^C}). \end{aligned} \quad (\text{submatrix-tech})$$

The first and third term have straight forward upper bounds. Now we need to consider the problem of finding an upper bound on $\text{Tr}((\mathbf{V}_S)^\top \mathbf{A}_{S,S^C} \mathbf{V}_{S^C})$.

Let S^* be the global optimal row-support set of rsPCA. Then

$$\begin{aligned} &\text{Tr}((\mathbf{V}_S)^\top \mathbf{A}_{S,S^C} \mathbf{V}_{S^C}) \\ &= \text{Tr} \left(((\mathbf{V}_{S \cap S^*})^\top (\mathbf{V}_{S \setminus S^*})^\top) \begin{pmatrix} \mathbf{A}_{S \cap S^*, S^C \cap S^*} & \mathbf{A}_{S \cap S^*, S^C \setminus S^*} \\ \mathbf{A}_{S \setminus S^*, S^C \cap S^*} & \mathbf{A}_{S \setminus S^*, S^C \setminus S^*} \end{pmatrix} \begin{pmatrix} \mathbf{V}_{S^C \cap S^*} \\ \mathbf{V}_{S^C \setminus S^*} \end{pmatrix} \right) \\ &= \text{Tr}((\mathbf{V}_{S \cap S^*})^\top \mathbf{A}_{S \cap S^*, S^C \cap S^*} \mathbf{V}_{S^C \cap S^*}). \end{aligned}$$

Since $\mathbf{V}^\top \mathbf{V} = \mathbf{I}^r$, then we have $\mathbf{V}_{S \cap S^*}^\top \mathbf{V}_{S \cap S^*} + \mathbf{V}_{S^C \cap S^*}^\top \mathbf{V}_{S^C \cap S^*} = \mathbf{I}^r$. Thus it is sufficient to consider the following optimization problem:

$$2 \max_{\mathbf{V}^1, \mathbf{V}^2} \text{Tr}((\mathbf{V}^1)^\top \mathbf{A}_{S \cap S^*, S^C \cap S^*} \mathbf{V}^2) \quad \text{s.t.} \quad (\mathbf{V}^1)^\top \mathbf{V}^1 + (\mathbf{V}^2)^\top \mathbf{V}^2 = \mathbf{I}^r,$$

We show in Proposition 2, proved in the appendix, that the above term is upper bounded by $\sqrt{r} \cdot \|\mathbf{A}_{(S \cap S^*), (S^C \cap S^*)}\|_F$.

Therefore, letting $\tilde{k} := |S \cap S^*|$ be the cardinality of the intersection, we can upper bound the right-hand side of (submatrix-tech) as

$$\text{opt}^{\mathcal{F}} \leq \text{ub}^{\text{CIP}}(\mathbf{A}_{S,S}; \tilde{k}) + \sqrt{r} \cdot \|\mathbf{A}_{S \cap S^*, S^C \cap S^*}\|_F + \text{Baseline-1}(\mathbf{A}_{S^C, S^C}; k - \tilde{k}),$$

where the first term $\text{ub}^{\text{CIP}}(\mathbf{A}_{S,S}; \tilde{k})$ is the optimal value obtained from CIP-impl with covariance matrix $\mathbf{A}_{S,S}$ and sparsity parameter \tilde{k} (if $\tilde{k} < r$, then reset $\tilde{k} = r$), and the the third term is the value of Baseline-1 obtained from \mathbf{A}_{S^C, S^C} with sparsity parameter $k - \tilde{k}$.

Since S^* is unknown, then the second term can be further upper bounded by

$$\|\mathbf{A}_{S \cap S^*, S^* \setminus S}\|_F \leq \sqrt{\|\mathbf{A}_{\{j_1\}, S^C}^{k-\tilde{k}}\|_2^2 + \dots + \|\mathbf{A}_{\{j_{\tilde{k}}\}, S^C}^{k-\tilde{k}}\|_2^2} =: \text{ub}(S; \tilde{k}; S^C; k - \tilde{k}),$$

where

$$\|\mathbf{A}_{\{j\}, S^C}^l\|_2^2 := \mathbf{A}_{j, i_1}^2 + \dots + \mathbf{A}_{j, i_l}^2 \text{ with } |\mathbf{A}_{j, i_1}| \geq \dots \geq |\mathbf{A}_{j, i_l}| \geq \dots \text{ for all } i \in S^C,$$

and $j_1, \dots, j_{\tilde{k}}$ are indices satisfying: $\|\mathbf{A}_{j_1, S^C}^{k-\tilde{k}}\|_2^2 \geq \dots \geq \|\mathbf{A}_{j_{\tilde{k}}, S^C}^{k-\tilde{k}}\|_2^2 \geq \dots$.

Since \tilde{k} is also not known, we arrive at our final upper bound $\text{ub}^{\text{sub-mat}}$ by considering all of its possibilities:

$$\begin{aligned} \text{opt}^{\mathcal{F}} &\leq \max_{\tilde{k}=0}^k \left\{ \text{ub}^{\text{CIP}}(\mathbf{A}_{S,S}; \tilde{k}) + \sqrt{r} \cdot \text{ub}(S; \tilde{k}; S^C; k - \tilde{k}) + \text{Baseline-1}(\mathbf{A}_{S^C, S^C}; k - \tilde{k}) \right\} \\ &=: \text{ub}^{\text{sub-mat}}. \end{aligned}$$

5.6.2 Times for larger instances

We set a more stringent time limit of 20 seconds for each CIP-impl used within the sub-matrix technique, since a number of these computations are required to compute $\text{ub}^{\text{sub-mat}}$. Again we did not set a time limit for the primal heuristic, an just note its running times as a function of the matrix size on Table 4.

size	500 × 500	1500 × 1500	2000 × 2000
running time	≤ 20 min	≤ 100 min	≤ 120 min

Table 4 Running time for primal heuristic

5.6.3 Results on larger instances

We compare the gap obtained by the upper bound $ub^{\text{sub-mat}}$ (CIP-impl plus sub-matrix technique) and compare it against that obtained by Baseline1 on the artificial and real instances with original sizes. These are reported on Tables 5 and 6.

On the spiked covariance matrix artificial instances we see that our dual bound $ub^{\text{sub-mat}}$ is typically orders of magnitude better than Baseline1, and is at most 0.35 for all instances. These results also illustrate that the sub-matrix ratio parameter can have a big impact on the bound obtained by the sub-matrix technique.

On the real instances, we see from Table 6 that on instances CovColon and Lymph our dual bound $ub^{\text{sub-mat}}$ performs slightly better than Baseline1 (except instance Lymph with parameters (3,10)), and the gaps are overall less than 0.39. However, on instances Reddit1500 and Reddit2000 our dual bound $ub^{\text{sub-mat}}$ vastly outperforms Baseline1 on all settings of parameters. We remark that these are the largest instances in the experiments, which attest the scalability of our proposed bound.

name	param (r, k) :	(2, 10)	(2, 20)	(2, 30)	(3, 10)	(3, 20)	(3, 30)
Artificial-10 500×500	$m = 1.5$	0.527	0.151	0.25	0.366	0.1	0.169
	$m = 2$	0.079	0.15	0.249	0.064	0.1	0.169
	$m = 2.5$	0.079	0.15	0.248	0.064	0.099	0.168
	$m = 5$	0.071	0.145	0.241	0.056	0.099	0.293
	$m = 10$	0.026	0.002	0.002	0.03	0.003	0.003
	Baseline1	3.522	4.309	4.403	2.101	2.625	2.688
Artificial-20 500×500	$m = 1.5$	2.397	0.566	0.268	1.629	0.384	0.186
	$m = 2$	0.455	0.179	0.266	0.317	0.127	0.185
	$m = 2.5$	0.606	0.178	0.265	0.463	0.126	0.184
	$m = 5$	0.097	0.176	0.261	0.078	0.124	0.346
	$m = 10$	0.073	0.014	0.009	0.139	0.013	0.008
	Baseline1	3.58	7.838	8.251	2.097	4.942	5.216
Artificial-30 500×500	$m = 1.5$	3.515	0.595	0.65	2.071	0.406	0.425
	$m = 2$	3.509	0.721	0.314	2.068	0.512	0.211
	$m = 2.5$	2.304	0.709	0.312	1.586	0.511	0.209
	$m = 5$	0.474	0.225	0.305	0.365	0.158	0.468
	$m = 10$	0.231	0.026	0.017	0.349	0.154	0.014
	Baseline1	3.519	7.626	11.68	2.074	4.82	7.508

Table 5 Gap values for artificial instances.

name	param (r, k) :	(2, 10)	(2, 20)	(2, 30)	(3, 10)	(3, 20)	(3, 30)
CovColon 500×500	$m = 1.5$	0.054	0.112	0.128	0.05	0.08	0.092
	$m = 2$	0.051	0.107	0.126	0.062	0.076	0.09
	$m = 2.5$	0.05	0.104	0.124	0.066	0.089	0.088
	$m = 5$	0.094	0.113	0.143	0.11	0.122	2.349
	$m = 10$	1.787	1.709	1.645	3.321	3.124	3.015
	Baseline1	0.063	0.118	0.133	0.049	0.086	0.097
Lymph 500×500	$m = 1.5$	0.09	0.27	0.41	0.064	0.174	0.315
	$m = 2$	0.078	0.267	0.406	0.103	0.171	0.312
	$m = 2.5$	0.104	0.264	0.403	0.155	0.194	0.309
	$m = 5$	0.236	0.268	0.388	0.2	0.296	2.698
	$m = 10$	2.105	1.738	1.548	4.489	3.894	3.447
	Baseline1	0.095	0.277	0.413	0.049	0.18	0.319
Reddit1500 1500×1500	$m = 1.5$	0.687	0.95	0.8	0.39	0.625	0.677
	$m = 2$	0.683	0.94	0.749	0.387	0.617	0.632
	$m = 2.5$	0.672	0.937	0.727	0.377	0.614	0.611
	$m = 5$	0.426	0.47	1.068	0.346	0.393	1.307
	$m = 10$	0.384	0.927	1.075	0.316	1.222	1.343
	Baseline1	0.695	0.962	1.199	0.396	0.635	0.848
Reddit2000 2000×2000	$m = 1.5$	0.845	1.408	0.76	0.556	1.026	0.667
	$m = 2$	0.837	1.4	0.664	0.549	1.019	0.585
	$m = 2.5$	0.827	1.396	0.601	0.541	1.016	0.538
	$m = 5$	0.456	0.436	1.52	0.395	0.381	1.311
	$m = 10$	0.298	0.866	2.234	0.266	1.289	1.41
	Baseline1	0.876	1.426	1.775	0.582	1.041	1.326

Table 6 Gap values for real instances.

6 Conclusion

In this paper, we proposed a scheme for producing good primal feasible solutions and dual bounds for rsPCA problem. The primal feasible solution is obtained from a monotonically improving heuristic for rsPCA problem. We showed that the solution produced by this algorithm are of very high quality by comparing the objective value of the solutions generated to upper bounds. These upper bounds are obtained using second order cone IP relaxation designed in this paper. We also presented theoretical guarantees (affine guarantee) on the quality of the upper bounds produced by the second order cone IP. The running-time for both the primal algorithm and the dual bounding heuristic are very reasonable (less than 2 hours for the 500×500 instances and less than 3.5 hours for the 2000×2000 instance). These problems are quite challenging and on some instances, we still need more techniques to close the gap. However, to the best of our knowledge, there is no comparable theoretical or computational results for solving model-free rsPCA.

References

1. Alizadeh, A.A., Eisen, M.B., Davis, R.E., Ma, C., Lossos, I.S., Rosenwald, A., Boldrick, J.C., Sabet, H., Tran, T., Yu, X., et al.: Distinct types of diffuse large b-cell lymphoma identified by gene expression profiling. *Nature* **403**(6769), 503 (2000)
2. Alon, U., Barkai, N., Notterman, D.A., Gish, K., Ybarra, S., Mack, D., Levine, A.J.: Broad patterns of gene expression revealed by clustering analysis of tumor and normal colon tissues probed by oligonucleotide arrays. *Proceedings of the National Academy of Sciences* **96**(12), 6745–6750 (1999)
3. Asteris, M., Papailiopoulos, D., Kyrillidis, A., Dimakis, A.G.: Sparse PCA via bipartite matchings. In: *Advances in Neural Information Processing Systems*, pp. 766–774 (2015)
4. Asteris, M., Papailiopoulos, D.S., Karystinos, G.N.: Sparse principal component of a rank-deficient matrix. In: *2011 IEEE International Symposium on Information Theory Proceedings*, pp. 673–677. IEEE (2011)
5. Attouch, H., Bolte, J., Redont, P., Soubeyran, A.: Proximal alternating minimization and projection methods for nonconvex problems: An approach based on the kurdyka-łojasiewicz inequality. *Mathematics of Operations Research* **35**(2), 438–457 (2010)
6. Berthet, Q., Rigollet, P.: Computational lower bounds for sparse pca. *arXiv preprint arXiv:1304.0828* (2013)
7. Bolte, J., Sabach, S., Teboulle, M.: Proximal alternating linearized minimization for nonconvex and nonsmooth problems. *Mathematical Programming* **146**(1-2), 459–494 (2014)
8. Boutsidis, C., Drineas, P., Magdon-Ismael, M.: Sparse features for PCA-like linear regression. In: *Advances in Neural Information Processing Systems*, pp. 2285–2293 (2011)
9. Burgel, P.R., Paillasseur, J., Caillaud, D., Tillie-Leblond, I., Chanez, P., Escamilla, R., Perez, T., Carré, P., Roche, N., et al.: Clinical COPD phenotypes: a novel approach using principal component and cluster analyses. *European Respiratory Journal* **36**(3), 531–539 (2010)
10. Cai, T., Ma, Z., Wu, Y.: Optimal estimation and rank detection for sparse spiked covariance matrices. *Probability theory and related fields* **161**(3-4), 781–815 (2015)
11. Cai, T.T., Ma, Z., Wu, Y., et al.: Sparse PCA: Optimal rates and adaptive estimation. *The Annals of Statistics* **41**(6), 3074–3110 (2013)
12. Chan, S.O., Papailiopoulos, D., Rubinstein, A.: On the approximability of sparse PCA. In: *Conference on Learning Theory*, pp. 623–646 (2016)
13. Chen, S., Ma, S., Xue, L., Zou, H.: An alternating manifold proximal gradient method for sparse PCA and sparse CCA. *arXiv preprint arXiv:1903.11576* (2019)
14. d’Aspremont, A., Bach, F., El Ghaoui, L.: Approximation bounds for sparse principal component analysis. *Mathematical Programming* **148**(1-2), 89–110 (2014)
15. d’Aspremont, A., Bach, F., Ghaoui, L.E.: Optimal solutions for sparse principal component analysis. *Journal of Machine Learning Research* **9**(Jul), 1269–1294 (2008)
16. d’Aspremont, A., Ghaoui, L.E., Jordan, M.I., Lanckriet, G.R.: A direct formulation for sparse PCA using semidefinite programming. In: *Advances in neural information processing systems*, pp. 41–48 (2005)
17. Del Pia, A.: Sparse PCA on fixed-rank matrices. http://www.optimization-online.org/DB_HTML/2019/07/7307.html (2019)
18. Deshpande, Y., Montanari, A.: Sparse PCA via covariance thresholding. *The Journal of Machine Learning Research* **17**(1), 4913–4953 (2016)
19. Dey, S.S., Mazumder, R., Wang, G.: A convex integer programming approach for optimal sparse pca. *arXiv preprint arXiv:1810.09062* (2018)
20. Erichson, N.B., Zheng, P., Manohar, K., Brunton, S.L., Kutz, J.N., Aravkin, A.Y.: Sparse principal component analysis via variable projection. *arXiv preprint arXiv:1804.00341* (2018)
21. Gallivan, K.A., Absil, P.: Note on the convex hull of the stiefel manifold. *Technical note* (2010)
22. Gu, Q., Wang, Z., Liu, H.: Sparse PCA with oracle property. In: *Advances in neural information processing systems*, pp. 1529–1537 (2014)
23. Hiriart-Urruty, J.B., Lemaréchal, C.: *Fundamentals of convex analysis*. Springer Science & Business Media (2012)

24. Jolliffe, I.T., Cadima, J.: Principal component analysis: a review and recent developments. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences* **374**(2065), 20150202 (2016)
25. Jolliffe, I.T., Trendafilov, N.T., Uddin, M.: A modified principal component technique based on the LASSO. *Journal of computational and Graphical Statistics* **12**(3), 531–547 (2003)
26. Journée, M., Nesterov, Y., Richtárik, P., Sepulchre, R.: Generalized power method for sparse principal component analysis. *Journal of Machine Learning Research* **11**(Feb), 517–553 (2010)
27. Kannan, R., Vempala, S.: Randomized algorithms in numerical linear algebra. *Acta Numerica* **26**, 95 (2017)
28. Kim, J., Tawarmalani, M., Richard, J.P.P.: Convexification of permutation-invariant sets and applications. *arXiv preprint arXiv:1910.02573* (2019)
29. Krauthgamer, R., Nadler, B., Vilenchik, D., et al.: Do semidefinite relaxations solve sparse PCA up to the information limit? *The Annals of Statistics* **43**(3), 1300–1322 (2015)
30. Lei, J., Vu, V.Q., et al.: Sparsistency and agnostic inference in sparse PCA. *The Annals of Statistics* **43**(1), 299–322 (2015)
31. Ma, S.: Alternating direction method of multipliers for sparse principal component analysis. *Journal of the Operations Research Society of China* **1**(2), 253–274 (2013)
32. Ma, T., Wigderson, A.: Sum-of-squares lower bounds for sparse pca. In: *Advances in Neural Information Processing Systems*, pp. 1612–1620 (2015)
33. Mackey, L.W.: Deflation methods for sparse PCA. In: *Advances in neural information processing systems*, pp. 1017–1024 (2009)
34. Magdon-Ismail, M.: NP-hardness and inapproximability of sparse PCA. *Information Processing Letters* **126**, 35–38 (2017)
35. Mitzenmacher, M., Upfal, E.: *Probability and computing: Randomization and probabilistic techniques in algorithms and data analysis*. Cambridge university press (2017)
36. Papaliopoulos, D., Dimakis, A., Korokythakis, S.: Sparse PCA through low-rank approximations. In: *International Conference on Machine Learning*, pp. 747–755 (2013)
37. Pietsch, A.: *Operator ideals*, vol. 16. Deutscher Verlag der Wissenschaften (1978)
38. Sigg, C.D., Buhmann, J.M.: Expectation-maximization for sparse and non-negative PCA. In: *Proceedings of the 25th international conference on Machine learning*, pp. 960–967. ACM (2008)
39. Steinberg, D.: *Computation of matrix norms with applications to robust optimization*. Research thesis, Technion-Israel University of Technology **2** (2005)
40. Tropp, J.A.: Column subset selection, matrix factorization, and eigenvalue optimization. In: *Proceedings of the twentieth annual ACM-SIAM symposium on Discrete algorithms*, pp. 978–986. SIAM (2009)
41. Tropp, J.A.: User-friendly tail bounds for sums of random matrices. *Foundations of computational mathematics* **12**(4), 389–434 (2012)
42. Vu, V., Lei, J.: Minimax rates of estimation for sparse PCA in high dimensions. In: *Artificial intelligence and statistics*, pp. 1278–1286 (2012)
43. Vu, V.Q., Cho, J., Lei, J., Rohe, K.: Fantope projection and selection: A near-optimal convex relaxation of sparse PCA. In: *Advances in neural information processing systems*, pp. 2670–2678 (2013)
44. Wang, G., Dey, S.: Upper bounds for model-free row-sparse principal component analysis. In: *Proceedings of the International Conference on Machine Learning* (2020)
45. Wang, Z., Lu, H., Liu, H.: Tighten after relax: Minimax-optimal sparse PCA in polynomial time. In: *Advances in neural information processing systems*, pp. 3383–3391 (2014)
46. Wolsey, L.A., Nemhauser, G.L.: *Integer and combinatorial optimization*, vol. 55. John Wiley & Sons (1999)
47. Yeung, K.Y., Ruzzo, W.L.: Principal component analysis for clustering gene expression data. *Bioinformatics* **17**(9), 763–774 (2001)
48. Yongchun Li, W.X.: Exact and approximation algorithms for sparse PCA. http://www.optimization-online.org/DB_HTML/2020/05/7802.html (2020)
49. Yuan, X.T., Zhang, T.: Truncated power method for sparse eigenvalue problems. *Journal of Machine Learning Research* **14**(Apr), 899–925 (2013)

- 1
2
3
4
5
6
7
8
9
50. Zhang, Y., d'Aspremont, A., El Ghaoui, L.: Sparse PCA: Convex relaxations, algorithms and applications. In: Handbook on Semidefinite, Conic and Polynomial Optimization, pp. 915–940. Springer (2012)
51. Zou, H., Hastie, T., Tibshirani, R.: Sparse principal component analysis. Journal of computational and graphical statistics **15**(2), 265–286 (2006)

10 Appendix

11 A Additional concentration inequalities

12 We need the standard multiplicative Chernoff bound (see Theorem 4.4 [35]).

13 **Lemma 3 (Chernoff Bound)** *Let X_1, \dots, X_n be independent random variables taking values in $[0, 1]$. Then for any $\delta > 0$ we have*

$$14 \Pr\left(\sum_i X_i > (1 + \delta)\mu\right) < \left(\frac{e}{1 + \delta}\right)^{(1 + \delta)\mu},$$

15 where $\mu = \mathbb{E} \sum_i X_i$.

16 We also need the one-sided Chebychev inequality, see for example Exercise 3.18 of [35].

17 **Lemma 4 (One-sided Chebychev)** *For any random variable X with finite first and second moments*

$$18 \Pr\left(X \leq \mathbb{E}X - t\right) \leq \frac{\text{Var}(X)}{\text{Var}(X) + t^2}.$$

19 B Greedy heuristic for rsPCA

20 B.1 Complete pseudocode

21 B.2 Proof of Lemma 2

22 By optimality of $\mathbf{V}_{S_t}^t$ we can see that $f(S_t) = f(S_t, \mathbf{V}_{S_t}^t)$ for all t . Thus, letting $\mathbf{G}_t := \mathbf{I}^k - \mathbf{V}_{S_t}^t (\mathbf{V}_{S_t}^t)^\top$ to simplify the notation, we have

$$\begin{aligned}
 23 f(S_{t-1}) &= f(S_{t-1}, \mathbf{V}_{S_{t-1}}^{t-1}) = \left\| \mathbf{G}_t (\mathbf{A}^{1/2})_{S_{t-1}} \right\|_F^2 + \sum_{j \in S_{t-1}^C} \left\| (\mathbf{A}^{1/2})_j \right\|_2^2 \\
 24 &= \left\| \mathbf{G}_t \mathbf{A}_{S_t}^{1/2} \right\|_F^2 + \sum_{j \in S_t^C} \left\| \mathbf{A}_j^{1/2} \right\|_2^2 + \underbrace{\Delta_{j^{\text{in}}}^{\text{in}} - \Delta_{j^{\text{out}}}^{\text{out}}}_{>0} \\
 25 &> \left\| \mathbf{G}_t \mathbf{A}_{S_t}^{1/2} \right\|_F^2 + \sum_{j \in S_t^C} \left\| \mathbf{A}_j^{1/2} \right\|_2^2 \\
 26 &= f(S_t, \mathbf{V}_{S_t}^t) = f(S_t).
 \end{aligned}$$

27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Algorithm 1 Modified greedy neighborhood search

Input: Covariance matrix \mathbf{A} , sparsity parameter k , number of maximum iterations T

Output: A feasible solution \mathbf{V} for rsPCA.

Initialize with $S_0 \subseteq [d]$

Compute eigendecomposition of A_{S_0} : $\mathbf{A}_{S_0, S_0} = \mathbf{U}_{S_0} \mathbf{A}_{S_0} \mathbf{U}_{S_0}^\top$, $\mathbf{V}_{S_0} = (\mathbf{U}_{S_0})_{\star, [r]}$

for $t = 1, \dots, T$ **do**

 Compute the leaving candidate $j^{\text{out}} := \arg \min_{j \in S_{t-1}} \Delta_j^{\text{out}}$

 Compute the entering candidate $j^{\text{in}} := \arg \max_{j \in S_{t-1}^c} \Delta_j^{\text{in}}$

if $\Delta_{j^{\text{in}}}^{\text{in}} > \Delta_{j^{\text{out}}}^{\text{out}}$ **then**

 Set $S_t := S_{t-1} - \{j^{\text{out}}\} + \{j^{\text{in}}\}$

 Compute the eigenvalue decomposition $(\mathbf{A}^{1/2})_{S_t} = \mathbf{U}_{S_t} \mathbf{A}_{S_t} \mathbf{U}_{S_t}^\top$

 Set $\mathbf{V}_{S_t}^t = (\mathbf{U}_{S_t})_{\star, [r]}$

else

Return the matrix \mathbf{V} where in rows S_{t-1} equals $\mathbf{V}_{S_{t-1}}^{t-1}$ (i.e., $\mathbf{V}_{S_{t-1}} = \mathbf{V}_{S_{t-1}}^{t-1}$) and in rows S_{t-1}^c equals zero

end if

end for

C Techniques for reducing the running time of CIP

In practice, we want to reduce the running time of CIP. Here are the techniques that we used to enhance the efficiency in practice.

C.1 Threshold

The first technique is to reduce the number of SOS-II constraints. Let λ_{TH} be a threshold parameter that splits the eigenvalues $\{\lambda_j\}_{j=1}^d$ of sample covariance matrix \mathbf{A} into two parts $J^+ = \{j : \lambda_j > \lambda_{\text{TH}}\}$ and $J^- = \{j : \lambda_j \leq \lambda_{\text{TH}}\}$. The objective function $\text{Tr}(\mathbf{V}^\top \mathbf{A} \mathbf{V})$ satisfies

$$\text{Tr}(\mathbf{V}^\top \mathbf{A} \mathbf{V}) = \sum_{j \in J^+} (\lambda_j - \lambda_{\text{TH}}) \sum_{i=1}^r g_{ji}^2 + \sum_{j \in J^-} (\lambda_j - \lambda_{\text{TH}}) \sum_{i=1}^r g_{ji}^2 + \lambda_{\text{TH}} \sum_{j=1}^d \sum_{i=1}^r g_{ji}^2,$$

in which the first term is convex, the second term is concave, and the third term satisfies

$$\lambda_{\text{TH}} \sum_{j=1}^d \sum_{i=1}^r g_{ji}^2 \leq r \lambda_{\text{TH}} \quad (\text{threshold-term})$$

due to $\sum_{j=1}^d \sum_{i=1}^r g_{ji}^2 \leq r$. Since maximizing a concave function is equivalent to convex optimization, we replace the second term by a new auxiliary variable s and the third term by its upper bound $r \lambda_{\text{TH}}$ such that

$$\text{Tr}(\mathbf{V}^\top \mathbf{A} \mathbf{V}) \leq \sum_{j \in J^+} (\lambda_j - \lambda_{\text{TH}}) \sum_{i=1}^r g_{ji}^2 - s + r \lambda_{\text{TH}} \quad (\text{threshold-tech})$$

where

$$s \geq \sum_{j \in J^-} \underbrace{(\lambda_{\text{TH}} - \lambda_j)}_{\geq 0} \sum_{i=1}^r g_{ji}^2 \quad (\text{s-var})$$

is a convex constraint. We select a value of λ_{TH} so that $|J^+| = 3$. Therefore, it is sufficient to construct a piecewise-linear upper approximation for the quadratic terms g_{ji}^2 in the first term with $j \in J^+$, i.e., constraint set $\text{PLA}([J^+] \times [r])$. We thus, greatly reduce the number of SOS-II constraints from $\mathcal{O}(d \times r)$ to $\mathcal{O}(|J^+| \times r)$, i.e. in our experimnts to $3r$ SOS-II constraints.

C.2 Cutting planes

Similar to classical integer programming, we can incorporate additional cutting planes to improve the efficiency.

Cutting plane for sparsity: The first family of cutting-planes is obtained as follows: Since $\|\mathbf{V}\|_0 \leq k$ and $\mathbf{v}_1, \dots, \mathbf{v}_r$ are orthogonal, by Bessel inequality, we have

$$\sum_{i=1}^r g_{ji}^2 = \sum_{i=1}^r (\mathbf{a}_j^\top \mathbf{v}_i)^2 = \mathbf{a}_j^\top \mathbf{V} \mathbf{V}^\top \mathbf{a}_j \leq \theta_j^2, \quad (\text{sparse-g})$$

$$\sum_{i=1}^r \xi_{ji} \leq \theta_j^2 \left(1 + \frac{r}{4N^2}\right). \quad (\text{sparse-xi})$$

We call these above cuts–sparse cut since θ_j is obtained from the row sparsity parameter k .

Cutting plane from objective value: The second type of cutting plane is based on the property: for any symmetric matrix, the sum of its diagonal entries are equal to the sum of its eigenvalues. Let $\mathbf{A}_{j_1, j_1}, \dots, \mathbf{A}_{j_k, j_k}$ be the largest k diagonal entries of the sample covariance matrix \mathbf{A} , we have

Proposition 1 *The following are valid cuts for rsPCA:*

$$\sum_{j=1}^d \lambda_j \sum_{i=1}^r g_{ji}^2 \leq \mathbf{A}_{j_1, j_1} + \dots + \mathbf{A}_{j_k, j_k}. \quad (\text{cut-g})$$

When the splitting points $\{\gamma_{ji}^\ell\}_{\ell=-N}^N$ in SOS-II are set to be $\gamma_{ji}^\ell = \frac{\ell}{N} \cdot \theta_j$, we have:

$$\begin{aligned} \sum_{j \in J^+} (\lambda_j - \lambda_{\text{TH}}) \sum_{i=1}^r \xi_{ji} - s + g\lambda_{\text{TH}} &\leq \mathbf{A}_{j_1, j_1} + \dots + \mathbf{A}_{j_k, j_k} + \sum_{j \in J^+} \frac{r(\lambda_j - \phi)\theta_j^2}{4N^2} \\ g &\geq \sum_{j=1}^d \sum_{i=1}^r g_{ji}^2. \end{aligned} \quad (\text{cut-xi})$$

C.3 Implemented version of CIP

Thus the implemented version of CIP is

$$\begin{aligned} \max \quad & \sum_{j \in J^+} (\lambda_j - \lambda_{\text{LB}}) \sum_{i=1}^r \xi_{ji} - s + r\lambda_{\text{LB}} \\ \text{s.t.} \quad & \mathbf{V} \in \mathcal{CR2} \\ & (g, \xi, \eta) \in \text{PLA}' \\ & (\text{s-var}), (\text{sparse-g}), (\text{sparse-xi}), (\text{cut-g}), (\text{cut-xi}) \end{aligned} \quad (\text{CIP-impl})$$

C.4 Submatrix technique

Proposition 2 Let $X \in \mathbb{R}^{m \times n}$ and let θ be defined as

$$\theta := 2 \max_{\mathbf{V}^1 \in \mathbb{R}^{m \times r}, \mathbf{V}^2 \in \mathbb{R}^{n \times r}} 2 \text{Tr} \left((\mathbf{V}^1)^\top \mathbf{X} \mathbf{V}^2 \right) \text{ s.t. } (\mathbf{V}^1)^\top \mathbf{V}^1 + (\mathbf{V}^2)^\top \mathbf{V}^2 = \mathbf{I}^r,$$

then $\theta \leq \sqrt{r} \|X\|_F$

Proof

$$\begin{aligned} & \max_{\mathbf{V}^1, \mathbf{V}^2} 2 \text{Tr} \left((\mathbf{V}^1)^\top \mathbf{X} \mathbf{V}^2 \right) \text{ s.t. } (\mathbf{V}^1)^\top \mathbf{V}^1 + (\mathbf{V}^2)^\top \mathbf{V}^2 = \mathbf{I}^r, \\ \Leftrightarrow & \max_{\mathbf{V}^1, \mathbf{V}^2} \text{Tr} \left((\mathbf{V}^1)^\top (\mathbf{V}^2)^\top \begin{pmatrix} 0 & \mathbf{X} \\ \mathbf{X}^\top & 0 \end{pmatrix} \begin{pmatrix} \mathbf{V}^1 \\ \mathbf{V}^2 \end{pmatrix} \right) \text{ s.t. } (\mathbf{V}^1)^\top \mathbf{V}^1 + (\mathbf{V}^2)^\top \mathbf{V}^2 = \mathbf{I}^r, \\ \Leftrightarrow & \max_{\mathbf{V}} \text{Tr} \left(\mathbf{V}^\top \begin{pmatrix} 0 & \mathbf{X} \\ \mathbf{X}^\top & 0 \end{pmatrix} \mathbf{V} \right) \text{ s.t. } \mathbf{V}^\top \mathbf{V} = \mathbf{I}^r. \end{aligned}$$

Note that the final maximization problem is equal to

$$\begin{aligned} & \max_{\mathbf{V}} \text{Tr} \left(\mathbf{V}^\top \begin{pmatrix} 0 & \mathbf{X} \\ \mathbf{X}^\top & 0 \end{pmatrix} \mathbf{V} \right) \text{ s.t. } \mathbf{V}^\top \mathbf{V} = \mathbf{I}^r \\ & \leq \sum_{i=1}^r \lambda_i \left(\begin{pmatrix} 0 & \mathbf{X} \\ \mathbf{X}^\top & 0 \end{pmatrix} \right), \end{aligned}$$

Next we verify that the eigenvalues of

$$\begin{pmatrix} 0 & X \\ X^\top & 0 \end{pmatrix}$$

are \pm singular values of X : Let $X = U \Sigma W^\top$. In particular, note that:

$$\begin{aligned} \begin{pmatrix} 0 & U \Sigma W^\top \\ W \Sigma U^\top & 0 \end{pmatrix} \begin{bmatrix} u_i \\ w_i \end{bmatrix} &= \begin{bmatrix} U \Sigma e_i \\ W \Sigma e_i \end{bmatrix} = \sigma_i(X) \begin{bmatrix} u_i \\ w_i \end{bmatrix} \\ \begin{pmatrix} 0 & U \Sigma W^\top \\ W \Sigma U^\top & 0 \end{pmatrix} \begin{bmatrix} u_i \\ -w_i \end{bmatrix} &= \begin{bmatrix} -U \Sigma e_i \\ W \Sigma e_i \end{bmatrix} = -\sigma_i(X) \begin{bmatrix} u_i \\ -w_i \end{bmatrix}. \end{aligned}$$

Therefore, we have

$$\sum_{i=1}^r \lambda_i \left(\begin{pmatrix} 0 & \mathbf{X} \\ \mathbf{X}^\top & 0 \end{pmatrix} \right) = \sum_{i=1}^r \sigma_i(\mathbf{X}) \leq \sqrt{r} \|\mathbf{X}\|_F.$$