# Using $\ell_1$-relaxation and integer programming to obtain dual bounds for sparse PCA

Santanu S. Dey

School of Industrial and Systems Engineering, Georgia Institute of Technology, santanu.dey@isye.gatech.edu

Rahul Mazumder

Operations Research Center, Massachusetts Institute of Technology, rahulmaz@mit.edu

Guanyi Wang

School of Industrial and Systems Engineering, Georgia Institute of Technology, gwang93@gatech.edu

Principal component analysis (PCA) is one of the most widely used dimensionality reduction tools in scientific data analysis. The PCA direction, given by the leading eigenvector of a covariance matrix, is a linear combination of all features with nonzero loadings—this impedes interpretability. Sparse principal component analysis (SPCA) is a framework that enhances interpretability by incorporating an additional sparsity requirement in the feature weights (factor loadings) while finding a direction that explains the maximal variation in the data. However, unlike PCA, the optimization problem associated with the SPCA problem is NP-hard. Most conventional methods for solving SPCA are heuristics with no guarantees such as certificates of optimality on the solution-quality via associated dual bounds. Dual bounds are available via standard semidefinite programming (SDP) based relaxations, which may not be tight and the SDPs are difficult to scale using off-the-shelf solvers. In this paper, we present a convex integer programming (IP) framework to derive dual bounds. At the heart of our approach is the so-called $\ell_1$-relaxation of SPCA. While the $\ell_1$-relaxation leads to convex optimization problems for $\ell_0$-sparse linear regression and relatives; it results in a non-convex optimization problem for the PCA problem. We first show that the $\ell_1$-relaxation gives tight multiplicative bound on SPCA. Then we show how to use standard integer programming techniques to further relax the $\ell_1$-relaxation into a convex IP, for which there are good commercial solvers. We present worst-case results on the quality of the dual bound provided by the convex IP. We empirically observe that the dual bounds are significantly better than worst-case performance, and are superior to the SDP bounds on some real-life instances. Moreover, solving the convex IP model using commercial IP solvers appears to scale much better than solving the SDP-relaxation using commercial solvers. To the best of our knowledge, we obtain the best dual bounds for real and artificial instances for SPCA problems involving covariance matrices of size up to $2000 \times 2000$.

*Key words*: $\ell_1$ relaxation, Dual bounds, Sparse principal component analysis

## 1. Introduction

Principal component analysis (PCA) is one of the most widely used dimensionality reduction methods in data science. Given a data matrix $Y \in \mathbb{R}^{m \times n}$ (with $m$ samples and $n$ features; and each feature is centered to have zero mean), PCA seeks to find a principal component (PC) direction

$x \in \mathbb{R}^n$ with $\|x\|_2 = 1$ that maximizes the variance of a weighted combination of features. Formally, this PC direction can be found by solving

$$\max_{\|x\|_2=1} x^\top A x \tag{PCA}$$

where $A \triangleq \frac{1}{m} Y^\top Y$ is the sample covariance matrix. An obvious drawback of PCA is that all the entries of $\hat{x}$ (an optimal solution to (PCA)) are (usually) nonzero, which leads to the PC direction being a linear combination of all features – this impedes interpretability [11, 23, 41]. In biomedical applications for example, when $Y$ corresponds to the gene-expression measurements for different samples, it is desirable to obtain a PC direction which involves only a handful of the features (e.g, genes) for interpretation purposes. In financial applications (e.g, $A$ may denote a covariance matrix of stock-returns), a sparse subset of stocks that are responsible for driving the first PC direction may be desirable for interpretation purposes. Indeed, in many scientific and industrial applications [1, 20, 36], for additional interpretability, it is desirable for the factor loadings to be sparse, i.e., few of the entries in $\hat{x}$ are nonzero and the rest are zero. This motivates the notion of a sparse principal component analysis (SPCA) [20, 23], wherein, in addition to maximizing the variance, one also desires the direction of the first PC to be sparse in the factor loadings. The most natural optimization formulation of this problem, modifies criterion PCA with an additional sparsity constraint on $x$ leading to:

$$\lambda^k(A) \triangleq \max_{\|x\|_2=1, \|x\|_0 \leq k} x^\top A x \tag{SPCA}$$

where $\|x\|_0 \leq k$, is equivalent to allowing at most $k$ components of $x$ to be nonzero. Unlike the PCA problem, the SPCA problem is NP-hard [12, 27].

Many heuristic algorithms have been proposed in the literature that use greedy methods [23, 40, 21, 18], alternating methods [38] and the related power methods [24]. However, conditions under which (some of) these computationally friendlier methods can be shown to work well, make very strong and often unverifiable assumptions on the problem data. Therefore, the performance of these heuristics (in terms of how close they are to an optimal solution of the SPCA problem) on a given dataset is not clear.

Since SPCA is NP-hard, there has been exciting work in the statistics community [4, 35] in understanding the statistical properties of convex relaxations (e.g., those proposed by [13] and variants) of SPCA. It has been established [4, 35] that the statistical performance of estimators available from convex relaxations are sub-optimal (under suitable modeling assumptions) when compared to estimators obtained by (optimally) solving SPCA—this further underlines the importance of creating tools to be able to solve SPCA to optimality.

Our main goal in this paper is to propose an integer programming framework that allows the computation of certificates of optimality via dual bounds, which make limited restrictive/unverifiable assumptions on the data. Dual bounds can also translate into suitable guarantees for statistical performance of the estimator—see for example, [28][Theorem 4] for results pertaining to approximate solutions for sparse regression settings[1]. To the best of our knowledge, the only published methods for obtaining dual bounds of SPCA are based on semidefinite programming (SDP) relaxations [15, 17, 18, 39] (see Appendix B for the SDP relaxation) and spectral methods involving a low-rank approximation of the matrix $A$ [30]. Both these approaches however, have some limitations. The SDP relaxation does not appear to scale easily (using off-the-shelf solver Mosek 8.0.0.60) for matrices with more than a few hundred rows/columns, while applications can be significantly larger. Indeed, even a relatively recent implementation based on the Alternating Direction Method of Multipliers for solving the SDP considers instances with $n \approx 200$ [26]. The spectral methods involving a low-rank approximation of $A$ proposed in [30] have a running time of $\mathcal{O}(n^d)$ where $d$ is the rank of the matrix—in order to scale to large instances, no more than a rank 2 approximation of the original matrix seems possible. The paper [3] presents a specialized branch and bound solver[2] to obtain solutions to the SPCA problem, but their method can handle problems with $n \approx 100$ – the approach presented here is different, and our proposal scales to problem instances that are much larger.

The methods proposed here are able to obtain approximate dual bounds of SPCA by solving convex integer programs and a related perturbed version of convex integer programs that are easier to solve. The dual bounds we obtain are incomparable to dual bounds based on the SDP relaxation, i.e. neither dominates the other, and the method appears to scale well to matrices up to sizes of $2000 \times 2000$.

## 2. Main results

In this paper, we use upper case letters such as $A, X$ to denote symmetric matrices. The $(i, j)$-th component of matrix $A$ is denoted as $[A]_{ij}$ or $A_{ij}$ in short. We use lower case letters such as $v, x$ for vectors, and denote the $i$-th component of a vector $v$ as $[v]_i$ or $v_i$ in short. We use upper case letter $I$ for set of indices. Given a vector where $v \in \mathbb{R}^n$ and $I \subseteq [n]$, we let $v_I \in \mathbb{R}^n$ to be the vector:

$$[v_I]_i = \begin{cases} v_i & i \in I \\ 0 & i \notin I \end{cases}$$

---

[1] In [28], estimators with certificates on dual bounds translate to simple modifications of error bounds that correspond to the global solution of the original nonconvex estimator.

[2] This paper is not available in the public domain at the time of writing this paper.

We use the usual notation $\|\cdot\|_1$, $\|\cdot\|_2$ for $\ell_1$, $\ell_2$ norm respectively for a given vector. Let $\|\cdot\|_0$ be the $\ell_0$ norm which denotes the number of non-zero components. Given a set $S$, we denote $\text{conv}(S)$ as the convex hull of $S$; given a positive integer $n$ we denote $\{1, \ldots, n\}$ by $[n]$; given a matrix $A$, we denote its trace by $\text{tr}(A)$. Given $n$ scalars $v_1, \ldots v_n$, $\text{diag}(v_1, \ldots, v_n)$ is the $n \times n$ matrix whose diagonal elements are $v_i$'s and the off-diagonal terms are equal to 0. We list all the notation used in this paper in Table 13.

Notice that the constraint $\|x\|_2 = 1, \|x\|_0 \leq k$ implies that $\|x\|_1 \leq \sqrt{k}$. Thus, one obtains the so-called $\ell_1$-norm relaxation of SPCA:

$$\text{OPT}_{\ell_1} \triangleq \max_{\|x\|_2 \leq 1, \|x\|_1 \leq \sqrt{k}} x^\top A x. \qquad (\ell_1\text{-relax})$$

The relaxation $\ell_1$-relax has two advantages:

(a) As shown in Theorem 1 below, $\ell_1$-relax gives a constant factor bound on SPCA,

(b) The feasible region is convex and all the nonconvexity is in the objective function.

We build on these two advantages: our convex IP relaxation is a further relaxation of $\ell_1$-relax (together with some implied linear inequalities for SPCA) which heavily use the fact that the feasible region of $\ell_1$-relax is convex. We require to use IP methods and construct the convex IP, since the objective of $\ell_1$-relax is non-convex. Thus, we use a combination of $\ell_1$-relax and IP methods to obtain strong dual bounds.

We note that $\ell_1$-relax is an important estimator in its own right [20, 36]—it is commonly used in the statistics/machine-learning community as one that leads to an eigenvector of $A$ with entries having a small $\ell_1$-norm (as opposed to a small $\ell_0$-norm). We emphasize that $\ell_1$-relaxation has never been used to computationally obtain dual bounds for SPCA. Indeed, to the best of our knowledge there has been no systematic study of the theoretical and empirical computational properties of the $\ell_1$-relaxation vis-à-vis SPCA.

The rest of this section is organized as follows: In Section 2.1, we present the constant factor bound on SPCA given by $\ell_1$-relax, improving upon some known results. In Section 2.2, we present the construction of our convex IP and prove results on the quality of bound provided. In Section 2.3, we discuss perturbing the original matrix in order to make the convex IP more efficiently solvable while still providing reasonable dual bounds. In Section 4, we present results from our computational experiments.

## 2.1. Quality of $\ell_1$-relaxation as a surrogate for the SPCA problem

The following theorem is an improved version of a result appearing in [34] (Exercise 10.3.7).

THEOREM 1. *The objective value $OPT_{\ell_1}$ is upper bounded by a multiplicative factor $\rho^2$ away from $\lambda^k(A)$, i.e., $\lambda^k(A) \leq OPT_{\ell_1} \leq \rho^2 \cdot \lambda^k(A)$ with $\rho \leq 1 + \sqrt{\frac{k}{k+1}}$.*

Proof of Theorem 1 is provided in Section 3. While we have improved upon the bound presented in [34], we do not know if this new bound is tight.

The approximation ratio $1 + \sqrt{\frac{k}{k+1}}$ from Theorem 1 yields an almost 100% gap (see formal definition of gap in Section 4) in the worst case. From a practitioners' viewpoint, a 100% gap is obviously far from ideal and would not be considered as "solving" the problem. However, as we shall see in Section 4, the $\ell_1$-relaxation does provide very good dual bounds in many instances. Moreover, as stated above the approximation ratio of $1 + \sqrt{\frac{k}{k+1}}$ is the best we can prove; however this bound may be significantly away from the actual bound.

Theorem 1 has implications regarding existence of polynomial-time algorithms to obtain a constant-factor approximation guarantee for $\ell_1$-relax. In particular, the proof of Theorem 1 implies that if one can obtain a solution for $\ell_1$-relax which is within a constant factor, say $\theta$, of $\text{OPT}_{\ell_1}$, then a solution for SPCA problem can be obtained, which is within a constant factor (at most $\theta \rho \approx 4\theta$) of $\lambda^k(A)$. Therefore, the $\ell_1$-relaxation is also inapproximable in general.

## 2.2. From $\ell_1$-relaxation to convex integer programming model

A classical integer programming approach to finding dual bounds of SPCA would be to go to an extended space involving the product of $x$-variables and include one binary variable per $x$-variable in order to model the $\ell_0$-norm constraint, resulting in a very large number of binary variables. In particular, a typical model could be of the form:

$$\max \quad \text{tr}(AX) \tag{1}$$

$$\text{s.t.} \ -z_i \le x_i \le z_i, \ i \in [n] \tag{2}$$

$$\sum_{j=1}^{n} z_i \le k \tag{3}$$

$$\|x\|_2 \le 1 \tag{4}$$

$$\begin{bmatrix} 1 & x^\top \\ x & X \end{bmatrix} \succeq 0 \tag{5}$$

$$\text{rank}\left(\begin{bmatrix} 1 & x^\top \\ x & X \end{bmatrix}\right) = 1 \tag{6}$$

$$z \in \{0,1\}^n. \tag{7}$$

It is easy to see that such a model is challenging due to (a) $n$ binary variables (b) "quadratic" increase in number of variables $(X)$ and (c) the presence of the rank constraint. Even with significant progress, it is well-known that solving such problems beyond $n$ being a few hundred variables is extremely challenging [5, 19]. Indeed, instances with an arbitrary quadratic objective and bound constraints cannot be generally solved (exactly) by modern state-of-the-art methods as soon as the number of variables exceed a hundred or so [10, 7].

This is how we address the challenges discussed above.

1. $n$ binary variables (a): the feasible region of $\ell_1$-relax is a convex set. Therefore, we do not have to include binary variables to model the $\ell_0$-norm constraint. We will use $\ell_1$-relax as our basic relaxation.

2. Quadratic increase in number of variables (b) and rank constraint (c): We do not use the $X$ variables to model the quadratic objective. Instead we upper bound the quadratic objective using piecewise linear function via integer programming techniques.

In other words, since the feasible region of $\ell_1$-relax is a convex set and takes care of challenge (a), we model/upper bound the objective function using IP techniques to deal with challenges (b) and (c). Specifically, we follow the following procedure:

**step-0**: By spectral decomposition, let $A = \sum_{i=1}^{n} \lambda_i v_i v_i^\top$ where $(\lambda_i)_{i=1}^n, (v_i)_{i=1}^n$ are unit norm orthogonal eigen-pairs. Then the objective function of $\ell_1$-relax is:

$$\sum_{i=1}^{n} \lambda_i (x^\top v_i)^2.$$

**step-1**: Assuming that $\lambda \leq \lambda^k(A)$, we have that $x^\top A x = x^\top (A - \lambda I)x + \lambda$ for $x$ such that $\|x\|_2 = 1$, where $I$ is the identity matrix. Therefore, if we split the eigenvalues into two sets as $\{i : \lambda_i > \lambda\}$ and $\{i : \lambda_i < \lambda\}$, the objective function can be represented as

$$\lambda + \sum_{i \in \{i:\lambda_i > \lambda\}} (\lambda_i - \lambda)(x^\top v_i)^2 + \sum_{i \in \{i:\lambda_i < \lambda\}} (\lambda_i - \lambda)(x^\top v_i)^2$$

where for each eigenvalue $\lambda_i$ that equals to $\lambda$, since $\lambda_i - \lambda = 0$, it does not contribute anything to objective function. Note that the first term is convex and the second term is concave. Since the objective is a maximizing, we need to deal with the first term. This idea of splitting the objective function into convex and concave part is a well-studied approach for attacking non-convex quadratic objective functions. See for example [6, 9] for use of some similar ideas.

**step-2**: For each index $i \in \{i : \lambda_i > \lambda\}$, replace $x^\top v_i$ with a single continuous variable $g_i$, and set $\theta_i \leftarrow \max\{x^\top v_i : \|x\|_2 \leq 1, \|x\|_0 \leq k\}$ (or $\theta_i \leftarrow \max\{x^\top v_i : \|x\|_2 \leq 1, \|x\|_1 \leq \sqrt{k}\}$ if we explicitly want a relaxation of $\ell_1$-relax) be an upper bound of $g_i$. Then for each $g_i$ with $i \in \{i : \lambda_i > \lambda\}$, construct a piecewise linear upper approximation $\xi_i$ for $g_i^2$. Such piecewise linear upper approximation is usually modelled via *special ordered sets of type 2* (SOS-2) constraints [29].

**step-3**: For $\sum_{i \in \{i:\lambda_i < \lambda\}} (\lambda_i - \lambda)(x^\top v_i)^2$, since $\lambda_i - \lambda < 0$, we obtain a convex constraint $\sum_{i \in \{i:\lambda_i < \lambda\}} -(\lambda_i - \lambda)(x^\top v_i)^2 \leq s$.

Therefore, a convex integer programming problem is obtained as follows:

$$
\begin{aligned}
\text{OPT}_{\text{convex-IP}} \triangleq \max \quad & \lambda + \sum_{i \in \{i : \lambda_i > \lambda\}} (\lambda_i - \lambda)\xi_i - s \\
\text{s.t.} \quad & \begin{cases} g_i = x^\top v_i \\ -\theta_i \le g_i \le \theta_i \end{cases} \hspace{4cm} i \in [n] \\
& \begin{cases} g_i = \sum_{j=-N}^{N} \gamma_i^j \eta_i^j \\ \xi_i = \sum_{j=-N}^{N} (\gamma_i^j)^2 \eta_i^j \\ (\eta_i^{-N}, \dots, \eta_i^N) \in \text{SOS-2} \end{cases} \hspace{1.3cm} i \in \{i : \lambda_i > \lambda\} \hspace{1cm} \text{(Convex-IP)} \\
& \begin{cases} \sum_{i=1}^{n} x_i^2 \le 1 \\ \sum_{i \in \{i:\lambda_i > \lambda\}} \left( \xi_i - \frac{\theta_i^2}{4N^2} \right) + \sum_{i \in \{i:\lambda_i \le \lambda\}} g_i^2 \le 1 \\ \sum_{i=1}^{n} y_i \le \sqrt{k} \\ y_i \ge x_i, \ y_i \ge -x_i, \ \forall i \in [n] \\ \sum_{i \in \{i:\lambda_i < \lambda\}} -(\lambda_i - \lambda)g_i^2 \le s \end{cases}
\end{aligned}
$$

**Notations and explanations of Convex-IP:**

*Variable $g_i$:* The first set of constraints

$$
\begin{cases} g_i = x^\top v_i \\ -\theta_i \le g_i \le \theta_i \end{cases}
$$

transfers $x^\top v_i$ into a single variable for each $i \in [n]$.

*Variable $\xi_i$:* Based on step-2 above, for each $i \in \{i : \lambda_i > \lambda\}$, the second set of constraints

$$
\begin{cases} g_i = \sum_{j=-N}^{N} \gamma_i^j \eta_i^j \\ \xi_i = \sum_{j=-N}^{N} (\gamma_i^j)^2 \eta_i^j \\ (\eta_i^{-N}, \dots, \eta_i^N) \in \text{SOS-2} \end{cases}
$$

forms $\xi_i$ as a piecewise-linear upper approximation of $g_i^2$. Let $2N+1$ be the number of splitting points of the domain $[-\theta_i, \theta_i]$ of variable $g_i$, where the set of splitting points $(\gamma_i^j)_{j=-N}^{N}$ satisfy

$$
-\theta_i = \gamma_i^{-N} < \dots \gamma_i^0 \ (=0) < \dots < \gamma_i^N = \theta_i.
$$

Without any prior information of the optimal solution, we partition the set $[-\theta_i, \theta_i]$ equally to minimize the (worst-case) upper bounds, i.e., by letting $(\gamma_i^j)_{j=-N}^{N} \leftarrow \left( \frac{j}{N} \cdot \theta_i \right)_{j=-N}^{N}$ be the value of $j^{\text{th}}$ splitting point. See Section D for details.

*Quadratic constraints:* The third set of constraints does the following: Since $v_i$'s are orthogonal, then $\sum_{i=1}^{n} x_i^2 \le 1$ implies $\sum_{i=1}^{n} g_i^2 \le 1$. Together with $\xi_i$ representing $g_i^2$, we can obtain the implied inequality:

$$
\sum_{i \in \{i:\lambda_i > \lambda\}} \xi_i + \sum_{i \in \{i:\lambda_i \le \lambda\}} g_i^2 \le 1 + \sum_{i \in \{i:\lambda_i > \lambda\}} \frac{\theta_i^2}{4N^2}
$$

The second term in the right-hand-side reflects the fact that $\xi_i$ is not exactly equal to $g_i^2$, but only a piecewise linear upper bound of $g_i^2$. Note that the exact value of the second term in the right-hand-side also depends on the way one splits the set $[-\theta_i, \theta_i]$, the value $\sum_{i \in \{i:\lambda_i > \lambda\}} \frac{\theta_i^2}{4N^2}$ in

above formula is obtained via splitting $[-\theta_i, \theta_i]$ equally, which can be shown as the minimum upper bounds without any prior idea of the optimal solution $x$ of SPCA or $\ell_1$-relax. See the proof in Section D for details. This constraint (cutting-plane) is not necessarily needed for a correct model – it is used since it helps improving the dual bound of the LP relaxation and significantly improves the running-time of the solver.

$\ell_1$ *constraints:* The fourth set of constraints (the fourth one within the curly brackets in Convex IP) introduce new variables $y_i$ to denote $|x_i|$ for $i = 1, \ldots, n$ and model the constraint

$$\sum_{i=1}^{n} |x_i| \leq \sqrt{k}.$$

*Convex constraint:* The final constraint

$$\sum_{i \in \{i : \lambda_i < \lambda\}} -(\lambda_i - \lambda)g_i^2 \leq s \qquad \text{(convex-constraint)}$$

is a convex constraint that we obtained in step-3 where $x^\top v_i$ is replaced by a variable $g_i$ since $g_i = x^\top v_i$.

We arrive at the following result:

PROPOSITION 1. *The optimal objective value $OPT_{convex\text{-}IP}$ of Convex-IP is an upper bound on the SPCA problem.*

Proposition 1 is formally verified in Appendix C.

Next combining the result of Theorem 1 with the quality of the approximation of the objective function of $\ell_1$-relax by Convex-IP, we obtain the following result:

PROPOSITION 2. *The optimal objective value $OPT_{convex\text{-}IP}$ of Convex-IP is upper bounded by*

$$OPT_{\text{convex-IP}} \leq \rho^2 \lambda^k(A) + \frac{1}{4N^2} \sum_{i \in \{i : \lambda_i > \lambda\}} (\lambda_i - \lambda)\theta_i^2.$$

A proof of Proposition 2 is presented in Appendix D.

Finally, let us discuss why we expect Convex-IP to be appealing from a computational viewpoint. Unlike typical integer programming approaches, the number of binary variables in Convex-IP is $(2N + 1) \cdot |\{i : \lambda_i > \lambda\}|$ which is usually significantly smaller than $n$. Indeed, heuristics for SPCA generally produce good values of $\lambda$, and in almost all experiments we found that $|\{i : \lambda_i > \lambda\}| \ll n$. Moreover, $N$ is a parameter we control. In order to highlight the "computational tractability" of Convex-IP, we formally state the following result:

PROPOSITION 3. *Assuming the number of splitting points $N$ and the size of set $\{i : \lambda_i > \lambda\}$ is fixed, the Convex-IP problem can be solved in polynomial time.*

Note that the convex integer programming method which is solvable in polynomial time, does not contradict the inapproxamability of the SPCA problem, since $OPT_{\text{convex-IP}}$ is upper bounded by the sum of $\rho^2 \lambda^k(A)$ and a term corresponding to the sample covariance matrix.

### 2.3. Improving the running time of Convex-IP

**2.3.1. Perturbation of the covariance matrix** $A$**:** In practice, we do the following (sequence of) perturbation on covariance matrix $A$ to reduce the running time of solving convex IP. Again let $\lambda$ (obtained from some heuristic method) be a lower bound on the $\lambda^k(A)$, let $A = \sum_{i=1}^n \lambda_i v_i v_i^\top$ be the spectral decomposition of $A$ with $\lambda_1 \geq \ldots \geq \lambda_n \geq 0$.

1. Set $\bar{\lambda} \triangleq \max\{\lambda_i : \lambda_i \leq \lambda\}$ (where $\lambda_1, \lambda_2, \ldots, \lambda_n$ are the eigenvalues of $A$). We assume $\bar{\lambda} < \lambda$. However, when $\bar{\lambda} \triangleq \max\{\lambda_i : \lambda_i \leq \lambda\} = \lambda$, one can apply Algorithm 1 to obtain a matrix $\bar{A} \succeq A$ such that none of the eigenvalues of $\bar{A}$ equals $\lambda$. We then replace $A$ by $\bar{A}$. Now letting $\lambda_1, \lambda_2, \ldots, \lambda_n$ to be the eigenvalues of (the updated) $A$ and $\bar{\lambda} \triangleq \max\{\lambda_i : \lambda_i \leq \lambda\}$, we obtain that $\bar{\lambda} < \lambda$ for $\bar{A}$.

---

**Algorithm 1** Perturbation of $A$

---

1: *Input*: Sample covariance matrix $A$ and $\lambda$.

2: *Output*: A perturbed sample covariance matrix $\bar{A}$ with distinct eigenvalues such that $\bar{A} \succeq A$ and none of the eigenvalues of $\bar{A}$ equals $\lambda$.

3: **function** PERTURBATION METHOD$(A, \lambda)$

4:     Compute spectral decomposition on $A$ as $A = V^\top \Lambda V$, where $\Lambda = \mathrm{diag}(\lambda_1, \ldots, \lambda_n)$. Let $\lambda_{i_1} > \cdots > \lambda = \lambda_{i_j} > \cdots \lambda_{i_p} \geq 0$ be all its distinct values of eigenvalues where $p \leq n$.

5:     Set $\Delta\lambda \leftarrow \min\{\lambda_{i_j} - \lambda_{i_{j+1}} \mid j = 1, \ldots, p-1\}$.

6:     Set $\bar{\Lambda} \leftarrow \Lambda + \mathrm{diag}\left(\frac{i-1}{n}\epsilon \mid i = n, \ldots, 1\right)$ with $\epsilon = \frac{1}{2}\Delta\lambda$.

7:     **return** $\bar{A} \leftarrow V^\top \bar{\Lambda} V$.

8: **end function**

---

2. Perturb the covariance matrix $A = \sum_{i=1}^n \lambda_i v_i v_i^\top$ by $\bar{A} = \sum_{i \in \{i:\lambda_i > \lambda\}} \lambda_i v_i v_i^\top + \sum_{i \in \{i:\lambda_i \leq \lambda\}} \bar{\lambda} v_i v_i^\top$. Note that the objective value $\mathrm{OPT}_{\mathrm{convex\text{-}IP}}(\bar{A})$ in Convex-IP is an upper bound on $\mathrm{OPT}_{\mathrm{convex\text{-}IP}}(A)$. This is because if $(x, y, g, \xi, \eta, s)$ is a feasible solution of Convex-IP, then the objective function value of Convex-IP corresponding to $\bar{A}$ is at least as large as that of $A$. Replace $A$ by $\bar{A}$.

3. Therefore, the convex constraint $\sum_{i \in \{i:\lambda_i \leq \lambda\}} -(\lambda_i - \lambda)g_i^2 \leq s$ in Convex-IP can be replaced by $\sum_{i \in \{i:\lambda_i \leq \lambda\}} -(\bar{\lambda} - \lambda)g_i^2 \leq s$, i.e., $\sum_{i \in \{i:\lambda_i \leq \lambda\}} g_i^2 \leq \frac{s}{\lambda - \bar{\lambda}}$.

4. Let $(\bar{x}, \bar{y}, \bar{g}, \bar{\xi}, \bar{\eta}, \bar{s})$ be an optimal solution for Convex-IP. Since the convex constraint achieves equality for any optimal solution of Convex-IP, i.e.,

$$\sum_{i \in \{i:\lambda_i \leq \lambda\}} -(\lambda - \bar{\lambda})\bar{g}_i^2 = \bar{s}$$

together with

$$\sum_{i=1}^n \bar{g}_i^2 = \sum_{i \in \{i:\lambda_i \leq \lambda\}} \bar{g}_i^2 + \sum_{i \in \{i:\lambda_i > \lambda\}} \bar{g}_i^2 \leq 1$$

$$1 \leq \sum_{i \in \{i : \lambda_i > \lambda\}} \bar{\xi}_i + \sum_{i \in \{i : \lambda_i \leq \lambda\}} \bar{g}_i^2 \leq 1 + \frac{1}{4N^2} \sum_{i \in \{i : \lambda_i > \lambda\}} \theta_i^2,$$

imply the following inequalities:

$$1 - \frac{\bar{s}}{\lambda - \bar{\lambda}} \leq \sum_{i \in \{i : \lambda_i > \lambda\}} \bar{\xi}_i \leq 1 + \frac{1}{4N^2} \sum_{i \in \{i : \lambda_i > \lambda\}} \theta_i^2 - \frac{\bar{s}}{\lambda - \bar{\lambda}},$$

$$\sum_{i \in \{i : \lambda_i > \lambda\}} \bar{g}_i^2 \leq 1 - \frac{\bar{s}}{\lambda - \bar{\lambda}}.$$

Thus a simplified convex IP corresponding to the perturbed covariance matrix is:

$$
\begin{aligned}
\text{OPT}_{\text{pert-convex-IP}} \triangleq \max \quad & \lambda + \sum_{i \in \{i : \lambda_i > \lambda\}} (\lambda_i - \lambda)\xi_i - s \\
\text{s.t.} \quad & \begin{cases} g_i = x^\top v_i \\ -\theta_i \leq g_i \leq \theta_i \end{cases} \qquad i \in \{i : \lambda_i > \lambda\} \\
& \begin{cases} g_i = \sum_{j=-N}^{N} \gamma_i^j \eta_i^j \\ \xi_i = \sum_{j=-N}^{N} (\gamma_i^j)^2 \eta_i^j \\ (\eta_i^{-N}, \ldots, \eta_i^N) \in \text{SOS-2} \end{cases} \qquad i \in \{i : \lambda_i > \lambda\} \\
& \begin{cases} \sum_{i=1}^{n} x_i^2 \leq 1 \\ \sum_{i \in \{i : \lambda_i > \lambda\}} g_i^2 \leq 1 - \frac{s}{\lambda - \lambda} \\ 1 - \frac{s}{\lambda - \lambda} \leq \sum_{i \in \{i : \lambda_i > \lambda\}} \xi_i \leq 1 + \sum_{i \in \{i : \lambda_i > \lambda\}} \frac{\theta_i^2}{4N^2} - \frac{s}{\lambda - \lambda} \end{cases} \\
& \begin{cases} \sum_{i=1}^{n} y_i \leq \sqrt{k} \\ y_i \geq x_i, \ y_i \geq -x_i, \ \forall i \in [n] \end{cases} \\
& v^\top y \leq b_{(v)}
\end{aligned}
$$

$$\text{(Pert-Convex-IP)}$$

where the quadratic constraints in Pert-Convex-IP are updated based on the discussion above and the final constraint $v^\top y \leq b_{(v)}$ represents the cutting planes that we add, see Proposition 5 for details.

PROPOSITION 4. *The optimal objective value $OPT_{Pert\text{-}Convex\text{-}IP}$ is upper bounded by*

$$OPT_{\text{Pert-Convex-IP}} \leq \rho^2 \lambda^k(A) + \rho^2(\bar{\lambda} - \lambda_{\min}(A)) + \frac{1}{4N^2} \sum_{i \in \{i : \lambda_i > \lambda\}} (\lambda_i - \lambda)\theta_i^2.$$

Note that in Pert-Convex-IP, we do not need the variables $g_i, i \in \{i : \lambda_i \leq \lambda\}$ which greatly reduces the number of variables since in general $|\{i : \lambda_i \geq \lambda\}| \ll n$. In practice, we note a significant reduction in running time, while the dual bound obtained from Pert-Convex-IP model remains reasonable. More details are presented in Section 4.

**2.3.2. Refining the splitting points** Since the Pert-Convex-IP model runs much faster than the Convex-IP model, we run the Pert-Convex-IP model iteratively. In each new iteration, we add one extra splitting point describing each $\xi_i$ function. In particular, once we solve the Pert-Convex-IP model, we add one splitting point at the optimal value of $g_i$.

### 2.3.3. Cutting planes

PROPOSITION 5. *Let* $x \in \mathbb{R}^n$. *Let* $|x_{i_1}| \geq |x_{i_2}| \geq \cdots \geq |x_{i_{n-1}}| \geq |x_{i_n}|$. *Then let* $v$ *be the vector:*

$$v_{i_j} = \begin{cases} |x_{i_j}| & \text{if } j \leq k \\ |x_{i_k}| & \text{if } j > k. \end{cases} \tag{8}$$

*Also let* $b_{(v)} := \|(v_{i_1}, v_{i_2}, v_{i_3}, \ldots, v_{i_k})\|_2$. *The inequality*

$$v^\top y \leq b_{(v)}, \tag{9}$$

*is a valid inequality for SPCA.*

The validity of this inequality is clear: If $(x, y)$ is a feasible point, then the support of $y$ is at most $k$ and $\|y\|_2 \leq 1$. Therefore, $v^\top y \leq \|(v_{i_1}, v_{i_2}, v_{i_3}, \ldots, v_{i_k})\|_2 = b_{(v)}$. Notice that this inequality is not valid for $\ell_1$-relax. Also see [25].

We add these inequalities at the end of each iteration for the model where the seeding $x$ for constructing $v$ is chosen to be the optimal solution of the previous iteration.

## 3. Proof of Theorem 1

Given a vector $v \in \mathbb{R}^n$, we denote the $j^{th}$ coordinate of $v$ as $v_j$, and for some $J \subseteq [n]$ we denote the projection of $v$ onto the coordinates in the index set $J$ as $v_J$. Define

$$S_k \triangleq \{x \in \mathbb{R}^n \mid \|x\|_2 \leq 1, \|x\|_0 \leq k\}, \tag{10}$$

$$T_k \triangleq \{x \in \mathbb{R}^n \mid \|x\|_2 \leq 1, \|x\|_1 \leq \sqrt{k}\}. \tag{11}$$

Note that any $x \in T_k$ can be represented as a nonnegative combination of points in $S_k$, i.e., $x = x^1 + \cdots + x^m$ and $x^i \in S_k$ for all $i$. Here we think of each $x^i$ as a projection onto some unique $k$ components of $x$ and setting the other components to *zero*. Let $y^i = \frac{x^i}{\|x^i\|_2}$, then $y^i \in S_k$. Now we have, $x = \sum_{i=1}^{m} \|x^i\|_2 \cdot y^i$, and therefore

$$\frac{1}{\sum_{i=1}^{m} \|x^i\|_2} x = \sum_{i=1}^{m} \frac{\|x^i\|_2}{\sum_{i=1}^{m} \|x^i\|_2} \cdot y^i. \tag{12}$$

Thus, if we scale $x \in T_k$ by $\|x^1\|_2 + \ldots + \|x^m\|_2$, then the resulting vector belongs to $\text{conv}(S_k)$. Since we want this scaling factor to be as small as possible, we solve the following optimization problem:

$$\min \|x^1\|_2 + \ldots + \|x^m\|_2 : \ x = x^1 + \ldots + x^m; \ x^i \in S_k, \forall i \in [m]. \qquad \text{(Bound)}$$

Without loss of generality, we assume that $x \geq 0$ and $x_1 \geq x_2 \geq \cdots \geq x_n \geq 0$. Let $x = \bar{v}^1 + \ldots + \bar{v}^m$ where $v^1, \ldots, v^m \in S_k$ is an optimal solution of Bound. The following proposition presents a result on an optimal solution of Bound.

PROPOSITION 6. *Let $I^1, \ldots, I^m$ be a collection of supports such that: $I^1$ indexes the $k$ largest (in absolute value) components in $x$, $I^2$ indexes the second $k$ largest (in absolute value) components in $x$, and so on (note that $m = \lceil \frac{n}{k} \rceil$). Then $I^1, \ldots, I^m$ is an optimal set of supports for Bound.*

*Proof.* We prove this result by the method of contradiction. Suppose we have an optimal representation as $x = \bar{v}^1 + \cdots \bar{v}^m$ — and without loss of generality, we assume that $\|\bar{v}^1\|_2 \geq \cdots \geq \|\bar{v}^m\|_2$. Let $\bar{I}^1, \ldots, \bar{I}^m$ be the set of supports of $\bar{v}^1, \ldots, \bar{v}^m$ respectively, where we assume that the indices within each support vector are ordered such that

$$(x_{\bar{I}^j})_1 \geq (x_{\bar{I}^j})_2 \geq \cdots \geq (x_{\bar{I}^j})_g$$

for all $j \in \{1, \ldots, m\}$ (note that $g = k$ if $j < m$).

Let $\bar{I}^p$ be the first support that is different from $I^p$, i.e., $\bar{I}^1 = I^1, \ldots, \bar{I}^{p-1} = I^{p-1}$ and $\bar{I}^p \neq I^p$. Let $I_q^p$ be the first index in $I^p$ that does not belong to $\bar{I}^p$ with $q \leq k$ since $\|\bar{I}^p\|_0 = k$. Therefore, $I_q^p$ must be in $\bar{I}^{p'}$ where $p' > p$. Note now that by construction of $I$ and our assumption on $\bar{I}$, we have that $(x_{I^p})_q \geq (x_{\bar{I}^p})_q \geq (x_{\bar{I}^p})_k$. Now we exchange the index $I_q^p$ in $\bar{I}^{p'}$ with $\bar{I}_k^p$ in $\bar{I}^p$. We have:

$$\sqrt{\|x_{\bar{I}^p}\|_2^2 + ((x_{I^p})_q)^2 - ((x_{\bar{I}^p})_k)^2} + \sqrt{\|x_{\bar{I}^{p'}}\|_2^2 + ((x_{\bar{I}^p})_k)^2 - ((x_{I^p})_q)^2} \leq \|x_{\bar{I}^p}\|_2 + \|x_{\bar{I}^{p'}}\|_2, \qquad (13)$$

which holds because $\|x_{\bar{I}^p}\|_2 \geq \|x_{\bar{I}^{p'}}\|_2$ and $((x_{I^p})_q)^2 - ((x_{\bar{I}^p})_k)^2 \geq 0$.

Now repeating the above step, we obtain the result. $\square$

Based on Proposition 6, for any fixed $x \in T_k$, we can find out an optimal solution of Bound in closed form. Now we would like to know, for which vector $x$, the scaling factor $\|v^1\|_2 + \ldots + \|v^m\|_2$ will be the largest. Let $\rho$ be obtained by solving the following optimization problem:

$$\begin{aligned}
\rho = \max_x \ & \|x_{I^1}\|_2 + \cdots + \|x_{I^m}\|_2 \\
\text{s.t.} \quad & x = x_{I^1} + \cdots + x_{I^m} \\
& \|x\|_2^2 = \|x_{I^1}\|_2^2 + \cdots + \|x_{I^m}\|_2^2 \leq 1 \qquad \text{(Approximation ratio)} \\
& \|x\|_1 = \|x_{I^1}\|_1 + \cdots + \|x_{I^m}\|_1 \leq \sqrt{k} \\
& x_1 \geq \cdots \geq x_n \geq 0.
\end{aligned}$$

Then we obtain

$$T_k \subseteq \rho \cdot \mathrm{Conv}\,(S_k). \qquad (14)$$

Although the optimal objective value of Approximation ratio is hard to compute exactly, we can still find an upper bound.

LEMMA 1. *The objective value $\rho$ of Approximation ratio is bounded from above by $1 + \sqrt{\frac{k}{k+1}}$.*

*Proof.* First consider the case when $n \leq 2k$. In this case, $m \leq 2$. Consider the optimization problem:

$$\theta = \max \quad u + v$$

$$\text{s.t. } u^2 + v^2 \leq 1$$

If we think of $\|x_{I^1}\|_2$ as $u$ and $\|x_{I^2}\|_2$ as $v$, then we see that the above problem is a relaxation of Approximation ratio and therefore $\theta = \sqrt{2}$ is an upper bound on $\rho$. Noting that $\sqrt{2} \leq 1 + \sqrt{\frac{k}{k+1}}$ for all $k \geq 1$, we have the result.

Now we assume that $n > 2k$ and consequently $m > 2$. From Approximation ratio, let $\|x_{I^1}\|_1 = t$ and $\|x_{I^1}\|_2 = \gamma$. Based on the standard relationship between $\ell_1$ and $\ell_2$ norm, we have

$$\gamma \leq t \leq \sqrt{k}\gamma.$$

Since each coordinate of $x_{I^2}$ is smaller in magnitude than the average coordinate of $x_{I^1}$, we have

$$\|x_{I^2}\|_2 \leq \sqrt{\left(\frac{\|x_{I^2}\|_1}{k}\right)^2 k} = \frac{t}{\sqrt{k}}. \tag{15}$$

Also note that an alternative bound is given by

$$\|x_{I^2}\|_2 \leq \sqrt{1 - \gamma^2}.$$

Using an argument similar to the one used to obtain (15), we obtain that

$$\sum_{i=3}^{m} \|x_{I^i}\|_2 \leq \sum_{i=2}^{m-1} \sqrt{\left(\frac{\|x_{I^i}\|_1}{k}\right)^2 k} = \frac{1}{\sqrt{k}} \sum_{i=2}^{m-1} \|x_{I^i}\|_1 \leq \frac{\sqrt{k} - t}{\sqrt{k}}.$$

Therefore we obtain

$$\sum_{i=1}^{m} \|x_{I^i}\|_2 = \|x_{I^1}\|_2 + \|x_{I^2}\|_2 + \sum_{i=3}^{m} \|x_{I^i}\|_2 \leq \gamma + \min\left\{\frac{t}{\sqrt{k}}, \sqrt{1-\gamma^2}\right\} + 1 - \frac{t}{\sqrt{k}}. \quad \text{(Upper-Bound)}$$

Now we consider two cases:

1. If $\frac{t}{\sqrt{k}} \geq \sqrt{1-\gamma^2}$, then Upper-Bound becomes $\gamma + \sqrt{1-\gamma^2} + 1 - \frac{t}{\sqrt{k}}$. Since $\gamma \geq \frac{t}{\sqrt{k}} \geq \sqrt{1-\gamma^2}$, $\gamma$ satisfies $\gamma \geq \frac{1}{\sqrt{2}}$. Moreover we have that $t \geq \gamma, t \geq \sqrt{k(1-\gamma^2)}$. Since $\gamma \leq \sqrt{k(1-\gamma^2)}$ iff $\gamma \leq \sqrt{\frac{k}{k+1}}$ we obtain two cases:

$$\gamma + \sqrt{1-\gamma^2} + 1 - \frac{t}{\sqrt{k}} \leq \begin{cases} \gamma + \sqrt{1-\gamma^2} + 1 - \sqrt{1-\gamma^2} & \text{if } \gamma \in \left[\frac{1}{\sqrt{2}}, \sqrt{\frac{k}{k+1}}\right] \\ \gamma + \sqrt{1-\gamma^2} + 1 - \frac{\gamma}{\sqrt{k}} & \text{if } \gamma \in \left[\sqrt{\frac{k}{k+1}}, 1\right] \end{cases}$$

$$\leq \begin{cases} 1 + \sqrt{\frac{k}{k+1}} \\ 1 + \sqrt{\frac{k}{k+1}} \end{cases} \tag{16}$$

where (i) the first inequality holds when $\gamma = \sqrt{\frac{k}{k+1}}$, (ii) the second inequality holds since the function $f(\gamma) = \gamma + \sqrt{1 - \gamma^2} + 1 - \frac{\gamma}{\sqrt{k}}$ achieves (local and global) maximum at point $\gamma = \sqrt{\frac{k+1-2\sqrt{k}}{2k+1-2\sqrt{k}}}$ which is less than $\sqrt{\frac{k}{k+1}}$ for $k = 1, 2, \ldots$, thus $f(\gamma) \leq \max\left\{ f\left(\sqrt{\frac{k}{k+1}}\right), f(1) \right\} = 1 + \sqrt{\frac{k}{k+1}}$ for part $\gamma \in \left[\sqrt{\frac{k}{k+1}}, 1\right]$.

2. If $\frac{t}{\sqrt{k}} \leq \sqrt{1 - \gamma^2}$, then Upper-Bound becomes $\gamma + 1$. Note now that $\frac{\gamma}{\sqrt{k}} \leq \frac{t}{\sqrt{k}} \leq \sqrt{1 - \gamma^2}$, implies that $\gamma$ satisfies $\gamma \leq \sqrt{\frac{k}{k+1}}$. Therefore, $1 + \gamma \leq 1 + \sqrt{\frac{k}{k+1}}$.

Therefore, this upper bound holds.  $\square$

Therefore, we can show Theorem 1 holds.

*Proof of Theorem 1.* Since $T_k \subseteq \rho \cdot \text{Conv}(S_k)$ with $\rho \leq 1 + \sqrt{\frac{k}{k+1}}$ and the objective function is maximizing a convex function, we obtain that $\lambda^k(A) \leq \text{OPT}_{\ell_1} \leq \rho^2 \cdot \lambda^k(A)$.  $\square$

# 4. Numerical experiments

In this section, we report results on our empirical comparison of the performances of Convex-IP method, Pert-Convex-IP method and the SDP relaxation method.

## 4.1. Hardware and Software

All numerical experiments are implemented on MacBookPro13 with 2 GHz Intel Core i5 CPU and 8 GB 1867 MHz LPDDR3 Memory. Convex-IPs were solved using Gurobi 7.0.2. SDPs were solved using Mosek 8.0.0.60.

## 4.2. Obtaining primal solutions

We used a heuristic, which is very similar to the truncated power method [38], but has some advantages over the truncated power method. Given $v \in \mathbb{R}^n$, let $I_k(v)$ be the set of indices corresponding to the top $k$ entries of $v$ (in absolute value).

We start with a random initialization $x^0$ such that $\|x^0\|_2 = 1$, and set $I^0 \leftarrow I_k(V^\top x^0)$ where $V$ is a square root of $A$, i.e. $A = V^\top V$. In the $i^{\text{th}}$ iteration, we update

$$I^i \leftarrow I_k(V^\top x^i), \ x^{i+1} \leftarrow \underset{\|x\|_2 = 1}{\arg\max} \ x^\top A_{I^i} x \tag{17}$$

where $A_I \in \mathbb{R}^{n \times n}$ is the matrix with $[A_I]_{i,j} = [A]_{i,j}$ for all $i, j \in I$ and $[A_I]_{i,j} = 0$ otherwise. It is easy to see that $x^1, x^2, \ldots$ satisfy the condition $\|x\|_0 \leq k$. Moreover, using the fact $A$ is a PSD matrix, it is easy to verify that $(x^{i+1})^\top A x^{i+1} \geq (x^i)^\top A x^i$ for all $i$. Therefore, in each iteration, the above heuristic method leads to an improved feasible solution for the SPCA problem.

Our method has two clear advantages over the truncated power method:

• We use standard and efficient numerical linear algebra methods to compute eigenvalues of small $k \times k$ matrices.

• The termination criteria used in our algorithm is also simple: if $I^i = I^{i'}$ for some $i' < i$, then we stop. Clearly, this leads to a finite termination criteria.

In practice, we stop using a stopping criterion based on improvement and number of iterations instead of checking $I^i = I^{i'}$. Details are presented in Algorithm 2.

---

**Algorithm 2** Primal Algorithm

---

1: *Input*: Sample covariance matrix $A$, cardinality constraint $k$, initial vector $x^0$.

2: *Output*: A feasible solution $x^*$ of SPCA, and its objective value.

3: **function** HEURISTIC METHOD$(A, k, x^0)$

4:     Start with an initial (randomized) vector $x^0$ such that $\|x^0\|_2 = 1$ and $\|x^0\|_0 \leq k$.

5:     Set the initial current objective value $\text{Obj} \leftarrow (x^0)^\top A x^0$.

6:     Set the initial past objective value $\tilde{\text{Obj}} \leftarrow 0$.

7:     Set the maximum number of iterations be $i^{\max}$.

8:     **while** $\text{Obj} - \tilde{\text{Obj}} > \epsilon$ and $i \leq i^{\max}$ **do**

9:         Set $\tilde{\text{Obj}} \leftarrow \text{Obj}$.

10:        Set $I^i \leftarrow I_k(V^\top x^i)$.

11:        Set $x^{i+1} \leftarrow \arg\max_{\|x\|_2=1} x^\top A_{I^i} x$.

12:        Set $\text{Obj} \leftarrow (x^{i+1})^\top A x^{i+1}$.

13:     **end while**

14:     **return** $x^*$ as the final $x$ obtained from while-loop, and Obj.

15: **end function**

---

We use the values of $\epsilon = 10^{-6}$ and $i^{\max} = 20$ in our experiments in Algorithm 2. We repeat this algorithm with multiple random initializations. We repeat 20 times and take the best solution. We emphasize that Algorithm 2 may not lead to a global solution of SPCA.

Our Algorithm may also be interpreted as a version of the "alternating method" used regularly as a heuristic for bilinear programs as the sparse PCA problem can be equivalently rewritten as $\max\{x^\top A y \mid \|x\|_2 = \|y\|_2 = 1, \|x\|_0 \leq k, \|y\|_0 \leq k\}$. We have compared our primal method to two standard heuristics for finding primal feasible solutions of the sparse PCA problems in the literature: truncated power method (TPM, [37]), generalized power method (GPM, [24]) with $\ell_0$-penalty. The performances of all these methods are quite similar to our method (in terms of primal objective function values) on the real instances; see details in Appendix I.

## 4.3. Implementation of Convex-IP model and Pert-Convex-IP model

### 4.3.1. Deciding $\lambda$, $N$

1. Deciding $\lambda$: The size of the set $\{i : \lambda_i > \lambda\}$ denoted by $I_{\text{pos}}$ plays an important role for the computational tractability of our method. So our algorithm inputs an initial value, $I_{\text{pos}}^{\text{ini}}$. From the primal heuristic, we obtain a lower bound $\text{LB}^{\text{primal}}$ on $\lambda^k(A)$. Let

$$\lambda_{i_1} \geq \lambda_{i_2} \geq \cdots \geq \lambda_{i_n}$$

be the eigenvalues of $A$. If $\lambda_{i_{I_{\text{pos}}^{\text{ini}}}} < \text{LB}^{\text{primal}}$, then we set $\lambda \triangleq \lambda_{i_{I_{\text{pos}}^{\text{ini}}}}$. On the other hand, if $\lambda_{i_{I_{\text{pos}}^{\text{ini}}}} > \text{LB}^{\text{primal}}$, then let $l$ be the smallest index such that $\lambda_{i_l} > \text{LB}^{\text{primal}}$ and we set $\lambda \triangleq \lambda_{i_l}$.

2. Deciding $N$: In practice, $\theta_i$ was found to be significantly smaller than 1. So we used a value of $N = 3$ in all our experiments.

**4.3.2. Final details** A total time of 7200 seconds were given to each instance for running the convex IP (any extra time reported in the tables is due to running time of singular value decomposition and primal heuristics). We have run all our experiments with $k = 10, 20$. For the Convex-IP method, we use: $(I_{\text{pos}}^{\text{ini}}, N) = (10, 3)$. For the Pert-Convex-IP method, we let "iter" denote the maximum number of iterations. We used three settings in our experiments:

$$(I_{\text{pos}}^{\text{ini}}, N, \text{iter}) \in \{(5, 3, 10), \ (10, 3, 3), \ (15, 3, 2)\}.$$

The overall algorithms using the Pert-Convex-IP model and the Convex-IP model are presented in Appendix G.

## 4.4. Data Sets

We conduct numerical experiments on two types of data sets. Details of these two types of data sets are presented in Appendix H.

- **Artificial data set:** Tables 4, 5, 6, 7, 8, 9 present results for artificial/synthetic datasets.
- **Real data set:** Tables 10, 11, 12 show results for real data sets.

## 4.5. Description of the rows/columns in the tables

Note that the labels for each of the columns in Tables 4, 5, 6, 7, 8, 9, 10, 11, 12 are as follows:

- **Case:** The first part is a name. '**Case 1**' or '**Case 2**' denotes the instance number. The second part is the format (size, cardinality) which denotes the number of columns/rows of the $A$ matrix and the right-hand-side of the $\ell_0$ constraint of the original SPCA problem.
- **LB-$\ell_0$:** denotes the lower bound on the SPCA problem obtained from the (heuristic) Algorithm 2 in Section 4.2.
- **#-$\lambda$:** denotes the size of set $\{i \,|\, \lambda_i > \text{LB-}\ell_0\}$ where $\lambda_i$ are the eigenvalues of the covariance matrix.

- **Convex-IP-$\ell_0$, Pert-Convex-IP$_0$:** denote the Convex-IP and the Pert-Convex-IP models.
- **SDP:** denotes the semidefinite programming relaxation solved using Mosek. In Appendix J, we compare the dual bounds by alternative methods [16] to solve the SDP-relaxation for the real instances. Our conclusion based on our implementation of other algorithms is that when Mosek solves the instance, the best dual bound is obtained from Mosek. For some slightly larger instances, other algorithms might produce dual bounds. Usually, these dual bounds are extremely poor in quality. Moreover, these other methods do not scale up to instances with $d \geq 1000$. Therefore, we have chosen to present results only from Mosek in Tables 4, 5, 6, 7, 8, 9, 10, 11, 12; and the remaining results are relegated to Appendix J.
- **UB:** denotes the upper bound obtained from current dual bound method (i.e., Convex-IP-$\ell_0$, Pert-Convex-IP$_0$, SDP).
- **gap:** denotes the approximation ratio (duality gap) obtained by the formula $\mathbf{gap} = \frac{\text{UB}-\text{LB-}\ell_0}{\text{LB-}\ell_0}$.
- **time:** denotes the total running time—we present the overall running time due to singular value decomposition, heuristic method to obtain primal solutions, and solvers (Gurobi, Mosek) used to solve integer programming (set to terminate within 7200 seconds).

The three rows corresponding to Pert-Convex-IP, corresponds to experiments with three settings: $(I_{\text{pos}}, N, \text{iter}) = \{(5,3,10),\ (10,3,3),\ (15,3,2)\}$.

### 4.6. Conclusions and summary of numerical experiments

Based on numerical results reported in Tables 4, 5, 6, 7, 8, 9, 10, 11, 12 we draw some preliminary observations:

1. **Size of instances solved:**
- SDP: Because of limitation of hardware and software, the SDP relaxation method does not solve instances with input matrix of size greater than or equal to $300 \times 300$.
- Convex-IP: The convex IP shows better scalability than the SDP relaxation and produces dual bounds for instances with input matrix of size up to $500 \times 500$.
- Pert-Convex-IP: The perturbed convex IP scales significantly better that the other methods. While we experimented with instances up to size $2000 \times 2000$, we believe this method will easily scale to larger instances, when $k = 10, 20$ with $(I_{\text{pos}}, N)$ being chosen appropriately.

2. **Quality of dual bound:**
- SDP vs Best of {Convex-IP, Pert-Convex-IP}: While on some instances SDP obtained better dual bounds, this was not the case for all instances. For example, on the 'controlling sparsity' random instances and both the real data sets Eisen-1 and Eisen-2, SDP bounds are weaker.
- Convex-IP vs Pert-Convex-IP: If the convex IP solved within the time limit, then usually the bound is better than that obtained for Pert-Convex-IP. In other cases, Pert-Convex-IP performs better as it is easy to solve and usually solves within 1 hour.

- Overall gaps for Best of {Convex-IP, Pert-Convex-IP}: Except for the random instances of type 'controlling sparsity' of size $1000 \times 1000$, and Lymphoma data set, in all other instances at least one method had a gap less that 10%.

- Cardinality 10 vs Cardinality 20: When the cardinality budget is allowed to increase, based on our numerical results, we can see that the running time of our Convex-IP and Pert-Convex-IP methods do not change a lot, since the parameter of cardinality $k$ of Convex-IP and Pert-Convex-IP method only influences the linear constraint $\sum_{i=1}^{n} y_i \leq \sqrt{k}$, which is more robust to changes in the value of the cardinality $k$ than typical cardinality constraint in interger programming.

3. **Comparison of different numbers of splitting points (parameter $N$):** We compare the performances of the Pert-Convex-IP$_0$ method under distinct initialization splitting points with $(I_{\text{pos}}, N_{\text{ini}}, \# \text{ of iterations}) = (5,1,1), (5,3,1), (5,5,1)$, see Table 1. We present results with just one round of iterations to clearly understand the effect of number of splitting points. We observe that the gap decreases when the number of splitting points increases. On the other hand, the running time increases with the number of splitting points incereasing. However increasing splitting points from 3 to 5 does not significantly improve the bounds.

**Table 1**     Comparison of distinct splitting points

| Instance \ Splitting points | LB | (5,1,1) | | (5,3,1) | | (5,5,1) | |
|---|---|---|---|---|---|---|---|
| | | gap | Time | gap | Time | gap | Time |
| Eisen-1 (79, 10) | 17.335 | 2.619 % | 2.762 | 0.588 % | 3.049 | **0.329 %** | 3.127 |
| Eisen-2 (118, 10) | 11.718 | 13.245 % | 5.738 | 4.736 % | 7.194 | **4.207 %** | 7.78 |
| Colon (500, 10) | 2641.229 | 30.652 % | 72.802 | 27.755% | 73.149 | **27.673 %** | 76.115 |
| Lymphoma (500, 10) | 6008.741 | 52.412 % | 95.561 | 43.956 % | 83.902 | **43.587 %** | 86.422 |
| Reddit (2000, 10) | 1052.934 | 8.548 % | 1628.128 | 4.136 % | 1450.775 | **3.999 %** | 1488.936 |

4. **Comparison between $\ell_1$-relaxation and original sparsity constraint:** To further illustrate why we prescribe the use of $\ell_1$ relaxation to obtain dual bounds of SPCA, we compare the following two models: (1) The Pert-Convex-IP model used in the paper; (2) The same "perturbed convex IP" where the $\ell_1$ constraint is replaced by a cardinality constraint (with the introduction

of binary variables), denoted as Model-with-$\ell_0$.

$$
\max \quad \lambda + \sum_{i\in\{i:\lambda_i>\lambda\}}(\lambda_i-\lambda)\xi_i - s
$$

s.t.
$$
\begin{cases}
g_i = x^\top v_i \\
-\theta_i \le g_i \le \theta_i
\end{cases}
\qquad i \in \{i:\lambda_i>\lambda\}
$$

$$
\begin{cases}
g_i = \sum_{j=-N}^{N}\gamma_i^j\eta_i^j \\
\xi_i = \sum_{j=-N}^{N}(\gamma_i^j)^2\eta_i^j \\
(\eta_i^{-N},\ldots,\eta_i^N) \in \text{SOS-2}
\end{cases}
\qquad i \in \{i:\lambda_i>\lambda\} \qquad \text{(Model-with-}\ell_0\text{)}
$$

$$
\begin{cases}
\sum_{i=1}^{n} x_i^2 \le 1 \\
\sum_{i\in\{i:\lambda_i>\lambda\}} g_i^2 \le 1 - \frac{s}{\lambda-\overline{\lambda}} \\
1 - \frac{s}{\lambda-\overline{\lambda}} \le \sum_{i\in\{i:\lambda_i>\lambda\}}\xi_i \le 1 + \sum_{i\in\{i:\lambda_i>\lambda\}}\frac{\theta_i^2}{4N^2} - \frac{s}{\lambda-\overline{\lambda}}
\end{cases}
$$

$$
\begin{cases}
\sum_{i=1}^{n} z_i \le k \\
z_i \ge x_i, z_i \ge -x_i, z_i \in \{0,1\}, \forall i \in [n]
\end{cases}
\qquad (\ell_0 \text{ constraint})
$$

We tested on the real-life data for $k = 10$ and $k = 20$ in Table 2, Table 3. All parameters $(I_{\text{pos}}, N_{\text{ini}}, \#\text{iter})$ are also listed in Table 2, Table 3 which are the same as the parameters that used in the Section 4.3.2 (except for $\#\text{iter} = 1$ here).

Table 2     Comparison: Real Instances, cardinality parameter $k = 10$

| (size, index) | $(I_{\text{pos}}, N_{\text{ini}}, \#\text{ iter})$ | Pert-Convex-IP | | Model-with-$\ell_0$ | |
|---|---|---|---|---|---|
| | | Gap | Time | Gap | Time |
| Eisen Data 1 (79) | (5, 3, 1) | 0.588 % | 2.86 | **0.392 %** | 8.591 |
| | (10, 3, 1) | 0.796 % | 3.863 | 0.525 % | 99.168 |
| | (15, 3, 1) | 0.865 % | 10.049 | 0.588 % | 685.519 |
| Eisen Data 2 (118) | (5, 3, 1) | 4.736 % | 6.576 | 4.48 % | 86.251 |
| | (10, 3, 1) | 2.364 % | 27.525 | 2.321 % | 2105.51 |
| | (15, 3, 1) | 1.997 % | 195.356 | **1.971 %** | 5935.205 |
| Matrix CovColon (500) | (5, 3, 1) | 27.755 % | 90.362 | 4.48 % | 86.251 |
| | (10, 3, 1) | 2.364 % | 27.525 | **2.321 %** | 2105.51 |
| | (15, 3, 1) | 5.349 % | 2610.972 | 11.51 % | 7288.835 |
| Matrix LymphomaCov (500) | (5, 3, 1) | 43.956 % | 87.159 | 47.93 % | 7305.024 |
| | (10, 3, 1) | 23.662 % | 355.236 | 39.431 % | 7289.135 |
| | (15, 3, 1) | **17.863 %** | 4224.933 | 39.526 % | 7309.047 |
| Reddit (2000) | (5, 3, 1) | 4.136 % | 1867.157 | 5.826 % | 8765.165 |
| | (10, 3, 1) | **3.446 %** | 1831.221 | 8.867 % | 8638.037 |
| | (15, 3, 1) | 3.523 % | 3726.841 | 10.356 % | 8542.98 |

Based on the Table 2 3, following conclusions can be obtained:

(a) For instances with relative small size ($\le 500$): the upper bounds (UB) obtained from Model-with-$\ell_0$ is a slightly better than the upper bounds (UB) from Pert-Convex-IP, but the running time used for Model-with-$\ell_0$ is much longer than Pert-Convex-IP.

(b) For instances with relative large size ($\ge 500$): both the upper bounds and the running time obtained from Pert-Convex-IP method are significantly better than those obtained from Model-with-$\ell_0$. In another words, the Pert-Convex-IP is more scalable.

(c) Effect of $k$: We see that for $k = 20$ the performance of Pert-Convex-IP method is even more dramatically better than that of Model-with-$\ell_0$. In fact, now Pert-Convex-IP beats Model-with-$\ell_0$

**Table 3** Comparison: Real Instances, cardinality parameter $k = 20$

| (size, index) | $(I_{\mathbf{pos}}, N_{\mathbf{ini}}, \# \text{ iter})$ | Pert-Convex-IP | | Model-with-$\ell_0$ | |
|---|---|---|---|---|---|
| | | Gap | Time | Gap | Time |
| Eisen Data 1 (79) | (5, 3, 1) | **0.559 %** | 3.183 | 1.298 % | 7204.468 |
| | (10, 3, 1) | 0.813 % | 20.568 | 2.985 % | 7204.059 |
| | (15, 3, 1) | 0.886 % | 1016.839 | 5.519 % | 7229.677 |
| Eisen Data 2 (118) | (5, 3, 1) | 1.837 % | 6.48 | 2.65 % | 8062.349 |
| | (10, 3, 1) | 1.18 % | 46.001 | 4.223 % | 7211.949 |
| | (15, 3, 1) | **1.087 %** | 443.759 | 3.664 % | 7205.331 |
| Matrix CovColon (500) | (5, 3, 1) | 17.014 % | 75.267 | 18.539 % | 7268.644 |
| | (10, 3, 1) | 6.528 % | 372.802 | 12.903 % | 7271.37 |
| | (15, 3, 1) | **6.066 %** | 7275.58 | 12.737 % | 7273.013 |
| Matrix LymphomaCov (500) | (5, 3, 1) | 24.042 % | 91.786 | 26.622 % | 7288.825 |
| | (10, 3, 1) | 14.498 % | 214.784 | 24.381 % | 7302.236 |
| | (15, 3, 1) | **11.811 %** | 3349.161 | 35.286 % | 8831.009 |
| Reddit (2000) | (5, 3, 1) | **4.286 %** | 4652.869 | 7.139 % | 8708.004 |
| | (10, 3, 1) | 4.288 % | 1677.933 | 9.647 % | 8546.823 |
| | (15, 3, 1) | 4.776 % | 4274.327 | 12.157 % | 8560.558 |

on quality of bound and time even for small ($\leq 500$) instances. Indeed, this is another nice property

of the $\ell_1$-relaxation, namely it handles larger values of $k$ more robustly.

**Table 4**     Spiked Covariance Recovery - Cardinality 10

| Case | LB-$\ell_0$ | #-$\lambda$ | Convex-IP-$\ell_0$ | | Pert-Convex-IP$_0$ | | SDP | |
|---|---|---|---|---|---|---|---|---|
| | | | gap | Time | gap | Time | gap | Time |
| **Case 1 (200, 10)** | 511.95 | 1 | 0.005 % | 380 | 0.007 % | 76 | **0.001 %** | 1277 |
| | | | | | 0.005 % | 230 | | |
| | | | | | 0.005 % | 1605 | | |
| **Case 2 (200, 10)** | 592.45 | 1 | 0.003 % | 469 | 0.006 % | 615 | **0.002 %** | 1458 |
| | | | | | 0.006 % | 236 | | |
| | | | | | 0.005 % | 325 | | |
| **Case 1 (300, 10)** | 414.04 | 1 | **0.027 %** | 1692 | 0.03 % | 642 | NaN | - |
| | | | | | 0.029 % | 407 | | |
| | | | | | 0.027 % | 796 | | |
| **Case 2 (300, 10)** | 568.56 | 1 | **0.011 %** | 1067 | 0.016 % | 82 | NaN | - |
| | | | | | 0.014 % | 493 | | |
| | | | | | 0.012 % | 942 | | |
| **Case 1 (400, 10)** | 478.24 | 1 | **0.025 %** | 2598 | 0.04 % | 793 | NaN | - |
| | | | | | 0.03% | 610 | | |
| | | | | | 0.03% | 1495 | | |
| **Case 2 (400, 10)** | 426.91 | 1 | **0.037 %** | 3374 | 0.06 % | 181 | NaN | - |
| | | | | | 0.05 % | 846 | | |
| | | | | | 0.04 % | 2137 | | |
| **Case 1 (500, 10)** | 256.82 | 1 | **0.164 %** | 7525 | 0.21 % | 1345 | NaN | - |
| | | | | | 0.18 % | 1512 | | |
| | | | | | 0.17 % | 3279 | | |
| **Case 2 (500, 10)** | 551.74 | 1 | **0.029 %** | 7196 | 0.04 % | 152 | NaN | - |
| | | | | | 0.04 % | 725 | | |
| | | | | | 0.03 % | 1694 | | |
| **Case 1 (1000, 10)** | 315.16 | 1 | NaN | - | 0.57 % | 1147 | NaN | - |
| | | | | | 0.52 % | 776 | | |
| | | | | | **0.53 %** | 3633 | | |
| **Case 2 (1000, 10)** | 383.44 | 1 | NaN | - | 0.34 % | 2745 | NaN | - |
| | | | | | **0.32 %** | 403 | | |
| | | | | | 0.34 % | 3643 | | |

**Table 5** Spiked Covariance Recovery - Cardinality 20

| Case | LB-$\ell_0$ | #-$\lambda$ | Convex-IP-$\ell_0$ | | Pert-Convex-IP$_0$ | | SDP | |
|---|---|---|---|---|---|---|---|---|
| | | | gap | Time | gap | Time | gap | Time |
| **Case 1 (200, 20)** | 516.756 | 1 | 2.05 % | 493 | **0.008 %** | 746 | - % | - |
| | | | | | 0.073 % | 3116 | | |
| | | | | | 0.573 % | 7214 | | |
| **Case 2 (200, 20)** | 593.651 | 1 | 0.98 % | 1847 | **0.005 %** | 323 | -% | - |
| | | | | | 0.006 % | 5992 | | |
| | | | | | 0.102 % | 7215 | | |
| **Case 1 (300, 20)** | 499.92 | 1 | 0.70 % | 1848 | **0.018 %** | 745 | -% | - |
| | | | | | 0.021 % | 4799 | | |
| | | | | | 0.399 % | 7230 | | |
| **Case 2 (300, 20)** | 600.553 | 1 | 1.13 % | 1771 | 0.014 % | 530 | -% | - |
| | | | | | **0.013 %** | 2964 | | |
| | | | | | 0.272 % | 7232 | | |
| **Case 1 (400, 20)** | 483.995 | 1 | 2.74 % | 6398 | **0.034 %** | 1186 | -% - | |
| | | | | | 0.168 % | 7262 | | |
| | | | | | 0.832 % | 7255 | | |
| **Case 2 (400, 20)** | 428.275 | 1 | 1.92 % | 7426 | **0.045 %** | 576 | -% | - |
| | | | | | 0.074 % | 6965 | | |
| | | | | | 0.53 % | 7251 | - | |
| **Case 1 (500, 20)** | 294.35 | 1 | 1.19 % | 7027 | **0.162 %** | 1341 | -% | - |
| | | | | | 0.165 % | 6087 | | |
| | | | | | 1.285 % | 7294 | | |
| **Case 2 (500, 20)** | 571.15 | 1 | 1.96 % | 4628 | **0.039 %** | 1862 | - % | - |
| | | | | | 0.2 % | 1935 | | |
| | | | | | 1.215 % | 3360 | | |
| **Case 1 (1000, 20)** | 414 | 1 | - % | - | 0.53 % | 3133 | - % | - |
| | | | | | **0.50 %** | 2760 | | |
| | | | | | **0.50 %** | 5844 | | |
| **Case 2 (1000, 20)** | 391.795 | 1 | - % | - | **0.311 %** | 4756 | -% | - |
| | | | | | 0.74 % | 3596 | | |
| | | | | | 2.906 % | 7516 | | |

**Table 6** Synthetic Example - Cardinality 10

| Case | LB-$\ell_0$ | #-$\lambda$ | Convex-IP-$\ell_0$ | | Pert-Convex-IP$_0$ | | SDP | |
|---|---|---|---|---|---|---|---|---|
| | | | gap | Time | gap | Time | gap | Time |
| **Case 1 (200, 10)** | 5634.143 | 3 | 11.884 % | 7205 | 0.14 % | 38 | **0.10 %** | 1092 |
| | | | | | 0.15 % | 16 | | |
| | | | | | 0.15 % | 186 | | |
| **Case 2 (200, 10)** | 7321.23 | 3 | 1.703 % | 7205 | 0.13 % | 23 | **0.09 %** | 1086 |
| | | | | | 0.13 % | 13 | | |
| | | | | | 0.12 % | 47 | | |
| **Case 1 (300, 10)** | 4157.46 | 3 | 51.072 % | 7210 | **0.27 %** | 83 | NaN | - |
| | | | | | 0.29 % | 21 | | |
| | | | | | 0.27 % | 486 | | |
| **Case 2 (300, 10)** | 5135.50 | 3 | 65.275 % | 7210 | 0.23 % | 62 | NaN | - |
| | | | | | **0.22 %** | 59 | | |
| | | | | | 0.23 % | 58 | | |
| **Case 1 (400, 10)** | 6519.37 | 3 | 55.308 % | 7219 | **0.22 %** | 98 | NaN | - |
| | | | | | 0.23 % | 23 | | |
| | | | | | **0.22 %** | 349 | | |
| **Case 2 (400, 10)** | 5942.05 | 3 | 45.396 % | 7218 | **0.36 %** | 56 | NaN | - |
| | | | | | 0.42 % | 29 | | |
| | | | | | 0.41 % | 364 | | |
| **Case 1 (500, 10)** | 5125.86 | 3 | 65.98 % | 7230 | 0.38 % | 149 | NaN | - |
| | | | | | 0.38 % | 44 | | |
| | | | | | **0.37 %** | 132 | | |
| **Case 2 (500, 10)** | 5545.85 | 3 | 48.328 % | 7230 | 0.39 % | 50 | NaN | - |
| | | | | | **0.38 %** | 30 | | |
| | | | | | **0.38 %** | 231 | | |
| **Case 1 (1000, 10)** | 5116.08 | 3 | NaN | - | 0.58 % | 257 | NaN | - |
| | | | | | **0.57 %** | 128 | | |
| | | | | | **0.57 %** | 1373 | | |
| **Case 2 (1000, 10)** | 6946.12 | 3 | NaN | - | 0.39 % | 323 | NaN | - |
| | | | | | 0.36 % | 129 | | |
| | | | | | **0.34 %** | 1167 | | |

Dey et al.: *Dual bounds for sparse PCA*

| | | | Convex-IP-$\ell_0$ | | Pert-Convex-IP$_0$ | | SDP | |
|---|---|---|---|---|---|---|---|---|
| **Case** | **LB-$\ell_0$** | **#-$\lambda$** | gap | Time | gap | Time | gap | Time |
| **Case 1 (200, 20)** | 11222.152 | 2 | 0.779 % | 7205 | **0.041 %** | 2391 | -% | - |
| | | | | | 0.042 % | 2178 | | |
| | | | | | 0.466 % | 3707 | | |
| **Case 2 (200, 20)** | 14588.507 | 2 | 0.503 % | 7205 | **0.032 %** | 1285 | -% | - |
| | | | | | 0.036 % | 2772 | | |
| | | | | | 0.479 % | 7212 | | |
| **Case 1 (300, 20)** | 8282.32 | 3 | 13.336 % | 7212 | **0.089 %** | 2745 | - % | - |
| | | | | | 0.159 % | 1386 | | |
| | | | | | 1.523 % | 7227 | | |
| **Case 2 (300, 20)** | 10233.583 | 3 | 4.182 % | 7210 | 0.078 % | 1835 | -% | - |
| | | | | | **0.07 %** | 99 | | |
| | | | | | 0.817 % | 7229 | | |
| **Case 1 (400, 20)** | 12976.349 | 3 | 55.172 % | 7219 | **0.08 %** | 2563 | -% | - |
| | | | | | 0.105 % | 5278 | | |
| | | | | | 4.288 % | 7248 | | |
| **Case 2 (400, 20)** | 11809.325 | 2 | 45.209 % | 7219 | 0.082 % | 4257 | -% | - |
| | | | | | 0.084 % | 6934 | | |
| | | | | | **0.08 %** | 485 | | |
| **Case 1 (500, 20)** | 10218.591 | 3 | 65.637 % | 7231 | **0.13 %** | 3882 | -% | - |
| | | | | | 0.142 % | 6568 | | |
| | | | | | 2.067 % | 7288 | | |
| **Case 2 (500, 20)** | 11032.377 | 3 | 48.034 % | 7229 | **0.114 %** | 6603 | -% | - |
| | | | | | 0.138 % | 2753 | | |
| | | | | | 4.88 % | 7280 | | |
| **Case 1 (1000, 20)** | 10193.919 | 3 | - % | - | 1.38 % | 303 | -% | - |
| | | | | | 1.358 % | 1707 | | |
| | | | | | **0.24 %** | 3257 | | |
| **Case 2 (1000, 20)** | 13867.929 | 3 | - % | - | 0.691 % | 318 | -% | - |
| | | | | | 0.674 % | 1927 | | |
| | | | | | **0.18 %** | 8807 | | |

**Table 7**   Synthetic Example- Cardinality 20

**Table 8**    Controlling Sparsity - Cardinality 10

| Case | LB-$\ell_0$ | #-$\lambda$ | Convex-IP-$\ell_0$ | | Pert-Convex-IP$_0$ | | SDP | |
|---|---|---|---|---|---|---|---|---|
| | | | gap | Time | gap | Time | gap | Time |
| **Case 1 (200, 10)** | 706 | 1 | **0.14 %** | 925 | 2.9 %<br>2.6 %<br>2.6 % | 117<br>340<br>3663 | 0.42 % | 1360 |
| **Case 2 (200, 10)** | 680 | 1 | **0.14 %** | 1195 | 3.53 %<br>3.38 %<br>3.53 % | 176<br>372<br>3672 | 1.2 % | 1148 |
| **Case 1 (300, 10)** | 972 | 1 | **1.4 %** | 1958 | 3.91 %<br>3.81 %<br>3.70 % | 135<br>453<br>3635 | NaN | - |
| **Case 2 (300, 10)** | 976 | 1 | **1.1 %** | 3007 | 3.79 %<br>3.48 %<br>3.69 % | 278<br>1558<br>3772 | NaN | - |
| **Case 1 (400, 10)** | 1239 | 1 | **1.3 %** | 7207 | 4.21 %<br>3.96 %<br>3.96 % | 769<br>699<br>3699 | NaN | - |
| **Case 2 (400, 10)** | 1207 | 1 | **1.6 %** | 7206 | 3.56 %<br>3.48%<br>3.40 % | 221<br>1894<br>3697 | NaN | - |
| **Case 1 (500, 10)** | 1498 | 1 | **2.1 %** | 12180 | 5.21 %<br>4.74 %<br>4.81 % | 1026<br>2881<br>3661 | NaN | - |
| **Case 2 (500, 10)** | 1498 | 1 | **2.1 %** | 13917 | 4.14 %<br>4.07 %<br>4.01 % | 251<br>1039<br>3783 | NaN | - |
| **Case 1 (1000, 10)** | 3948 | 1 | - | - | 59.7 %<br>53.3 %<br>**49.5 %** | 2206<br>8318<br>3600 | NaN | - |
| **Case 2 (1000, 10)** | 4002 | 1 | NaN | - | 58.1 %<br>51.0 %<br>**47.6 %** | 3270<br>8356<br>3600 | NaN | - |

**Table 9**     Controlling Sparsity - Cardinality 20

| Case | LB-$\ell_0$ | #-$\lambda$ | Convex-IP-$\ell_0$ | | Pert-Convex-IP$_0$ | | SDP | |
|---|---|---|---|---|---|---|---|---|
| | | | gap | Time | gap | Time | gap | Time |
| **Case 1 (200, 20)** | 1341.432 | 1 | 0.97 % | 277 | 0.01 % | 1434 | -% | - |
| | | | | | **0.009 %** | 4726 | | |
| | | | | | 0.735 % | 2554 | | |
| **Case 2 (200, 20)** | 1287.45 | 1 | 1.63 % | 332 | 0.009 % | 887 | -% | - |
| | | | | | **0.008 %** | 2847 | | |
| | | | | | 1.22 % | 1971 | | |
| **Case 1 (300, 20)** | 1839.578 | 1 | 1.25 % | 1019 | **0.551 %** | 1932 | -% | - |
| | | | | | 0.636 % | 4854 | | |
| | | | | | 7.027 % | 7280 | | |
| **Case 2 (300, 20)** | 1849.485 | 1 | 0.192 % | 2217 | **0.19 %** | 897 | -% | - |
| | | | | | 0.796 % | 7229 | | |
| | | | | | 4.287 % | 7226 | | |
| **Case 1 (400, 20)** | 2339.441 | 1 | **1.45 %** | 907 | 2.140 % | 4343 | -% | - |
| | | | | | 5.47 % | 7265 | | |
| | | | | | 9.847 % | 7248 | | |
| **Case 2 (400, 20)** | 2273.785 | 1 | **2.34 %** | 3106 | 3.572 % | 3059 | -% | - |
| | | | | | 5.864 % | 5164 | | |
| | | | | | 10.537 % | 7249 | | |
| **Case 1 (500, 20)** | 2870.013 | 1 | **2.34 %** | 2773 | 3.376 % | 6013 | -% | - |
| | | | | | 4.077 % | 10870 | | |
| | | | | | 5.572 % | 7285 | | |
| **Case 2 (500, 20)** | 2832.149 | 1 | **2.37 %** | 3015 | 3.539 % | 5011 | -% | - |
| | | | | | 5.087 % | 7293 | | |
| | | | | | 5.063 % | 7283 | | |
| **Case 1 (1000, 20)** | 7535.996 | 1 | -% | - | 31.656 % | 7851 | -% | - |
| | | | | | 27.151 % | 721 | | |
| | | | | | **25.326 %** | 7518 | | |
| **Case 2 (1000, 20)** | 7759.88 | 1 | - % | - | 29.393 % | 311 | -% | - |
| | | | | | 25.230 % | 809 | | |
| | | | | | **23.433 %** | 7510 | | |

**Table 10**     First six sparse principal components of Pitprops

| Cardinality | LB-$\ell_0$ | Convex-IP-$\ell_0$ | | Pert-Convex-IP | | SDP | |
|---|---|---|---|---|---|---|---|
| | | gap | Time | gap | Time | gap | Time |
| **Cardinality 5** | 3.406 | 3.2 % | 0.40 | 6.0 % | 0.34 | **1.5 %** | 3.70 |
| **Cardinality 2** | 1.882 | 1.4 % | 0.23 | 3.6 % | 0.34 | **0 %** | 2.49 |
| **Cardinality 2** | 1.364 | 3.8 % | 0.30 | 7.6 % | 0.85 | **1.0 %** | 2.69 |
| **Cardinality 1** | 1 | 1.8 % | 0.75 | 3.5 % | 1.02 | **0 %** | 2.40 |
| **Cardinality 1** | 1 | 2.2 % | 0.30 | 3.6 % | 0.61 | **0 %** | 2.42 |
| **Cardinality 1** | 1 | 1.2 % | 0.30 | 2.1 % | 0.51 | **0 %** | 2.32 |
| **Sum of above** | 9.652 | 2.5 % | 2.28 | 4.8 % | 3.67 | **0.7 %** | 16.02 |

**Table 11**    Biological and Internet Data - Cardinality 10

| Case | LB-$\ell_0$ | #-$\lambda$ | Convex-IP-$\ell_0$ | | Pert-Convex-IP$_0$ | | SDP | |
|---|---|---|---|---|---|---|---|---|
| | | | gap | Time | gap | Time | gap | Time |
| **Eisen-1 (79, 10)** | 17.33 | 1 | 0.3 % | 4.6 | **0.12 %** | 63 | 2.2 % | 15 |
| | | | | | 0.17 % | 113 | | |
| | | | | | 0.4 % | 412 | | |
| **Eisen-2 (118, 10)** | 11.71 | 1 | **1.4 %** | 96 | 4.10 % | 69 | 2.0 % | 52 |
| | | | | | 2.13 % | 139 | | |
| | | | | | 1.70 % | 385 | | |
| **Colon (500, 10)** | 2641 | 1 | 14.7 % | 9000 | 27.7 % | 708 | NaN | - |
| | | | | | 9.58 % | 1181 | | |
| | | | | | **6.89 %** | 353 | | |
| **Lymphoma (500, 10)** | 6008 | 3 | 20.7 % | 3723 | 41 % | 610 | NaN | - |
| | | | | | 21 % | 1526 | | |
| | | | | | **17 %** | 2808 | | |
| **Reddit (2000, 10)** | 1052 | 1 | NaN | - | 3.59 % | 5663 | NaN | - |
| | | | | | **2.142 %** | 8584 | | |
| | | | | | 3.615 % | 4318 | | |

**Table 12**    Biological and Internet Data - Cardinality 20

| Case | LB-$\ell_0$ | #-$\lambda$ | Convex-IP-$\ell_0$ | | Pert-Convex-IP$_0$ | | SDP | |
|---|---|---|---|---|---|---|---|---|
| | | | gap | Time | gap | Time | gap | Time |
| **Eisen-1 (79, 20)** | 17.719 | 1 | 1.30 % | 742 | **0.062 %** | 450 | 2.37% | 13 |
| | | | | | 0.102 % | 7928 | | |
| | | | | | 0.333 % | 7205 | | |
| **Eisen-2 (118, 20)** | 19.323 | 1 | 2.02 % | 64 | 1.309 % | 283 | 2.28% | 53 |
| | | | | | **0.502 %** | 904 | | |
| | | | | | 1.294 % | 7206 | | |
| **Colon (500, 20)** | 4255.694 | 1 | 15.3 % | 7230 | 16.537 % | 4510 | - % | - |
| | | | | | **5.77 %** | 2931 | | |
| | | | | | 5.89 % | 7286 | | |
| **Lymphoma (500, 20)** | 9082.158 | 2 | 18.7 % | 7239 | 22.569 % | 1677 | - % | - |
| | | | | | 12.3 % | 1442 | | |
| | | | | | **11.81 %** | 3721 | | |
| **Reddit (2000, 20)** | 1119.046 | 1 | - % | - | **4.256 %** | 7920 | - % | - |
| | | | | | 4.288 % | 1677 | | |
| | | | | | 4.776 % | 4274 | | |

## 5. Acknowledgements

### References

[1] Genevera I Allen and Mirjana Maletić-Savatić. Sparse non-negative generalized PCA with applications to metabolomics. *Bioinformatics*, 27(21):3029–3035, 2011.

[2] Shrey Bagroy, Ponnurangam Kumaraguru, and Munmun De Choudhury. A social media based index of mental well-being in college campuses. In *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems*, CHI '17, pages 1634–1646, New York, NY, USA, 2017. ACM.

[3] Lauren Berk and Dimitris Bertsimas. Certifiably optimal sparse principal component analysis. *technical report*, 2016.

[4] Quentin Berthet and Philippe Rigollet. Optimal detection of sparse principal components in high dimension. *The Annals of Statistics*, 41(4):1780–1815, 2013.

[5] Daniel Bienstock. Computational study of a family of mixed-integer quadratic programming problems. *Mathematical Programming*, 74(2):121–140, 1996.

[6] Immanuel M Bomze and Gabriele Eichfelder. Copositivity detection by difference-of-convex decomposition and $\omega$-subdivision. *Mathematical Programming*, 138(1-2):365–400, 2013.

[7] Pierre Bonami, Oktay Günlük, and Jeff Linderoth. Solving box-constrained nonconvex quadratic programs. *Optimization online*, pages 26–76, 2016.

[8] Samuel Burer and Renato DC Monteiro. Local minima and convergence in low-rank semidefinite programming. *Mathematical Programming*, 103(3):427–444, 2005.

[9] Samuel Burer and Anureet Saxena. Old wine in a new bottle: The MILP road to MIQCP. *Optimization Online*, 2009.

[10] Samuel Burer and Dieter Vandenbussche. Globally solving box-constrained nonconvex quadratic programs with semidefinite-based finite branch-and-bound. *Computational Optimization and Applications*, 43(2):181–195, 2009.

[11] Jorge Cadima and Ian T Jolliffe. Loading and correlations in the interpretation of principle compenents. *Journal of Applied Statistics*, 22(2):203–214, 1995.

[12] Siu On Chan, Dimitris Papailiopoulos, and Aviad Rubinstein. On the worst-case approximability of sparse PCA. *arXiv preprint arXiv:1507.05950*, 2015.

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

[14] Alexandre d'Aspremont, Francis R Bach, and Laurent El Ghaoui. Full regularization path for sparse principal component analysis. In *Proceedings of the 24th international conference on Machine learning*, pages 177–184. ACM, 2007.

[15] Alexandre d'Aspremont, Laurent E Ghaoui, Michael I Jordan, and Gert R Lanckriet. A direct formulation for sparse PCA using semidefinite programming. In *Advances in neural information processing systems*, pages 41–48, 2005.

[16] Marianna De Santis, Franz Rendl, and Angelika Wiegele. Using a factored dual in augmented lagrangian methods for semidefinite programming. *Operations Research Letters*, 46(5):523–528, 2018.

[17] Alexandre d'Aspremont, Francis Bach, and Laurent El Ghaoui. Approximation bounds for sparse principal component analysis. *Mathematical Programming*, 148(1-2):89–110, 2014.

[18] Alexandre d'Aspremont, Francis Bach, and Laurent El Ghaoui. Optimal solutions for sparse principal component analysis. *Journal of Machine Learning Research*, 9(Jul):1269–1294, 2008.

[19] Antonio Frangioni and Claudio Gentile. SDP diagonalizations and perspective cuts for a class of nonseparable miqp. *Operations Research Letters*, 35(2):181–185, 2007.

[20] Trevor Hastie, Robert Tibshirani, and Martin Wainwright. *Statistical learning with sparsity*. CRC press, 2015.

[21] Yunlong He, Renato DC Monteiro, and Haesun Park. An algorithm for sparse PCA based on a new sparsity control criterion. In *Proceedings of the 2011 SIAM International Conference on Data Mining*, pages 771–782. SIAM, 2011.

[22] JNR Jeffers. Two case studies in the application of principal component analysis. *Applied Statistics*, pages 225–236, 1967.

[23] Ian T Jolliffe, Nickolay T Trendafilov, and Mudassir Uddin. A modified principal component technique based on the lasso. *Journal of computational and Graphical Statistics*, 12(3):531–547, 2003.

[24] Michel Journée, Yurii Nesterov, Peter Richtárik, and Rodolphe Sepulchre. Generalized power method for sparse principal component analysis. *Journal of Machine Learning Research*, 11(Feb):517–553, 2010.

[25] Jinhak Kim. *Cardinality Constrained Optimization Problems*. PhD thesis, Purdue University, West Lafayette, Indiana, 8 2016.

[26] Shiqian Ma. Alternating direction method of multipliers for sparse principal component analysis. *Journal of the Operations Research Society of China*, 1(2):253–274, Jun 2013.

[27] Malik Magdon-Ismail. NP-hardness and inapproximability of sparse PCA. *Information Processing Letters*, 126:35–38, 2017.

[28] Rahul Mazumder and Peter Radchenko. The discrete dantzig selector: Estimating sparse linear models via mixed integer linear optimization. *IEEE Transactions on Information Theory*, 63(5):3053–3075, 2017.

[29] George L Nemhauser and Laurence A Wolsey. *Integer and Combinatorial Optimization. Interscience Series in Discrete Mathematics and Optimization*. 1988.

[30] Dimitris Papailiopoulos, Alexandros Dimakis, and Stavros Korokythakis. Sparse PCA through low-rank approximations. In *International Conference on Machine Learning*, pages 747–755, 2013.

[31] James W Pennebaker, Martha E Francis, and Roger J Booth. Linguistic inquiry and word count: LIWC 2001. *Mahway: Lawrence Erlbaum Associates*, 71(2001):2001, 2001.

[32] Koustuv Saha and Munmun De Choudhury. Modeling stress with social media around incidents of gun violence on college campuses. *Proc. ACM Hum.-Comput. Interact.*, 1(CSCW):92:1–92:27, December 2017.

[33] Yla R Tausczik and James W Pennebaker. The psychological meaning of words: LIWC and computerized text analysis methods. *Journal of language and social psychology*, 29(1):24–54, 2010.

[34] Roman Vershynin. *High-Dimensional Probability An Introduction with Applications in Data Science*. Draft, 2016.

[35] Tengyao Wang, Quentin Berthet, and Richard J Samworth. Statistical and computational trade-offs in estimation of sparse principal components. *The Annals of Statistics*, 44(5):1896–1930, 2016.

[36] DM. Witten, R. Tibshirani, and T. Hastie. A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis. *Biostatistics*, 10(3):515–534, 2009.

[37] Xiao-Tong Yuan and Tong Zhang. Truncated power method for sparse eigenvalue problems. *arXiv preprint arXiv:1112.2679*, 2011.

[38] Xiao-Tong Yuan and Tong Zhang. Truncated power method for sparse eigenvalue problems. *Journal of Machine Learning Research*, 14(Apr):899–925, 2013.

[39] Youwei Zhang, Alexandre d'Aspremont, and Laurent El Ghaoui. Sparse PCA: Convex relaxations, algorithms and applications. In *Handbook on Semidefinite, Conic and Polynomial Optimization*, pages 915–940. Springer, 2012.

[40] Zhenyue Zhang, Hongyuan Zha, and Horst Simon. Low-rank approximations with sparse factors I: Basic algorithms and error analysis. *SIAM Journal on Matrix Analysis and Applications*, 23(3):706–727, 2002.

[41] Hui Zou, Trevor Hastie, and Robert Tibshirani. Sparse principal component analysis. *Journal of computational and graphical statistics*, 15(2):265–286, 2006.

## Appendix A: Notation

<p align="center">**Table 13**    Notation</p>

| Notation | Definition |
|---|---|
| $Y$ | data matrix of size $Y \in \mathbb{R}^{m \times n}$ |
| $A$ | sample covariance matrix $A = \frac{1}{m} Y^\top Y$ |
| $\|\cdot\|_0, \|\cdot\|_1, \|\cdot\|_2$ | $\ell_0, \ell_1, \ell_2$ norm |
| $k$ | sparsity parameter of sparse PCA problem |
| $\lambda^k(A)$ | optimal value of $\max_{\|x\|_0 \leq k, \|x\|_2 \leq 1} x^\top A x$ |
| $\mathrm{conv}(S)$ | convex hull of set $S$ |
| $[n]$ | short notation of index set $\{1, \ldots, n\}$ |
| $\mathrm{diag}(v)$ | diagonal matrix generated from a given vector $v$ |
| $\mathrm{tr}(A)$ | trace of a matrix $A$ |
| $\mathrm{OPT}_{\ell_1}$ | optimal value of $\max_{\|x\|_1 \leq \sqrt{k}, \|x\|_2 \leq 1} x^\top A x$ |
| $\rho$ | multiplicative approximation ratio between sparse PCA and its $\ell_1$ relaxation |
| $\{\lambda_i, v_i\}_{i=1}^n$ | eigenpair of covariance matrix $A$ |
| $\{g_i\}_{i=1}^n$ | continuous variable $g_i := x^\top v_i$ |
| $\{\theta_i\}_{i=1}^n$ | upper bound of $g_i$ defined as $\theta_i = \max\{x^\top v_i : \|x\|_2 \leq 1, \|x\|_0 \leq k\}$ |
| $\{\gamma_i^j\}_{j=-N}^N$ | splitting points of interval $[-\theta_i, \theta_i]$ for each $i$ |
| $\{\xi_i\}_{i=1}^n$ | piecewise linear upper approximation of $g_i^2$ |
| $s$ | upper bound of $\sum_{i \in \{i : \lambda_i < \lambda\}} -(\lambda_i - \lambda)(x^\top v_i)^2$ |
| $2N+1$ | number of splitting points for interval $[-\theta_i, \theta_i]$ for each $i \in \{i : \lambda_i > \lambda\}$ |
| $\bar{\lambda}$ | $\bar{\lambda} := \max\{\lambda_i : \lambda_i \leq \lambda\}$ |
| $\{\lambda_{i_j}\}_{j=1}^p$ | $\lambda_{i_1} \geq \cdots \geq \lambda_{i_p} \geq 0$ distinct values of eigenvalues of $A$ |
| $\Delta\lambda$ | eigenvalue gap $\Delta\lambda = \min\{\lambda_{i_j} - \lambda_{i_{j+1}}\}$ for $j = 1, \ldots, p-1$ |
| $\bar{A}$ | perturbed covariance matrix of $A$ |
| $(\bar{x}, \bar{y}, \bar{g}, \bar{\xi}, \bar{\eta}, \bar{s})$ | optimal solution for convex-IP |
| $\mathrm{OPT}_{\text{convex-IP}}$ | optimal value of convex integer programming model |
| $\mathrm{OPT}_{\text{pert-convex-IP}}$ | optimal value of perturbed convex integer programming model |
| $b_{(v)}$ | parameter used for cutting planes defined in Section 2.3.3 |
| $S_k$ | feasible region of sparse PCA with sparsity parameter $k$ |
| $T_k$ | $\ell_1$ relaxation of sparse PCA with sparsity parameter $k$ |
| $I_{\text{pos}}$ | the size of set $\{i : \lambda_i > \lambda\}$ |
| $I_{\text{pos}}^{\text{ini}}$ | initial input of $\{i : \lambda_i > \lambda\}$ |
| iter | number of iterations used for perturbed convex IP method |

## Appendix B: SDP relaxation

The SPCA problem $\max_{\|x\|_2 = 1, \|x\|_0 \leq k} x^\top A x$ is equivalent to a nonconvex problem:

$$\max \ \mathrm{tr}(AX)$$

$$\text{s.t. } \mathrm{tr}(X) = 1, \|X\|_0 \leq k^2, X \succeq 0, \mathrm{rank}(X) = 1.$$

Further relaxing this by replacing its rank and cardinality constraints with $\mathbf{1}^\top |X| \mathbf{1} \leq k$ gives the standard SDP relaxation:

$$\max \ \mathrm{tr}(AX)$$

$$\text{s.t. } \operatorname{tr}(X) = 1, \mathbf{1}^\top |X| \mathbf{1} \le k, X \succeq 0. \tag{SDP}$$

## Appendix C: Proof of Proposition 1

Proof of Proposition 1: Let $x^* = (x_i^*)_{i=1}^n$ be an optimal solution of SPCA. Then set

$$
\begin{cases}
g_i^* & \leftarrow (x^*)^\top v_i, & i \in [n], \\
((\eta_i^{-N})^*, \dots, (\eta_i^N)^*) & \leftarrow (\eta_i^{-N}, \dots, \eta_i^N) \in \text{SOS-2 and } \sum_{j=-N}^N \gamma_i^j (\eta_i^j)^* = g_i^*, & i \in \{i : \lambda_i > \lambda\}, \\
\xi_i^* & \leftarrow \sum_{j=-N}^N (\gamma_i^j)^2 \eta_i^j, & i \in \{i : \lambda_i > \lambda\}, \\
y_i^* & \leftarrow |x_i^*|, & i \in [n], \\
s_i^* & \leftarrow \sum_{i \in \{i : \lambda_i \le \lambda\}} -(\lambda_i - \lambda) g_i^*.
\end{cases}
$$

Note that the above solution $(x^*, y^*, g^*, \xi^*, \eta^*, s^*)$ is a feasible solution for Convex-IP. This is easy to verify for all the constraints except the constraint $\sum_{i \in \{i : \lambda_i > \lambda\}} \xi_i + \sum_{i \in \{i : \lambda_i \le \lambda\}} g_i^2 \le 1 + \frac{1}{4N^2} \sum_{i \in \{i : \lambda_i > \lambda\}} \theta_i^2$. Note that to verify this constraint, it is sufficient to verify that $\xi_i \le g_i^2 + \frac{1}{4N^2} \theta_i^2$ for $i \in \{i : \lambda_i > \lambda\}$. This is easily verified based on the size of the discretization and the structure of SOS-2 constraints.

Moreover, the objective value of feasible solution $(x^*, y^*, g^*, \xi^*, \eta^*, s^*)$ is

$$
\begin{aligned}
\lambda + \sum_{i \in \{i : \lambda_i > \lambda\}} (\lambda_i - \lambda) \xi_i^* - s^* &\geq \lambda + \sum_{i \in \{i : \lambda_i > \lambda\}} (\lambda_i - \lambda)(g_i^*)^2 - s^* \\
&= \lambda + \sum_{i \in \{i : \lambda_i > \lambda\}} (\lambda_i - \lambda)((x^*)^\top v_i)^2 + \sum_{i \in \{i : \lambda_i \le \lambda\}} (\lambda_i - \lambda)((x^*)^\top v_i)^2 \\
&= \lambda + \sum_{i=1}^n (\lambda_i - \lambda)((x^*)^\top v_i)^2.
\end{aligned}
$$

Note that the optimal solution $x^*$ of SPCA has property $\|x^*\|_2 = 1$ and $\sum_{i=1}^n v_i v_i^\top = I_n$. Then $\lambda + \sum_{i=1}^n (\lambda_i - \lambda)((x^*)^\top v_i)^2 = (x^*)^\top A x^* = \lambda^k(A)$. Therefore, $\text{OPT}_{\text{convex-IP}} \ge \lambda^k(A)$.

## Appendix D: Proof of Proposition 2

Proof of Proposition 2: Let $(\bar{x}, \bar{y}, \bar{g}, \bar{\xi}, \bar{\eta}, \bar{s})$ be an optimal solution for Convex-IP. Its optimal value then satisfies the following:

$$
\begin{aligned}
\text{OPT}_{\text{convex-IP}} &= \lambda + \sum_{i \in \{i : \lambda_i > \lambda\}} (\lambda_i - \lambda) \bar{\xi}_i - \bar{s} \\
&= \lambda + \sum_{i \in \{i : \lambda_i > \lambda\}} (\lambda_i - \lambda) \left( \bar{\xi}_i - \bar{g}_i^2 + \bar{g}_i^2 \right) - \bar{s} \\
&= \lambda + \sum_{i \in \{i : \lambda_i > \lambda\}} (\lambda_i - \lambda) \left( \bar{\xi}_i - \bar{g}_i^2 \right) + \sum_{i \in \{i : \lambda_i > \lambda\}} (\lambda_i - \lambda) \bar{g}_i^2 - \bar{s}.
\end{aligned}
$$

Since variable $s$ satisfies $\sum_{i \in \{i : \lambda_i \le \lambda\}} -(\lambda_i - \lambda) g_i^2 \le s$, to maximize the objective function, $\bar{s}$ should be equivalent to $\sum_{i \in \{i : \lambda_i \le \lambda\}} -(\lambda_i - \lambda) \bar{g}_i^2$, then the above formula can be represented as

$$
\lambda + \sum_{i \in \{i : \lambda_i > \lambda\}} (\lambda_i - \lambda) \left( \bar{\xi}_i - \bar{g}_i^2 \right) + \sum_{i \in \{i : \lambda_i > \lambda\}} (\lambda_i - \lambda) \bar{g}_i^2 - \bar{s}
$$

$$\begin{aligned}
=\lambda + & \sum_{i\in\{i:\lambda_i>\lambda\}}(\lambda_i-\lambda)\left(\bar{\xi}_i-\bar{g}_i^2\right) + \sum_{i\in\{i:\lambda_i>\lambda\}}(\lambda_i-\lambda)\bar{g}_i^2 + \sum_{i\in\{i:\lambda_i\leq\lambda\}}(\lambda-\lambda)\bar{g}_i^2 \\
= & \sum_{i\in\{i:\lambda_i>\lambda\}}(\lambda_i-\lambda)\left(\bar{\xi}_i-\bar{g}_i^2\right) + \left(\lambda+\sum_{i=1}^{n}(\lambda_i-\lambda)\bar{g}_i^2\right).
\end{aligned} \tag{18}$$

By previous results, $\lambda+\sum_{i=1}^{n}(\lambda_i-\lambda)\bar{g}_i^2 = \bar{x}^{\top}A\bar{x}$. Note that due to the $\ell_2$−norm constraint $\|x\|_2 \leq 1$ and the $\ell_1$−norm constraint present in Convex-IP problem, we have $\bar{x} \in T_k = \{x \in \mathbb{R}^n : \|x\|_2 \leq 1, \|x\|_1 \leq \sqrt{k}\} \subseteq \rho \cdot \mathrm{Conv}(S_k)$. Therefore $\bar{x}^{\top}A\bar{x}$ is upper bounded by the value $\rho^2 \cdot \lambda^k(A)$.

To upper bound the first term in (18), since $g_i = \sum_{j=-N}^{N}\gamma_i^j\eta_i^j$, $\xi_i = \sum_{j=-N}^{N}(\gamma_i^j)^2\eta_i^j$ for $i \in \{i : \lambda_i > \lambda\}$ and the SOS-2 construction enforces that there are at most two *active* continuous SOS-2 variables $\eta_i^j, \eta_i^{j+1}$ such that $\eta_i^j + \eta_i^{j+1} = 1$ with $\eta_i^j, \eta_i^{j+1} \geq 0$ and the other SOS-2 variables are all zeros, then

$$\begin{aligned}
\xi_i - g_i^2 & = \sum_{j=-N}^{N}(\gamma_i^j)^2\eta_i^j - \left(\sum_{j=-N}^{N}\gamma_i^j\eta_i^j\right)^2 \\
& = (\gamma_i^j)^2\eta_i^j + (\gamma_i^{j+1})^2\eta_i^{j+1} - \left(\gamma_i^j\eta_i^j + \gamma_i^{j+1}\eta_i^{j+1}\right)^2 & \text{for } \eta_i^j, \eta_i^{j+1} \text{ active} \\
& = (\gamma_i^{j+1}-\gamma_i^j)^2\eta_i^j(1-\eta_i^j) & \text{via } \eta_i^j + \eta_i^{j+1} = 1 \\
& \leq \max_{j=-N,\dots,N-1}(\gamma_i^{j+1}-\gamma_i^j)^2 \cdot \frac{1}{4}
\end{aligned}$$

where in all possible partition of $[-\theta_i, \theta_i]$, the evenly partition of $[-\theta_i, \theta_i]$ achieves the minimum value of $\max_{j=-N,\dots,N-1}(\gamma_i^{j+1}-\gamma_i^j)^2 = \frac{\theta_i^2}{N^2}$. Hence (18) can be upper bounded as follows:

$$\begin{aligned}
\mathrm{OPT}_{\text{convex-IP}} & = \sum_{i\in\{i:\lambda_i>\lambda\}}(\lambda_i-\lambda)\left(\bar{\xi}_i-\bar{g}_i^2\right) + \left(\lambda+\sum_{i=1}^{n}(\lambda_i-\lambda)\bar{g}_i^2\right) \\
& \leq \frac{1}{4N^2}\sum_{i\in\{i:\lambda_i>\lambda\}}(\lambda_i-\lambda)\theta_i^2 + \rho^2 \cdot \lambda^k(A).
\end{aligned}$$

## Appendix E: Appendix: Proof of Proposition 3

Proof of Proposition 3: Given the heuristic lower bound $\lambda$, the number of splitting points $N$, the size of set $I_{\text{pos}} = |\{i : \lambda_i > \lambda\}|$, for each $i \in \{i : \lambda_i > \lambda\}$, there are at most $2N$ possible choices of *active* SOS-2 variables, i.e.,

$$\eta_i^j, \eta_i^{j+1} > 0, \text{ for } j = -N,\dots,0,\dots,N-1.$$

Thus there are at most $(2N)^{|I_{\text{pos}}|}$ choices of *active* SOS-2 variables for a Convex-IP problem. For a fixed value of *active* SOS-2 variables, the Convex-IP problem reduces to be a continuous convex optimization problem which can be solved exactly within polynomial time, say $T$. Thus the Convex-IP can be solved within $(2N)^{|I_{\text{pos}}|} \cdot T$.

## Appendix F:  Proof of Proposition 4

Proof of Proposition 4: Based on Proposition 2, we have

$$\text{OPT}_{\text{Pert-Convex-IP}} \leq \rho^2 \lambda^k(\bar{A}) + \frac{1}{4N^2} \sum_{i \in \{i : \lambda_i > \lambda\}} (\lambda_i - \lambda)\theta_i^2.$$

Note that $\bar{A} - A = \sum_{i \in \{i : \lambda_i \leq \lambda\}} (\bar{\lambda} - \lambda_i) v_i v_i^\top$. Therefore,

$$\begin{aligned}
\rho^2 \lambda^k(\bar{A}) &= \rho^2 \lambda^k \left( A + (\bar{A} - A) \right) \\
&\leq \rho^2 \lambda^k(A) + \rho^2 \lambda^k(\bar{A} - A) \\
&\leq \rho^2 \lambda^k(A) + \rho^2(\bar{\lambda} - \lambda_{\min}(A)).
\end{aligned}$$

## Appendix G:  Convex-IP Method and Pert-Convex-IP Method

Algorithm 3 presents all the details of the convex IP solved. Algorithm 4 presents all the details of the Pert-Convex-IP solved.

---

**Algorithm 3** Convex-IP Method

---

1: *Input*: Sample covariance matrix $A$, cardinality constraint $k$, size of set $\{i : \lambda_i > \lambda\}$ we desire, number of one branch splitting points $N$.

2: *Output*: Lower and upper bound of SPCA or $\ell_1$-relax based on the choice of $\theta_i$.

3: **function** CONVEX-IP METHOD$(A, k, I_{\text{pos}}, N)$

4:     Set lower bound and warm starting point $(\text{LB}, \bar{x}) \leftarrow$ HEURISTIC METHOD$(A, k, x^0)$.

5:     Set parameter $\lambda_{I_{\text{pos}+1}} \leq \lambda \leq \text{LB}$ if possible, otherwise set $\lambda \leftarrow \text{LB}$.

6:     Set splitting points $\gamma_i^j$ as above based on $N$ and the choice of $\theta_i$, see Section 2.2 [2.2] .

7:     To warm start, add additional splitting points based on the point $\bar{x}$.

8:     Add cutting-plane (9) to the model based on the choice of $\theta_i$.

9:     Run Convex-IP problem.

10:     Set UB $\leftarrow$ Convex-IP if running to the optimal, or the current dual bound obtained from Convex-IP.

11:     **return** LB, UB.

12: **end function**

---

## Appendix H:  Description of Data Sets

### H.1.  Artificial Data Sets

We first conduct numerical experiments on three types of artificial data sets, denoted as the spiked covariance recovery from the paper [30], the synthetic example from the paper [41], and the controlling sparsity case from the paper [15]. A description of each of these three types of instances is presented below:

---

**Algorithm 4** Pert-Convex-IP Method

1:  *Input*: Sample covariance matrix $A$, cardinality constraint $k$, size of set $\{i : \lambda_i > \lambda\}$ we desire, number of one branch splitting points $N$, maximum number of iterations iter.

2:  *Output*: Lower and upper bound of SPCA or $\ell_1$-relax based on the choice of $\theta_i$.

3:  **function** PERT-CONVEX-IP METHOD$(A, k, I_{\text{pos}}, N, \text{iter})$

4:      Set lower bound and warm starting point $(\text{LB}, \bar{x}) \leftarrow$ HEURISTIC METHOD$(A, k, x^0)$.

5:      Set parameter $\lambda_{I_{\text{pos}+1}} \leq \lambda \leq \text{LB}$ if possible, otherwise set $\lambda \leftarrow \text{LB}$.

6:      Set parameter $\bar{\lambda} \triangleq \max\{\lambda_i : \lambda_i \leq \lambda\} < \lambda$ if possible.

7:      Set splitting points $\gamma_i^j$ as above based on $N$ and the choice of $\theta_i$, see Section 2.2 [2.2].

8:      To warm start, add additional splitting points based on the point $\bar{x}$.

9:      **while** current iteration does not exceed the maximum number of iterations iter or time limit is not up **do**

10:          Run Pert-Convex-IP problem.

11:          Set UB $\leftarrow$ Pert-Convex-IP if running to the optimal, or the current dual bound obtained from Pert-Convex-IP.

12:          Set $\hat{x} \leftarrow$ current feasible solution obtained from Pert-Convex-IP.

13:          Add additional splitting points based on solution obtained in solving Pert-Convex-IP problem.

14:          Add cutting-plane (9) to the model based on the choice of $\theta_i$.

15:      **end while**

16:      **return** LB, UB.

17: **end function**

---

**H.1.1. Spiked covariance recovery**   Consider a covariance matrix $\Sigma$, which has two sparse eigenvectors with dominated eigenvalues and the rest eigenvector are unconstrained with small eigenvalues. Let the first two dominant eigenvectors $v_1, v_2$ of $\Sigma$ be:

$$[v_1]_i = \begin{cases} \frac{1}{\sqrt{10}} & i = 1, \dots, 10, \\ 0 & \text{otherwise} \end{cases}, \qquad [v_2]_i = \begin{cases} \frac{1}{\sqrt{10}} & i = 11, \dots, 20, \\ 0 & \text{otherwise} \end{cases}, \qquad (19)$$

with the eigenvalues corresponding to the first two dominant eigenvectors be $\lambda_1 \gg 1$ and $\lambda_2 \gg 1$, and the remaining eigenvalues be 1. For example, in our numerical experiments, set $\Sigma \leftarrow 399 \cdot v_1 v_1^\top + 299 \cdot v_2 v_2^\top + I$.

We have four distinct settings under the spiked covariance recovery case. Let $n$ be the number of features, i.e., the size of the sample covariance matrix of our numerical cases. Let $m$ be the number of samples we generated. We set $n = \{200, 300, 400, 500, 1000\}$ and $m = \{50\}$. Therefore, under each setting of $n$, we generate $m$ random samples $x_i \sim N(0, \Sigma)$, and get our sample covariance matrix

$\hat{\Sigma} = \frac{1}{50} \sum_{i=1}^{50} x_i x_i^\top$. In Table 4, for each setting, we repeat the experiment for 2 times (case 1, case 2), and compare the dual bounds obtained from all three methods.

**H.1.2. Synthetic Example** Given $n$, let $n_1, n_2, n_3 \in \left\{ \lceil \frac{n}{3} \rceil, \lfloor \frac{n}{3} \rfloor \right\}$ such that $n_1 + n_2 + n_3 = n$. Let $\mathbf{0}_{p \times q}$ be the matrix of all zeros with size $p \times q$. Let $\mathbf{1}_p$ be the vector of all ones with length $p$.

Then:

$$\Sigma = \begin{pmatrix} 290 \cdot \mathbf{1}_{n_1} \mathbf{1}_{n_1}^\top + I_{n_1} & \mathbf{0}_{n_1 \times n_2} & -87 \cdot \mathbf{1}_{n_1} \mathbf{1}_{n_3}^\top \\ \mathbf{0}_{n_2 \times n_1} & 300 \cdot \mathbf{1}_{n_2} \mathbf{1}_{n_2}^\top + I_{n_2} & 277.5 \cdot \mathbf{1}_{n_2} \mathbf{1}_{n_3}^\top \\ -87 \cdot \mathbf{1}_{n_3} \mathbf{1}_{n_1}^\top & 277.5 \cdot \mathbf{1}_{n_3} \mathbf{1}_{n_2}^\top & 582.7875 \cdot \mathbf{1}_{n_3} \mathbf{1}_{n_3}^\top + I_{n_3} \end{pmatrix}. \tag{20}$$

In our experiments, we set $n = \{200, 300, 400, 500, 1000\}$, and generate $m = 50$ samples such that $x_i \sim N(0, \Sigma)$. Again, the sample empirical covariance matrix is $\hat{\Sigma} = \frac{1}{50} \sum_{i=1}^{50} x_i x_i^\top$. In Table 6, for each setting of $n$, we repeat the experiment twice (case 1, case 2), and compare dual bounds obtained from all three methods.

**H.1.3. Controlling Sparsity** Like the spiked covariance recovery case, the covariance matrix $\Sigma$ of controlling sparsity case can also be represented as the summation of a term generated by sparse eigenvector with dominated eigenvalue and the remaining part with small eigenvalues. Generate a $n \times n$ matrix $U$ with uniformly distributed coefficients in $[0, 1]$ which can be seen as white noise. Let $v \in \{0, 1\}^n$ be a sparse vector with $\|v\|_0 \leq k$. We then form a test matrix $\Sigma = U^\top U + \sigma v v^\top$, where $\sigma$ is the signal-to-noise ratio and is set to 15.

In our experiments, we set $n = \{200, 300, 400, 500, 1000\}$ and generate $m = 50$ samples $x_i \sim N(0, \Sigma)$ for $i = 1, \ldots, 50$. Therefore the sample empirical covariance matrix is $\hat{\Sigma} = \frac{1}{50} \sum_{i=1}^{50} x_i x_i^\top$. In Table 8, for each setting of $n$, we repeat the experiment twice (case 1, case 2), and compare dual bounds obtained from all three methods.

## H.2. Real Data Sets

We conduct numerical experiments on three types of real data sets, the benchmark pitprops data from [22], biological data from [14, 30, 38] and large-scale data collected from internet.

**H.2.1. Pitprops Data** The PitProps data set in [22] (consisting of 180 observations with 13 measured variables) has been a standard benchmark to evaluate algorithms for sparse PCA.

Based on previous work, we also consider the first six $k-$sparse principal components. Note the $i$-th $k-$sparse principal component $x^i$ is obtained by solving $\arg\max_{\|x\|_2 = 1, \|x\|_0 \leq k} x^\top A^i x$ where $A^1 \leftarrow A$ and $A^i \leftarrow (I - x^{i-1}(x^{i-1})^\top) A^{i-1} (I - x^{i-1}(x^{i-1})^\top)$ for $i = 2, \ldots, 6$. Table 10 lists the six extracted sparse principal direction with cardinality setting $5 - 2 - 2 - 1 - 1 - 1$.

**H.2.2. Biological Data** In Table 11 we present numerical experiments on four biological data sets. The first two biological data sets (Eisen-1, Eisen-2) are from [38]. The Colon cancer data set is from Alon et al. (1999). The Lymphoma data set is from Alizadeh et al. (2000).

**H.2.3. Large-scale Internet Data**   In Table 11 we also present numerical experiments on internet dataset. This dataset is constructed out of textual posts shared on the popular social media Reddit. Based on prior work [2, 32], the archive of all public Reddit posts shared on Google's Big Query was utilized to obtain a set of 3292 posts from the subreddit r/stress from December 2010 to January 2017. The r/stress community allows individuals to self-report and disclose their stressful experiences and is a support community. For example, two (paraphrased) post excerpts say: "Feel like I am burning out (again...) Help: what do I do?"; and "How do I calm down when I get triggered?". The community is also heavily moderated; hence these 3292 posts were considered to be indicative of actual stress. [32].

Then on this collected set of posts, standard text-based feature extraction techniques were applied per post, starting with cleaning the data (stopword elimination, removal of noisy words, stemming), and then building a language model with the n-grams in a post ($n=2$). The outcomes of this language model provided us with 1950 features, after including only the top most statistically significant features. Additionally, the psycholinguistic lexicon Linguistic Inquiry and Word Count (LIWC) [31] was leveraged to obtain features aligning with 50 different empirically validated psychological categories, such as positive affect, negative affect, cognition, and function words. These features have been extensively validated in prior work to be indicative of stress and similar psychological constructs [33]. Our final dataset matrix comprised 3092 rows, corresponding to the 3092 posts, and 2000 features in all.

The purpose of testing the sparse PCA technique on this dataset is to identify those features that are theoretically guaranteed to be the most salient in describing the nature of stress expressed in a post. In turn, these salient features could be utilized by a variety of stakeholders like clinical psychologists, and community moderators and managers to gain insights into stress-related phenomenon as well as to direct interventions as appropriate.

The final $A$ matrix can be found on the website:

https://www2.isye.gatech.edu/ sdey30/publications.html

## Appendix I: Comparison with Existing Primal Heuristics for Lower Bounds

In this section, we compare our method Algorithm 2 for obtaining good primal feasible solutions with two standard heuristics methods for sparse PCA in the literature: truncated power method (TPM, [37]), generalized power method (GPM, [24]) with $\ell_0$-penalty. See Table 14 for a comparison on all the real instances. As we can see, all the methods produce solutions with more or less the same objective function values.

**Table 14**     Compare with existing primal methods

| INSTANCE | SPCA-PRIMAL (OUR METHOD) | | TPM | | GPM | |
|---|---|---|---|---|---|---|
| | LB | TIME | LB | TIME | LB | TIME |
| PITPROPS $k=5$ | **3.406** | 0.1 | **3.406** | 0.0 | **3.406** | 0.1 |
| EISEN-1 $k=10$ | **17.335** | 0.0 | **17.335** | 0.0 | **17.335** | 2.3 |
| EISEN-2 $k=10$ | **11.718** | 0.0 | **11.718** | 0.0 | 11.605 | 4.1 |
| COVCOLON $k=10$ | **2641.228** | 0.4 | **2641.228** | 0.4 | **2641.228** | 59.7 |
| LYMP $k=10$ | **5911.412** | 0.3 | **5911.412** | 0.2 | 5753.563 | 81.4 |
| REDDIT $k=10$ | **1052.020** | 7.4 | **1052.020** | 4.5 | **1052.020** | 1881.4 |

## Appendix J: Comparison with Existing Methods for Dual Bounds

In this section, we compare the performance of our convex integer program method with (1) Mosek, in our experience one of the best commercial implementations of SDP solvers; and (2) two variants of the approach presented in [16], which uses the main idea of [8]. The variants are listed as follows:

1. **DADAL:** Directly using code available online from [16]: Dual Alternating Direction Augmented Lagrangian (DADAL) method can be used to find out the upper bounds of the SDP problem. In order to use the freely available implementation, the DADAL method requires the remodeling of the original problem into the following standard format:

$$\min \langle \boldsymbol{A}, \boldsymbol{X} \rangle \text{ s.t } \mathcal{A}(\boldsymbol{X}) = \boldsymbol{b}, \ \boldsymbol{X} \succeq \boldsymbol{0}.$$

Thus to find the dual bounds of the sparse PCA with covariance matrix of size $d$, we need to (1) add additional auxiliary variables for inequality constraints, (2) reformulate the variables into a p.s.d. matrix. For the step-(1), the original sparse PCA problem can be formulated in the following fashion:

$$\min \ \langle -\boldsymbol{A}, \boldsymbol{X} \rangle \qquad\qquad\qquad\qquad \text{(SDP-equality)}$$
$$\text{s.t.} \ \langle \boldsymbol{I}_d, \boldsymbol{X} \rangle + \mu_1 = 1$$
$$\langle \boldsymbol{I}_{d^2}, \text{diag}(\boldsymbol{Y}) \rangle + \mu_2 = k$$
$$\langle \boldsymbol{E}_{ij}^+, \boldsymbol{X} \oplus \text{diag}(\boldsymbol{Y}) \rangle + \gamma_{ij}^+ = 0, \ \forall \ ij$$
$$\langle \boldsymbol{E}_{ij}^-, \boldsymbol{X} \oplus \text{diag}(\boldsymbol{Y}) \rangle + \gamma_{ij}^- = 0, \ \forall \ ij$$
$$\boldsymbol{X}, \text{diag}(\boldsymbol{Y}), \text{diag}(\boldsymbol{\gamma}^+), \text{diag}(\boldsymbol{\gamma}^-), \text{diag}(\mu) \succeq \boldsymbol{0}$$

where $\oplus$ is the direct sum of two matrices, i.e., $\boldsymbol{A} \oplus \boldsymbol{B} := \begin{pmatrix} \boldsymbol{A} & 0 \\ 0 & \boldsymbol{B} \end{pmatrix}$, the matrix $\text{diag}(\boldsymbol{Y})$ is a short notation of $\text{diag}(\text{vec}(\boldsymbol{Y}))$ with $\text{vec}(\boldsymbol{Y})$ the vectorization of matrix $\boldsymbol{Y}$, and the matrix $\boldsymbol{E}_{ij}^+, \boldsymbol{E}_{ij}^-$ are

$$\boldsymbol{E}_{ij}^+ := \begin{pmatrix} E_{ij} & 0 \\ 0 & -\text{diag}(\text{vec}(E_{ij})) \end{pmatrix}, \ \boldsymbol{E}_{ij}^- := \begin{pmatrix} -E_{ij} & 0 \\ 0 & -\text{diag}(\text{vec}(E_{ij})) \end{pmatrix}, \qquad \forall i, j \in [d] \times [d]$$

with $E_{ij} \in \mathbb{R}^{d \times d}$ the standard basis matrix (i.e., the component $(i, j)$ equals to 1, and the rest components equal to 0). Rewrite the variables of SDP-equality into a p.s.d. matrix

$$\tilde{X} := \begin{pmatrix} X & & & & \\ & \mathrm{diag}(Y) & & & \\ & & \mathrm{diag}(\gamma^+) & & \\ & & & \mathrm{diag}(\gamma^-) & \\ & & & & \mu \end{pmatrix} \in \mathbb{R}^{(d+3d^2+2) \times (d+3d^2+2)}.$$

For the step-(2), the SDP-equality can be further transferred into the standard SDP format as follows:

$$\min \ \langle -A \oplus \mathbf{0}_{d^2} \oplus \mathbf{0}_{d^2} \oplus \mathbf{0}_{d^2} \oplus \mathbf{0}_2, \tilde{X} \rangle \qquad \text{(standard-SDP)}$$

$$\text{s.t.} \ \langle I_d \oplus \mathbf{0}_{d^2} \oplus \mathbf{0}_{d^2} \oplus \mathbf{0}_{d^2} \oplus \mathrm{diag}(1, 0), \tilde{X} \rangle = 1$$

$$\langle \mathbf{0}_d \oplus I_{d^2} \oplus \mathbf{0}_{d^2} \oplus \mathbf{0}_{d^2} \oplus \mathrm{diag}(0, 1), \tilde{X} \rangle = k$$

$$\langle (E_{ij}^+ + E_{ij}^+) \oplus (\mathrm{diag}(\mathrm{vec}(E_{ij})) + \mathrm{diag}(\mathrm{vec}(E_{ji}))) \oplus \mathbf{0}_{d^2} \oplus \mathbf{0}_2, \tilde{X} \rangle = 0, \ \forall i \geq j$$

$$\langle (E_{ij}^- + E_{ij}^-) \oplus \mathbf{0}_{d^2} \oplus (\mathrm{diag}(\mathrm{vec}(E_{ij})) + \mathrm{diag}(\mathrm{vec}(E_{ji}))) \oplus \mathbf{0}_2, \tilde{X} \rangle = 0, \ \forall i \geq j$$

$$\tilde{X} \succeq \mathbf{0}$$

with the size of variable matrix $n = d + 3d^2 + 2$ and the number of linear constraints $m = 2 + d \times (d + 1)$. The code of DADAL method is downloaded from the author's [16] homepage [3].

2. **DADAL-SPCA:** A DADAL-SPCA method designed by us (which uses the main ideas of the DADAL method) works specifically for the sparse PCA problem. As we have seen above, using the standard code of DADAL involves increasing dimension to $(d + 3d^2 + 2)^2$ which appears to be quiet inefficient for solving the standard SDP relaxation of sparse PCA. Therefore we alternatively pursued the following approach: Consider the primal and dual SDP relaxation of sparse PCA,

$$
\begin{array}{ll}
\text{Primal} := \min_{X,Y} & \langle -A, X \rangle \\
\text{s.t.} & \langle I, X \rangle \leq 1 \quad (\mu_1 \geq 0) \\
& \langle \mathbf{1}\mathbf{1}^\top, Y \rangle \leq k \quad (\mu_2 \geq 0) \\
& Y \geq X \quad (\gamma^+ \geq 0) \\
& Y \geq -X \quad (\gamma^- \geq 0) \\
& X \succeq \mathbf{0} \quad (Z \succeq \mathbf{0})
\end{array}
\qquad
\begin{array}{ll}
\text{Dual} := \max & -\mu_1 - \mu_2 k \\
\text{s.t.} & \mu_1 I + \gamma^+ - \gamma^- - A - Z = \mathbf{0} \\
& \mu_2 \mathbf{1}\mathbf{1}^\top - \gamma^+ - \gamma^- = \mathbf{0} \\
& Z \succeq \mathbf{0} \\
& \mu_1, \mu_2, \gamma^+, \gamma^- \geq 0
\end{array}
$$

with its augmented Lagrangian

$$\mathcal{L}_\sigma(\mu, \gamma, Z; X, Y) := -\mu_1 - \mu_2 k + \langle M_1, X \rangle + \langle M_2, Y \rangle - \frac{\sigma}{2} \|M_1\|_F^2 - \frac{\sigma}{2} \|M_2\|_F^2,$$

where $M_1, M_2$ are defined as

$$M_1 := \mu_1 I + \gamma^+ - \gamma^- - A - Z,$$

$$M_2 := \mu_2 \mathbf{1}\mathbf{1}^\top - \gamma^+ - \gamma^-.$$

---

[3] https://www.math.aau.at/or/Software/

We initialize $\boldsymbol{X}^0, \boldsymbol{Y}^0, \boldsymbol{Z}^0$ as follows: Compute eigenvalue decomposition of $\boldsymbol{A} = \boldsymbol{V} \boldsymbol{\Lambda}_A \boldsymbol{V}^\top$, let $\boldsymbol{v}_1$ be the leading eigenvector of $\boldsymbol{V}$ with respect to the largest eigenvalue. Set

$$\boldsymbol{X}^0 \leftarrow \boldsymbol{v}_1 \boldsymbol{v}_1^\top,$$

$$\boldsymbol{Y}^0 \leftarrow |\boldsymbol{X}^0|,$$

$$\boldsymbol{Z}^0 \leftarrow \boldsymbol{0},$$

along with the starting augmented Lagrangian parameter $\sigma^0$. In $(k+1)$-th iteration, update each variable based on the following rule which is similar as the DADAL method proposed in [16].

$$\boldsymbol{\mu}^{k+1}, \boldsymbol{\gamma}^{k+1} \leftarrow \underset{\boldsymbol{\mu} \geq 0, \boldsymbol{\gamma} \geq 0}{\arg\max} \mathcal{L}_{\sigma^k}(\boldsymbol{\mu}, \boldsymbol{\gamma}, \boldsymbol{Z}^k; \boldsymbol{X}^k, \boldsymbol{Y}^k)$$

$$\boldsymbol{Z}^{k+1} \leftarrow \left( -\frac{\boldsymbol{X}^k}{\sigma^k} + \mu_1^{k+1} \boldsymbol{I} + (\boldsymbol{\gamma}^+)^{k+1} - (\boldsymbol{\gamma}^-)^{k+1} - \boldsymbol{A} \right)_{\succeq 0}$$

$$\boldsymbol{X}^{k+1} \leftarrow -\sigma \cdot \left( -\frac{\boldsymbol{X}^k}{\sigma^k} + \mu_1^{k+1} \boldsymbol{I} + (\boldsymbol{\gamma}^+)^{k+1} - (\boldsymbol{\gamma}^-)^{k+1} - \boldsymbol{A} \right)_{\preceq 0}$$

$$\boldsymbol{Y}^{k+1} \leftarrow |\boldsymbol{X}^{k+1}|$$

Update $\sigma$ based on Algorithm 1 in [16]

where $(\boldsymbol{A})_{\succeq 0}, (\boldsymbol{A})_{\preceq 0}$ denote the positive semi-definite, negative semi-definite part of symmetric matrix $\boldsymbol{A}$. That is: Let $\boldsymbol{A} = \boldsymbol{U} \boldsymbol{\Sigma} \boldsymbol{U}^\top$ be its eigenvalue decomposition. Represent $\boldsymbol{\Sigma} = \boldsymbol{\Sigma}^+ + \boldsymbol{\Sigma}^-$ where $\boldsymbol{\Sigma}_{ii}^+ = \max\{\boldsymbol{\Sigma}_{ii}, 0\}$ and $\boldsymbol{\Sigma}_{ii}^- = \min\{\boldsymbol{\Sigma}_{ii}, 0\}$, then

$$(\boldsymbol{A})_{\succeq 0} := \boldsymbol{U} \boldsymbol{\Sigma}^+ \boldsymbol{U}^\top,$$

$$(\boldsymbol{A})_{\preceq 0} := \boldsymbol{U} \boldsymbol{\Sigma}^- \boldsymbol{U}^\top.$$

REMARK 1. The way we update our dual variables (and primal variables) in each iteration, there is no guarantee that the dual variables satisfy the equality constraints in the dual, namely,

$$\boldsymbol{M}_1 := \mu_1 \boldsymbol{I} + \boldsymbol{\gamma}^+ - \boldsymbol{\gamma}^- - \boldsymbol{A} - \boldsymbol{Z} = 0,$$

$$\boldsymbol{M}_2 := \mu_2 \boldsymbol{1} \boldsymbol{1}^\top - \boldsymbol{\gamma}^+ - \boldsymbol{\gamma}^- = 0.$$

Therefore, it is not true that we can always obtain exact dual bounds from every iteration. We store the dual bounds of iterations where the equality constraints are satisfied within a tolerance of 0.01, i.e.,

$$\|\boldsymbol{M}_1\|_F + \|\boldsymbol{M}_2\|_F \leq 0.01.$$

Moreover, after the final iteration, we add one more step by solving the following *linear program,*

$$\mu^{\text{final}}, \boldsymbol{\gamma}^{\text{final}} := \arg\max_{\mu, \boldsymbol{\gamma}} \quad -\mu_1 - \mu_2 k$$
$$\text{s.t.} \quad \mu_1 \boldsymbol{I} + \boldsymbol{\gamma}^+ - \boldsymbol{\gamma}^- - \boldsymbol{A} - \boldsymbol{Z}^{\text{final}} = \boldsymbol{0},$$
$$\mu_2 \boldsymbol{1}\boldsymbol{1}^\top - \boldsymbol{\gamma}^+ - \boldsymbol{\gamma}^- = \boldsymbol{0}, \quad \text{(final-dual)}$$
$$\mu_1, \mu_2, \boldsymbol{\gamma}^+, \boldsymbol{\gamma}^- \geq 0,$$

where $\boldsymbol{Z}^{\text{final}} \succeq 0$ is the dual variable obtained in the final step of DADAL-SPCA. It is easy to observe that $(\mu^{\text{final}}, \boldsymbol{\gamma}^{\text{final}}, \boldsymbol{Z}^{\text{final}})$ is a dual feasible solution, and therefore a dual bound can be obtained from this dual feasible solution.

*Stopping criteria:* The stopping criteria includes three conditions. Meeting any of the criteria stops the DADAL-SPCA algorithm.

  (a) The maximum number of iteration is set to be 200.

  (b) The stopping criteria quantity $\delta$ proposed in Algorithm 1 [16] is set to be 0.001, i.e., at the end of each iteration, we compute the primal and dual infeasibility errors as follows:

$$r_P := \frac{\max\{\text{Tr}(X) - 1, 0\} + \max\{\langle \boldsymbol{1}\boldsymbol{1}^\top, \boldsymbol{Y}\rangle - k, 0\}}{1 + \sqrt{1 + k^2}},$$
$$r_D := \frac{\|\boldsymbol{M}_1\|_F + \|\boldsymbol{M}_2\|_F}{1 + \|\boldsymbol{A}\|_F},$$

and set $\delta := \max\{r_P, r_D\}$.

  (c) Since there is no closed form solution of the following updating step:

$$\boldsymbol{\mu}^{k+1}, \boldsymbol{\gamma}^{k+1} \leftarrow \underset{\boldsymbol{\mu} \geq 0, \boldsymbol{\gamma} \geq 0}{\arg\max} \mathcal{L}_{\sigma^k}(\boldsymbol{\mu}, \boldsymbol{\gamma}, \boldsymbol{Z}^k; \boldsymbol{X}^k, \boldsymbol{Y}^k),$$

we use commercial solver Gurobi (called via Python) to solve this quadratic programming sub-problem in each iteration. For small instances (i.e., $d < 500$, Pitprops, Eisen-1, Eisen-2), the total time limit given for Gurobi solver is 3600 seconds (1 hour); and for middle-size instance (i.e., $d = 500$, CovColon, Lymp), the total time limit given for Gurobi solver is 7200 seconds (2 hours), and for large instance (i.e., $d = 2000$, Reddit), the total time limit given for Gurobi solver is 18000 seconds (5 hours).

Algorithm 5 is the pseudocode of finding dual bounds using DADAL-SPCA.

The gap obtained by DADAL-SPCA as described above with various values of $\sigma^0$ is reported in Table 15.

The "Time" column in Table 15 denotes the total running time used for the DADAL-SPCA method. We can see that the "Time" of CovColon, Lymp reported in Table 15 are greater than time limit for solver, since additional time are required to implement the other four updating steps in each iteration. The out of memory (O.M.) for Reddit instance is due to the memory limitation to load Reddit instance $d = 2000$ for the update step

$$\boldsymbol{\mu}^{k+1}, \boldsymbol{\gamma}^{k+1} \leftarrow \underset{\boldsymbol{\mu} \geq 0, \boldsymbol{\gamma} \geq 0}{\arg\max} \mathcal{L}_{\sigma^k}(\boldsymbol{\mu}, \boldsymbol{\gamma}, \boldsymbol{Z}^k; \boldsymbol{X}^k, \boldsymbol{Y}^k).$$

---

**Algorithm 5** Dual Bound DADAL-SPCA

---

1: *Input*: Covariance matrix $\boldsymbol{A}$, sparsity parameter $k$, maximum number of iteration $T_{\max}$, total time limit for solver $T_{\text{total}}$, starting Lagrangian augmented parameter $\sigma^0$.

2: *Output*: Dual bound of sparse PCA.

3: **function** DUAL BOUND METHOD($\boldsymbol{A}, k, T_{\max}, T_{\text{total}}$)

4:     Compute eigenvalue decomposition on $\boldsymbol{A}$, let $\boldsymbol{v}_1$ be its leading eigenvector.

5:     Initialize $\boldsymbol{X} \leftarrow \boldsymbol{v}_1 \boldsymbol{v}_1^\top, \boldsymbol{Y} \leftarrow |\boldsymbol{X}|, \boldsymbol{Z} \leftarrow \boldsymbol{0}^{d \times d}, (\mu_1, \mu_2) \leftarrow (0, 0), \boldsymbol{\gamma}^\pm \leftarrow \boldsymbol{0}^{d \times d}$.

6:     Run DADAL-SPCA with stopping criteria described above with starting Lagrangian augmented parameter $\sigma^0 \in \{0.001, 0.01, 0.1, 1\}$, and return $\text{UB}^{\text{DADAL-SPCA}}$.

7:     Solve final-dual for a dual bound $\text{UB}^{\text{final-dual}}$.

8:     **return** $\text{UB} \leftarrow \min\{\text{UB}^{\text{final-dual}}, \text{UB}^{\text{DADAL-SPCA}}\}$.

9: **end function**

---

**Table 15**    DADAL-SPCA under different starting augmented Lagrangian parameter $\sigma^0$.

| INSTANCE \ $\sigma^0$ | LB | $\sigma^0 = 0.001$ | | $\sigma^0 = 0.01$ | | $\sigma^0 = 0.1$ | | $\sigma^0 = 1$ | |
|---|---|---|---|---|---|---|---|---|---|
| | | GAP % | TIME | GAP % | TIME | GAP % | TIME | GAP % | TIME |
| PITPROPS $k=5$ | 3.406 | 3.96 | 6 | 1.79 | 5 | 1.70 | 2 | **1.64** | 3 |
| EISEN-1 $k=10$ | 17.33 | 2.23 | 270 | **2.19** | 225 | 11.07 | 294 | 39.10 | 288 |
| EISEN-2 $k=10$ | 11.71 | 2.32 | 1053 | 2.37 | 610 | **2.08** | 898 | 2.12 | 897 |
| COVCOLON $k=10$ | 2641 | 14.16 | 7492 | **13.51** | 7281 | 19.05 | 7369 | 26.82 | 7301 |
| LYMP $k=10$ | 6008 | **29.67** | 7339 | 34.79 | 7331 | 46.84 | 7367 | 59.09 | 7373 |
| REDDIT $k=10$ | 1052 | - | O.M. | - | O.M. | - | O.M. | - | O.M. |

We tried to solve the final-dual linear program for Reddit instance, but the LP did not solve in 5 hours. (This LP has order $d^2$ variables, whereas the number of variables of convex integer program is order $dI_{\text{pos}}N$ and $I_{\text{pos}}N \ll d$ in this instance.)

To complete the comparison, we also list the comparison between our model in paper and DADAL, DADAL-SPCA, Mosek in Table 16.

**Table 16**    Compare with existing SDP methods

| INSTANCE | LB | MODEL-IN-PAPER | | DADAL [16] | | DADAL-SPCA (BEST) | | MOSEK | |
|---|---|---|---|---|---|---|---|---|---|
| | | GAP % | TIME | GAP % | TIME | GAP % | TIME | GAP % | TIME |
| PITPROPS $k=5$ | 3.406 | 3.26 | 0.4 | 82.43 | 593 | 1.64 | 3 | **1.52** | 5 |
| EISEN-1 $k=10$ | 17.33 | **0.115** | 63 | - | O.M. | 2.19 | 225 | 2.19 | 15 |
| EISEN-2 $k=10$ | 11.71 | **1.71** | 385 | - | O.M. | 2.08 | 898 | 1.96 | 52 |
| COVCOLON $k=10$ | 2641 | **2.37** | 28 | - | O.M. | 13.51 | 7281 | - | O.M. |
| LYMP $k=10$ | 6008 | **17.86** | 4225 | - | O.M. | 29.67 | 7339 | - | O.M. |
| REDDIT $k=10$ | 1052 | **2.24** | 8584 | - | O.M. | - | O.M. | - | O.M. |

Based on Table 16, we observe that the SDP-relaxation solved by Mosek produces the best bounds for the small instances (Pitprops, Eisen-1, Eisen-2), while DADAL-SPCA is able to produce bounds for Pitprops, Eisen-1, Eisen-2, CovColon, and Lymp. However, as we can see, except for Pitprops, the best dual bounds are obtained by solving convex IP model of this paper.