

The Sparse Principal Component Analysis Problem: Optimality Conditions and Algorithms

Amir Beck · Yakov Vaisbourd

Received: date / Accepted: date

Abstract Sparse principal component analysis addresses the problem of finding a linear combination of the variables in a given data set with a sparse coefficients vector that maximizes the variability of the data. This model enhances the ability to interpret the principal components, and is applicable in a wide variety of fields including genetics and finance, just to name a few.

We suggest a necessary coordinate-wise-based optimality condition, and show its superiority over the stationarity-based condition that is commonly used in the literature, and which is the basis for many of the algorithms designed to solve the problem. We devise algorithms that are based on the new optimality condition, and provide numerical experiments that support our assertion that algorithms, which are guaranteed to converge to stronger optimality conditions, perform better than algorithms that converge to points satisfying weaker optimality conditions.

Keywords optimality conditions · principal component analysis · sparsity constrained problems · stationarity · numerical methods

1 Introduction

Principal component analysis (PCA) is a well known data-analytic technique that linearly transforms a given set of data to some equivalent representation. This transformation is defined in such a manner that any variable in the new representation, called a *principal component* (PC), expresses most of the variance in the data, which is not expressed by the PCs that precede it. The

Amir Beck (corresponding author)
Faculty of Industrial Engineering and Management, Technion, Haifa, Israel
E-mail: becka@ie.technion.ac.il

Yakov Vaisbourd
Faculty of Industrial Engineering and Management, Technion, Haifa, Israel
E-mail: yakovv@campus.technion.ac.il

linear combination defining each of the PCs is given by a coefficients (also termed *loadings*) vector. In terms of the covariance (or correlation) matrix of the data, the coefficients vector of the k -th PC is the eigenvector that corresponds to the k -th largest eigenvalue [1]. One major drawback of PCA is that commonly the coefficients vectors are dense, i.e., each PC is a linear combination of much, if not most, of the original variables, which causes a difficulty in interpreting the obtained PCs. This disadvantage encouraged a wide interest in the sparsity constrained version of PCA, which imposes an additional constraint, enforcing the coefficients vector not to exceed some predetermined sparsity level s .

Enforcing sparsity on the coefficients vector is commonly acceptable in some applications. For example, in the exploration of micro-array gene expression patterns, PCA is employed in order to classify different tissues according to their gene expression. It is also desirable that such discrimination can be executed by utilizing only a small subset of the genes, thus encouraging sparse solutions [2]. The desire to obtain interpretable coefficients vectors is not the only reason to favor the sparse PCA model. For example, some financial applications will prefer sparse solutions in order to reduce transaction costs [3]. Clearly, incorporating an additional sparsity constraint will provide a PC that, generally, does not explain all of the variance which is explained by the regular PC; nevertheless, in such applications, this sacrifice is acceptable with respect to the obtained benefits. We refer to this formulation as the sparsity constrained formulation, and it is merely one of several alternative formulations considered in the literature. The common alternatives are the result of treating the sparsity term, or its relaxation, by a penalty approach. The sparse PCA problem is a difficult non-convex problem, and can be optimally solved only for small scale problems by performing exhaustive or a branch and bound search over all possible support sets [4]. Thus, in order to handle large scale problems, the algorithms proposed in the literature are seeking to find an approximate solution. One of the first methods, suggested by Cadima and Jolliffe [5], is to threshold the smallest, in absolute value, elements of the dominant eigenvector. Unfortunately, this remarkably simple approach is known to frequently provide poor results. In [4] Moghaddam et al. proposed several greedy methods. An advantage of these methods is that they produce a full path of solutions (i.e., a solution for each of the values of sparsity level up to s), but the necessity to perform a large amount of eigenvalue computations at each step render them quite computationally expensive. In [6], d'Aspremont et al. proposed an approximate greedy approach that obviates the necessity to perform most of the eigenvalue computations by evaluating a lower bound on the eigenvalues, which results in a substantial reduction of computation time. Another approach presented by d'Aspremont et al. in [6], and earlier in [7], is to consider a semidefinite programming formulation with a rank constraint for some of the relaxed and/or penalized models of PCA. These equivalent formulations are still hard non-convex problems, and thus a relaxed model is solved and an approximate solution is derived for the original problem. The algorithms used to solve the SDP relaxations are not

applicable for large scale problems, rendering this approach as non-scalable. In [8], encouraged by the LASSO approach suggested for regression [9], Jolliffe et al. proposed the absolute value norm constrained formulation under the name SCoTLass (simplified component technique LASSO), which is a relaxation of the sparsity constrained problem. In practice, the numerical study was conducted on the penalized version by implementing the projected gradient algorithm. The relaxed model was further considered in the literature. An alternating minimization scheme to solve the constrained formulation was proposed in [10]. Another work that addressed the constrained formulation was motivated by the expectation maximization algorithm for probabilistic PCA [11]. Even though the work addressed the constrained formulation, the sequence generated by the method in [11] is guaranteed to be s -sparse. Penalized versions were also considered extensively. In [12] Zou et al. formulated the sparse PCA as a regression-type model, where the i -th principal component was approximated by the linear combination of the original variables. A LASSO and ridge penalties are imposed on the coefficients vector forming the elastic net model that generalizes the LASSO [13] and an alternating minimization algorithm, called SPCA, was proposed. In [14] Shen and Huang proposed several iterative schemes to solve the penalized versions via regularized SVD. These methods were considered further in [15], where a gradient scheme was proposed and a convergence analysis, that was missing in [14], was also provided.

Recently, Luss and Teboulle showed in [16] that the seemingly different methods proposed in [15, 14, 11, 17, 10, 12] are some particular realizations of the conditional gradient algorithm with unit step-size. The work [16] proposed a unified algorithmic framework which they refer to as ConGradU (the well known conditional gradient with unit step size) and established convergence results, showing that the algorithm produces a point satisfying some necessary first order optimality criteria. Some novel schemes are provided. One of them addresses directly the sparsity constrained formulation of sparse PCA.

As already mentioned, none of the methods listed above can guarantee to produce an optimal solution. In addition, the sparse PCA problem does not seem to possess a verifiable necessary and sufficient global optimality condition, and hence, in general, there is no efficient way to check if a given vector is the global optimal solution¹. Therefore, the comparison of the methods in the literature is based solely on numerical experiments without providing any theoretical justification for the advantage of a certain method over the others. However, most of the algorithms just listed will produce a solution that satisfies some necessary optimality condition. In a recent work, Beck and Eldar [18] employed some of the aforesaid conditions in order to provide an insight regarding the success of the corresponding algorithms. Under the framework of minimizing a continuously differentiable function subject to a sparsity constraint, several necessary optimality conditions were presented.

¹ In [6] the authors suggested a sufficient optimality condition.

The relations between the different optimality conditions were established, showing that some of the conditions are stronger (that is, more restrictive) than others. An extension to problems over sparse symmetric sets was considered in [19]. In this paper, we adopt this methodology in order to establish a hierarchy between two necessary optimality conditions for the sparsity constrained sparse PCA problem. The first condition that we consider is a well known first order condition, that was originally presented in the context of the sparse PCA problem in [16]. We will refer to it as the *complete (co) stationarity* condition. Much of the existing algorithms in the literature are actually guaranteed to converge to a co-stationary point. The second condition, which we call *coordinate-wise (CW) maximality*, is a generalization of one of the conditions considered in [18], and it essentially states that the function value cannot be improved by making changes of at most two coordinates.

In the following section we will explicitly define the conditions under consideration. In Section 3, we will establish the relation between the conditions, showing that the CW-maximality condition is stronger (that is, more restrictive) than co-stationarity. In Section 4, we will introduce algorithms that produce points satisfying the aforementioned conditions and finally, in Section 5, we will provide a numerical study on simulated and real life data that supports our assertion that algorithms that correspond to stronger conditions are more likely to provide better results.

2 Necessary Optimality Conditions

Throughout the paper, we consider the following sparsity constrained problem:

$$\max\{f(\mathbf{x}) : \mathbf{x} \in S\}, \quad (\text{P})$$

where f is a continuously differentiable convex function over \mathbb{R}^n and

$$S := \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\|_2 \leq 1, \|\mathbf{x}\|_0 \leq s\},$$

with $\|\cdot\|_0$ being the so-called l_0 -norm defined by $\|\mathbf{x}\|_0 := |\{i : x_i \neq 0\}|$ ². As a special case, when the objective function is chosen as $f(\mathbf{x}) = \mathbf{x}^T \mathbf{A} \mathbf{x}$, where \mathbf{A} is a given positive semidefinite matrix, problem (P) amounts to the *l_0 -constrained sparse PCA model*:

$$\max\{\mathbf{x}^T \mathbf{A} \mathbf{x} : \mathbf{x} \in S\}. \quad (\text{SPCA})$$

In PCA applications, \mathbf{A} usually stands for the covariance matrix of the data.

In this section, we will present two necessary optimality conditions for the general model (P). Although our main motivation is to study the sparse PCA problem, we will nonetheless consider the general model (P), since our results are also applicable in this general setting.

² Note that the l_0 -norm is not actually a norm since it does not satisfy the absolute homogeneity property.

Prior to presenting the optimality conditions, we will introduce in the following subsection some notation and definitions that will be used in our analysis.

2.1 Notation and Definitions

A subvector of a given vector $\mathbf{x} \in \mathbb{R}^n$ corresponding to a set of indices $T \subseteq \{1, 2, \dots, n\}$ is denoted by \mathbf{x}_T . Similarly, we will denote the subvector of the gradient $\nabla f(\mathbf{x})$ corresponding to the indices in T by $\nabla_T f(\mathbf{x})$. The sign of a given $\alpha \in \mathbb{R}$ is denoted by $\text{sgn}(\alpha)$ and is equal to 1 for $\alpha \geq 0$ and -1 for $\alpha < 0$. The support set of some arbitrary vector \mathbf{x} will be denoted by $I_1(\mathbf{x}) = \{i : x_i \neq 0\}$ and its complement by $I_0(\mathbf{x}) = \{i : x_i = 0\}$. For a given vector $\mathbf{x} \in \mathbb{R}^n$ and an integer $s \in \{1, 2, \dots, n-1\}$, we will define $M_s(\mathbf{x})$ to be the s -th largest absolute value component in \mathbf{x} . For such \mathbf{x} and s , we will define the sets $I_{>}(\mathbf{x}, s)$, $I_{= }(\mathbf{x}, s)$ and $I_{<}(\mathbf{x}, s)$ as follows:

$$\begin{aligned} I_{>}(\mathbf{x}, s) &:= \{i : |x_i| > M_s(\mathbf{x})\}, \\ I_{= }(\mathbf{x}, s) &:= \begin{cases} \{i : |x_i| = M_s(\mathbf{x})\}, & \|\mathbf{x}\|_0 \geq s, \\ \emptyset, & \|\mathbf{x}\|_0 < s, \end{cases} \\ I_{<}(\mathbf{x}, s) &:= \begin{cases} \{i : |x_i| < M_s(\mathbf{x})\}, & \|\mathbf{x}\|_0 \geq s, \\ \{i : x_i = 0\}, & \|\mathbf{x}\|_0 < s. \end{cases} \end{aligned}$$

We will also define the set $I_{\geq}(\mathbf{x}, s) := I_{>}(\mathbf{x}, s) \cup I_{= }(\mathbf{x}, s)$ and the set $I_{\leq}(\mathbf{x}, s) := I_{<}(\mathbf{x}, s) \cup I_{= }(\mathbf{x}, s)$. Obviously, the sets $I_{>}(\mathbf{x}, s)$, $I_{= }(\mathbf{x}, s)$ and $I_{<}(\mathbf{x}, s)$ form a partition of $\{1, 2, \dots, n\}$. Furthermore, when $\|\mathbf{x}\|_0 < s$, we have that $I_{>}(\mathbf{x}, s) = I_1(\mathbf{x})$, $I_{= }(\mathbf{x}, s) = \emptyset$ and $I_{<}(\mathbf{x}, s) = I_0(\mathbf{x})$.

The sets defined above possess some convenient and elementary properties which are given in Lemma 2.1 below. Since all the properties stated in the lemma are rather simple consequences of the definition of the sets $I_{>}(\mathbf{x}, s)$, $I_{= }(\mathbf{x}, s)$, $I_{<}(\mathbf{x}, s)$, the proof is omitted.

Lemma 2.1 1. If $\mathbf{x} \neq \mathbf{0}$, then $I_{\geq}(\mathbf{x}, s) \neq \emptyset$.

2. If $|I_{\geq}(\mathbf{x}, s)| < s$ then $x_j = 0$ for all $j \in I_{<}(\mathbf{x}, s)$.

3. For any $i \in I_{>}(\mathbf{x}, s)$, $j \in I_{= }(\mathbf{x}, s)$ and $k \in I_{<}(\mathbf{x}, s)$, it holds that $|x_i| > |x_j| > |x_k|$.

We will frequently use the notation

$$R_s(\mathbf{x}) := \{T : I_{>}(\mathbf{x}, s) \subseteq T \subseteq I_{\geq}(\mathbf{x}, s), |T| = \min\{s, |I_{\geq}(\mathbf{x}, s)|\}$$

for the set containing all the subsets of indices corresponding to the nonzero s largest in absolute value components of a given vector \mathbf{x} . When $\|\mathbf{x}\|_0 \leq s$, there are no more than s nonzero elements in \mathbf{x} , and the above definition actually amounts to $R_s(\mathbf{x}) = \{I_1(\mathbf{x})\}$. However, when $\|\mathbf{x}\|_0 > s$, there might be more than one set of indices corresponding to the s largest absolute value

components of \mathbf{x} . For example, consider the vector $\mathbf{x} = (3, 2, 1, 1, 1, 0, 0)^T$ and the sparsity level $s = 3$. Then,

$$R_3(\mathbf{x}) = \{\{1, 2, 3\}, \{1, 2, 4\}, \{1, 2, 5\}\}.$$

On the other hand, in the following examples, the set contains a single subset:

$$R_3((0, -5, 4, -3, 2, 0)^T) = \{\{2, 3, 4\}\}, R_3((0, 0, 4, -3, 0, 0)^T) = \{\{3, 4\}\}.$$

The *hard thresholding* operator maps a vector $\mathbf{x} \in \mathbb{R}^n$ to the set of vectors that are generated by keeping the s largest absolute value components of \mathbf{x} and setting all the others to zeros. This operator, which we denote by H_s , is formally defined by

$$H_s(\mathbf{x}) := \bigcup_{T \in R_s(\mathbf{x})} \{\mathbf{y} : \mathbf{y}_T = \mathbf{x}_T, \mathbf{y}_{\bar{T}} = \mathbf{0}\}.$$

Thus, for example,

$$H_3((3, 2, 1, 1, 1, 0, 0)^T) = \{(3, 2, 1, 0, 0, 0, 0)^T, (3, 2, 0, 1, 0, 0, 0)^T, (3, 2, 0, 0, 1, 0, 0)^T\}.$$

2.2 Complete (co) - Stationarity

The first condition that we consider was presented for the sparse PCA problem in [16]. We refer to it as the complete (co) stationarity condition.

Definition 2.1 (co-stationarity) Let \mathbf{x} be a feasible solution of (P). Then, \mathbf{x} is called a co-stationary point of (P) over S if and only if it satisfies:

$$\langle \nabla f(\mathbf{x}), \mathbf{v} - \mathbf{x} \rangle \leq 0 \quad \forall \mathbf{v} \in S.$$

This is probably the most elementary first order condition for constrained differentiable optimization problems. The work [16] provided a unified framework for several algorithms designed to solve different formulations of sparse PCA. Actually, [16] considered the co-stationarity condition over a general nonempty and compact set instead of S , and for this general case, the following proposition, which was originally established in [20], was recalled. This result follows from the convexity of the objective function.

Proposition 2.1 Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a continuous differentiable and convex function over \mathbb{R}^n , and let C be a nonempty and compact set. If \mathbf{x} is a global maximum of f over C , then \mathbf{x} is a co-stationary point over C , meaning that $\langle \nabla f(\mathbf{x}), \mathbf{v} - \mathbf{x} \rangle \leq 0$ for any $\mathbf{v} \in C$.

2.3 CW-Maximality

The second necessary optimality condition that we will consider is coordinate-wise maximality. This optimality condition is in fact a type of a local optimality condition, stating that a given point \mathbf{x} is a minimizer over a neighbourhood consisting of all feasible points, that are different by at most two coordinates. We will denote the corresponding neighbourhood by

$$S_2(\mathbf{x}) := \{\mathbf{z} : \|\mathbf{z} - \mathbf{x}\|_0 \leq 2, \mathbf{z} \in S\}.$$

The formal definition of a CW-maximum point follows.

Definition 2.2 (CW-maximum point) Let \mathbf{x} be a feasible solution of (P). Then, \mathbf{x} is called a coordinate-wise (CW) maximum point of (P) if and only if $f(\mathbf{x}) \geq f(\mathbf{z})$ for every $\mathbf{z} \in S_2(\mathbf{x})$.

Obviously, CW-maximality, by its definition, is a necessary optimality condition.

Proposition 2.2 *Let \mathbf{x} be an optimal solution to (P). Then, \mathbf{x} is an CW-maximum point.*

Instead of considering the neighbourhood $S_2(\mathbf{x})$ in the definition of CW-maximality (Definition 2.2), we could have alternatively considered larger neighbourhoods consisting of vectors that differ from \mathbf{x} by at most k coordinates for some $2 \leq k \leq s$:

$$S_k(\mathbf{x}) := \{\mathbf{z} : \|\mathbf{z} - \mathbf{x}\|_0 \leq k, \mathbf{z} \in S\}.$$

A similar optimality condition over such a neighbourhood can be defined, and clearly since $S_t(\mathbf{x}) \subseteq S_k(\mathbf{x})$ for any $t \leq k$, considering neighbourhoods that differ by a larger amount of coordinates will result in stronger optimality conditions. Note that the amount of comparisons required in order to verify that a vector $\mathbf{x} \in \mathbb{R}^n$ with a full support ($I_1(\mathbf{x}) = s$) is CW-maximal ($k = 2$) is $O(s \cdot n)$, while changing the neighbourhood to S_3 will increase the amount of comparisons to $O(s \cdot n^2)$. Hence, considering such a stronger optimality condition has a substantial computational price. Keeping in mind that we seek scalable conditions and algorithms, we restrict the discussion to the case $k = 2$.

3 Optimality Conditions Hierarchy

Our main result in this section is that CW-maximality is a stronger (that is, more restrictive) optimality condition than co-stationarity. This result also has an impact on the performance of the corresponding algorithms in the sense that, loosely speaking, algorithms that are only guaranteed to converge to a co-stationary point are less likely to produce the optimal solution of the problem than algorithms that are guaranteed to converge to a CW-maximal point. In Section 5, we will show that the numerical results support this assertion.

3.1 Technical Preliminaries

We will begin by providing some auxiliary technical results that will be used in order to establish the main result. Lemma 3.1 is a trivial result, that follows directly from the Cauchy-Schwarz inequality (see also Lemma 4.1 in [16]).

Lemma 3.1 *Suppose that $\mathbf{0} \neq \mathbf{q} \in \mathbb{R}^d$ and $\rho > 0$. Then, the optimal solution of the optimization problem*

$$\max_{\mathbf{x} \in \mathbb{R}^d} \{\mathbf{q}^T \mathbf{x} : \|\mathbf{x}\|_2 \leq \rho\}, \quad (\text{QCLP})$$

is given by $\mathbf{x}^ = \rho \frac{\mathbf{q}}{\|\mathbf{q}\|_2}$ with the optimal value of $\rho \|\mathbf{q}\|_2$.*

The following simple lemma is an extension of Proposition 4.3 from [16].

Lemma 3.2 *Assume that $\mathbf{0} \neq \mathbf{p} \in \mathbb{R}^n$. Then, the set of optimal solutions of the optimization problem*

$$\max_{\mathbf{x} \in \mathbb{R}^n} \{\mathbf{p}^T \mathbf{x} : \|\mathbf{x}\|_0 \leq s, \|\mathbf{x}\|_2 \leq 1\}, \quad (\text{S-QCLP})$$

is given by

$$\mathbf{X}^*(\mathbf{p}, s) := \left\{ \frac{\mathbf{x}}{\|\mathbf{x}\|_2} : \mathbf{x} \in H_s(\mathbf{p}) \right\},$$

with the optimal value of $\|\mathbf{p}_T\|_2$, where $T \in R_s(\mathbf{p})$.

Proof We can write (S-QCLP) as

$$\max_{\substack{T \subseteq \{1, \dots, n\} \\ |T| \leq s}} \max_{\mathbf{x} \in \mathbb{R}^n} \{\mathbf{p}^T \mathbf{x} : \|\mathbf{x}\|_2 \leq 1, I_1(\mathbf{x}) \subseteq T\}. \quad (1)$$

According to Lemma 3.1, for each $T \subseteq \{1, \dots, n\}$ satisfying $|T| \leq s$, the optimal value of the inner optimization problem is $\|\mathbf{p}_T\|_2$, and if $\mathbf{p}_T \neq \mathbf{0}$, then a solution \mathbf{x}^* to the inner optimization problem is given by

$$\mathbf{x}_T^* = \frac{\mathbf{p}_T}{\|\mathbf{p}_T\|_2}, \quad \mathbf{x}_{\bar{T}}^* = \mathbf{0}. \quad (2)$$

The problem (1) thus reduces to

$$\max_{\substack{T \subseteq \{1, \dots, n\} \\ |T| \leq s}} \|\mathbf{p}_T\|_2. \quad (3)$$

Obviously, when $\|\mathbf{p}\|_0 \geq s$, the optimal solutions of the latter problem are all the sets containing the indices of components corresponding to the s largest absolute values in \mathbf{p} , and when $\|\mathbf{p}\|_0 < s$, the unique optimal solution is $I_1(\mathbf{p})$. Thus, the set of all optimal solutions of (3) is $R_s(\mathbf{p})$. Noting that $\mathbf{p}_T \neq \mathbf{0}$ for any $T \in R_s(\mathbf{p})$, we conclude that the optimal solutions of (S-QCLP) are given by (2) with T being any set in $R_s(\mathbf{p})$, which are exactly the members of $\mathbf{X}^*(\mathbf{p}, s)$. \square

Our final technical lemma states that, if a given vector $\tilde{\mathbf{x}}$ is *not* an optimal solution of the problem of maximizing a linear function over the unit norm, then there must be two indices $i \neq j$ for which the subvector $\tilde{\mathbf{x}}_{\{i,j\}}$ is also *not* an optimal solution for the problem restricted to the variables x_i, x_j (while fixing all the other variables). This lemma is rather simple, but will play a key role in the proof of the main result.

Lemma 3.3 *Let $\mathbf{q} \in \mathbb{R}^d$ and $\rho > 0$. Suppose that $\tilde{\mathbf{x}}$ satisfies $\|\tilde{\mathbf{x}}\|_2 \leq \rho$, and that it is not an optimal solution of (QCLP). Then, there exist indices $i, j (i \neq j)$ such that $\tilde{\mathbf{x}}_{\{i,j\}}$ is not the optimal solution of*

$$\max_{\mathbf{x}_{\{i,j\}} \in \mathbb{R}^2} \left\{ \mathbf{q}_{\{i,j\}}^T \mathbf{x}_{\{i,j\}} : \|\mathbf{x}_{\{i,j\}}\|_2 \leq \left(\rho^2 - \sum_{l \neq i,j} \tilde{x}_l^2 \right)^{1/2} \right\}. \quad (2\text{-QCLP}_{\{i,j\}})$$

Proof Since $\tilde{\mathbf{x}}$ is not the optimal solution of (QCLP), we obtain that $\mathbf{q} \neq \mathbf{0}$ (since otherwise, if $\mathbf{q} = \mathbf{0}$, all feasible points are also optimal). Thus, the set $I_1(\mathbf{q})$ is nonempty. We will split the analysis into two cases.

- If $\|\tilde{\mathbf{x}}\|_2 < \rho$, then take any $i \in I_1(\mathbf{q})$ and $j \neq i$, and we can write

$$\|\tilde{\mathbf{x}}_{\{i,j\}}\|_2 < \left(\rho^2 - \sum_{l \neq i,j} \tilde{x}_l^2 \right)^{1/2},$$

which together with $\mathbf{q}_{\{i,j\}} \neq \mathbf{0}$ (since $i \in I_1(\mathbf{q})$) implies that $\tilde{\mathbf{x}}_{\{i,j\}}$ is not the optimal solution of (2-QCLP $_{\{i,j\}}$), since we have, by Lemma 3.1, that the constraint at the optimal solution must be active.

- If, on the other hand, $\|\tilde{\mathbf{x}}\|_2 = \rho$, then assume in contradiction that for each $i \neq j$ the vector $\tilde{\mathbf{x}}_{\{i,j\}}$ is the optimal solution of (2-QCLP $_{\{i,j\}}$). Take some $i \in I_1(\mathbf{q})$. For any $j \in I_0(\mathbf{q})$, we know that $\tilde{\mathbf{x}}_{\{i,j\}}$ is the optimal solution of (2-QCLP $_{\{i,j\}}$) and thus, according to Lemma 3.1 (employed on the problem (2-QCLP $_{\{i,j\}}$)), it must in particular satisfy $\tilde{x}_j = 0$, that is, $j \in I_0(\tilde{\mathbf{x}})$. To summarize,

$$\tilde{x}_j = 0 \text{ for any } j \in I_0(\mathbf{q}). \quad (4)$$

Now, for any $j \in I_1(\mathbf{q})$, according to Lemma 3.1, $\tilde{\mathbf{x}}_{\{i,j\}}$ must satisfy

$$\tilde{x}_i = \frac{q_i}{\|(q_i, q_j)^T\|_2} (\tilde{x}_i^2 + \tilde{x}_j^2)^{1/2}, \quad (5)$$

where here we used the fact that $\rho^2 - \sum_{l \neq i,j} \tilde{x}_l^2 = \tilde{x}_i^2 + \tilde{x}_j^2$. Squaring both sides of (5), we obtain that it is equivalent to $q_j^2 \tilde{x}_i^2 = q_i^2 \tilde{x}_j^2$, and hence

$$\tilde{x}_j^2 = \frac{q_j^2}{q_i^2} \tilde{x}_i^2 \quad \text{for any } j \in I_1(\mathbf{q}).$$

By (4), $\tilde{x}_j = 0$ whenever $j \in I_0(\mathbf{q})$, and we can therefore write

$$\tilde{x}_j^2 = \frac{q_j^2}{q_i^2} \tilde{x}_i^2, \quad j = 1, 2, \dots, n.$$

Summing over $j = 1, 2, \dots, n$, and using the fact that $\|\tilde{\mathbf{x}}\|_2^2 = \rho^2$, it follows that

$$\sum_{j=1}^n \tilde{x}_i^2 \frac{q_j^2}{q_i^2} = \rho^2,$$

implying that

$$\tilde{x}_i^2 = \rho^2 \frac{q_i^2}{\|\mathbf{q}\|_2^2},$$

which combined with the fact that $\text{sgn}(\tilde{x}_i) = \text{sgn}(q_i)$ (see (5)), yields

$$\tilde{x}_i = \rho \frac{q_i}{\|\mathbf{q}\|_2}.$$

Since we actually proved the latter for an arbitrary $i \in I_1(\mathbf{q})$, and since $\tilde{x}_i = 0$ for any $i \in I_0(\mathbf{q})$ (see (4)), it follows that

$$\mathbf{x} = \rho \frac{\mathbf{q}}{\|\mathbf{q}\|_2},$$

in contradiction to the assumption that $\tilde{\mathbf{x}}$ is not an optimal solution of (QCLP). \square

The following corollary is a direct consequence of Lemmas 3.2 and 3.3.

Corollary 3.1 *Let $\tilde{\mathbf{x}} \in S$. If $\tilde{\mathbf{x}}$ is not an optimal solution to (S-QCLP) and $I_1(\tilde{\mathbf{x}}) \subseteq T$ for some $T \in R_s(\mathbf{p})$, then there exist indices $i, j \in T (i \neq j)$ such that $\tilde{\mathbf{x}}_{\{i,j\}}$ is not an optimal solution of (2-QCLP $_{\{i,j\}}$).*

Proof Assume that $|T| = k$. Since $\tilde{\mathbf{x}}$ is not an optimal solution of (S-QCLP), it follows by Lemma 3.2 that $\tilde{\mathbf{x}}_T \neq \frac{\mathbf{p}_T}{\|\mathbf{p}_T\|}$, which implies that $\tilde{\mathbf{x}}_T$ is not the optimal solution of the restricted problem

$$\min_{\mathbf{y} \in \mathbb{R}^k} \{ \mathbf{p}_T^T \mathbf{y} : \|\mathbf{y}\|_2 \leq \rho \}.$$

Therefore, invoking Lemma 3.3 with $d = k$, $\mathbf{q} = \mathbf{p}_T$, it follows that there exist indices $i, j \in T (i \neq j)$ such that $\tilde{\mathbf{x}}_{i,j}$ is not an optimal solution of (2-QCLP $_{\{i,j\}}$). \square

3.2 Co-Stationarity vs. CW-Maximality

The main result of this paper is given in the following theorem, which establishes the superiority of the CW-maximality condition over the co-stationarity condition.

Theorem 3.1 *Let \mathbf{x} be a CW-maximum point of problem (P). Then, \mathbf{x} is a co-stationary point of (P).*

Proof Let \mathbf{x} be a CW-maximum point of (P). Assume by contradiction that \mathbf{x} is not a co-stationary point. This means that there exists a vector $\mathbf{v} \in S$ such that

$$\nabla f(\mathbf{x})^T(\mathbf{v} - \mathbf{x}) > 0. \quad (6)$$

We will show that we can find a vector $\mathbf{z} \in S_2(\mathbf{x})$ such that

$$\nabla f(\mathbf{x})^T(\mathbf{z} - \mathbf{x}) > 0. \quad (7)$$

This will imply a contradiction to the CW-maximality of \mathbf{x} by the following simple argument: since f is a convex function, we have

$$f(\mathbf{z}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{z} - \mathbf{x}),$$

which combined with (7) implies that

$$f(\mathbf{z}) > f(\mathbf{x}),$$

which is an obvious contradiction to the CW-maximality of \mathbf{x} .

Since \mathbf{x} satisfies (6), we obviously have $\nabla f(\mathbf{x}) \neq \mathbf{0}$. Let $X^*(\nabla f(\mathbf{x}), s)$ be the set of optimal solutions of (S-QCLP) with $\mathbf{p} = \nabla f(\mathbf{x})$ and let $\mathbf{x}^* \in X^*(\nabla f(\mathbf{x}), s)$ be some particular solution. Then,

$$\nabla f(\mathbf{x})^T \mathbf{x}^* \geq \nabla f(\mathbf{x})^T \mathbf{v} > \nabla f(\mathbf{x})^T \mathbf{x},$$

and thus $\mathbf{x} \notin X^*(\nabla f(\mathbf{x}), s)$.

Suppose that there exists some l for which $\nabla_l f(\mathbf{x}) \cdot x_l < 0$ (and in particular $l \in I_1(\mathbf{x})$). Define \mathbf{z} as:

$$j = 1, \dots, n \quad z_j := \begin{cases} -x_l, & j = l, \\ x_j, & \text{otherwise.} \end{cases}$$

$\mathbf{z} \in S_2(\mathbf{x})$ and $\nabla f(\mathbf{x})^T(\mathbf{z} - \mathbf{x}) > 0$ since

$$\nabla f(\mathbf{x})^T(\mathbf{z} - \mathbf{x}) = -2 \cdot \nabla_l f(\mathbf{x}) \cdot x_l > 0.$$

We have thus shown in this case the desired contradiction. From now on, we will therefore consider the case where $\nabla_i f(\mathbf{x}) \cdot x_i \geq 0$ for all $i = 1, \dots, n$. Consider the following cases:

1. $I_1(\mathbf{x}) \not\subseteq I_{\geq}(\nabla f(\mathbf{x}), s)$.

Obviously, there is some $h \in I_1(\mathbf{x}) \cap I_{<}(\nabla f(\mathbf{x}), s)$.

We will consider the following subcases:

- 1.1. If $|I_{\geq}(\nabla f(\mathbf{x}), s)| < s$, then $\nabla_h f(\mathbf{x}) = 0$ (by Lemma 2.1, part 2), and since $\nabla f(\mathbf{x}) \neq \mathbf{0}$, we conclude, using Lemma 2.1 (part 1), that there is some $l \in I_{\geq}(\nabla f(\mathbf{x}), s)$. Define \mathbf{z} as:

$$j = 1, \dots, n \quad z_j := \begin{cases} \operatorname{sgn}(\nabla_l f(\mathbf{x})) \cdot (x_h^2 + x_l^2)^{1/2} & j = l, \\ 0 & j = h, \\ x_j & \text{otherwise.} \end{cases}$$

Obviously $\mathbf{z} \in S_2(\mathbf{x})$, and in addition $\nabla f(\mathbf{x})^T(\mathbf{z} - \mathbf{x}) > 0$ since

$$\begin{aligned} \nabla f(\mathbf{x})^T(\mathbf{z} - \mathbf{x}) &= \nabla_l f(\mathbf{x}) \cdot \operatorname{sgn}(\nabla_l f(\mathbf{x})) \cdot (x_h^2 + x_l^2)^{1/2} - \nabla_l f(\mathbf{x}) \cdot x_l \\ &= \left| \nabla_l f(\mathbf{x}) \right| \cdot (x_h^2 + x_l^2)^{1/2} - \nabla_l f(\mathbf{x}) \cdot x_l \\ &= \left| \nabla_l f(\mathbf{x}) \right| \cdot (x_h^2 + x_l^2)^{1/2} - \left| \nabla_l f(\mathbf{x}) \right| \cdot |x_l| \quad (\nabla_l f(\mathbf{x}) \cdot x_l \geq 0) \\ &= \left| \nabla_l f(\mathbf{x}) \right| \cdot ((x_h^2 + x_l^2)^{1/2} - |x_l|) > 0 \quad (\nabla_l f(\mathbf{x}) \neq 0, x_h \neq 0). \end{aligned}$$

- 1.2. If $|I_{\geq}(\nabla f(\mathbf{x}), s)| \geq s$, then there is some $l \in I_{\geq}(\nabla f(\mathbf{x}), s)$ such that $l \notin I_1(\mathbf{x})$. Otherwise $I_{\geq}(\nabla f(\mathbf{x}), s) \subseteq I_1(\mathbf{x})$, and since $|I_{\geq}(\nabla f(\mathbf{x}), s)| \geq s$ and $|I_1(\mathbf{x})| \leq s$, we have that $I_{\geq}(\nabla f(\mathbf{x}), s) = I_1(\mathbf{x})$, contradicting our assumption that $I_1(\mathbf{x}) \not\subseteq I_{\geq}(\nabla f(\mathbf{x}), s)$. We will define \mathbf{z} as:

$$j = 1, \dots, n \quad z_j := \begin{cases} \operatorname{sgn}(\nabla_l f(\mathbf{x})) \cdot |x_h| & j = l, \\ 0 & j = h, \\ x_j & \text{otherwise.} \end{cases}$$

Clearly, $\mathbf{z} \in S_2(\mathbf{x})$. In addition, $\nabla f(\mathbf{x})^T(\mathbf{z} - \mathbf{x}) > 0$ since:

$$\begin{aligned} \nabla f(\mathbf{x})^T(\mathbf{z} - \mathbf{x}) &= \nabla_l f(\mathbf{x}) \cdot \operatorname{sgn}(\nabla_l f(\mathbf{x})) \cdot |x_h| \\ &\quad - \nabla_h f(\mathbf{x}) \cdot x_h \\ &= \left| \nabla_l f(\mathbf{x}) \right| \cdot |x_h| - \left| \nabla_h f(\mathbf{x}) \right| \cdot |x_h| \quad (\nabla_h f(\mathbf{x}) \cdot x_h \geq 0) \\ &= (\left| \nabla_l f(\mathbf{x}) \right| - \left| \nabla_h f(\mathbf{x}) \right|) \cdot |x_h| > 0, \end{aligned}$$

where the last inequality holds since $x_h \neq 0$ and the indices l and h are such that $l \in I_{\geq}(\nabla f(\mathbf{x}), s)$ and $h \in I_{<}(\nabla f(\mathbf{x}), s)$, thus according to Lemma 2.1 (part 3) $|\nabla_l f(\mathbf{x})| > |\nabla_h f(\mathbf{x})|$.

2. $I_1(\mathbf{x}) \subseteq I_{\geq}(\nabla f(\mathbf{x}), s)$

Now we will consider the following subcases:

- 2.1. If $I_1(\mathbf{x}) \subseteq T$ for some $T \in R_s(\nabla f(\mathbf{x}))$, then since $\mathbf{x} \notin X^*(\nabla f(\mathbf{x}), s)$, it follows that according to Corollary 3.1, there exist indices $h, l \in T$ such that

$$\hat{\mathbf{x}} := \operatorname{argmax}_{\mathbf{y} \in \mathbb{R}^2} \left\{ \nabla_{\{h,l\}} f(\mathbf{x})^T \mathbf{y} : \|\mathbf{y}\|^2 \leq 1 - \sum_{i \neq h,l} x_i^2 \right\}$$

satisfies

$$\nabla_{\{h,l\}} f(\mathbf{x})^T \hat{\mathbf{x}} > \nabla_{\{h,l\}} f(\mathbf{x})^T \mathbf{x}_{\{h,l\}}. \quad (8)$$

Since $|T| \leq s$ and $\|\hat{\mathbf{x}}\|_2^2 \leq 1 - \sum_{i \neq h,l} x_i^2$, the vector

$$j = 1, \dots, n \quad z_j := \begin{cases} \hat{x}_1, & j = h, \\ \hat{x}_2, & j = l, \\ x_j, & \text{otherwise,} \end{cases}$$

is in $S_2(\mathbf{x})$, and satisfies by (8) that $\nabla f(\mathbf{x})^T(\mathbf{z} - \mathbf{x}) > 0$.

2.2. If $I_1(\mathbf{x}) \not\subseteq T$ for all $T \in R_s(\nabla f(\mathbf{x}))$, then:

- Take $h \in I_1(\mathbf{x})$ such that $h \notin T$ for some $T \in R_s(\nabla f(\mathbf{x}))$. Since $I_1(\mathbf{x}) \subseteq I_{\geq}(\nabla f(\mathbf{x}), s)$, it follows that $h \in I_{\geq}(\nabla f(\mathbf{x}), s)$. Moreover, since $I_{>}(\nabla f(\mathbf{x}), s) \subseteq T$ and $h \notin T$, we have that $h \notin I_{>}(\nabla f(\mathbf{x}), s)$, implying that $h \in I_{= }(\nabla f(\mathbf{x}), s)$. Thus, $h \in I_{= }(\nabla f(\mathbf{x}), s) \cap I_1(\mathbf{x})$.
- $I_{>}(\nabla f(\mathbf{x}), s) \not\subseteq I_1(\mathbf{x})$. To show this, note that otherwise, $I_{>}(\nabla f(\mathbf{x}), s) \subseteq I_1(\mathbf{x})$, and since $I_1(\mathbf{x}) \subseteq I_{\geq}(\nabla f(\mathbf{x}), s)$ and $|I_1(\mathbf{x})| \leq s$, we obtain that $|I_1(\mathbf{x})| \leq \min\{s, |I_{\geq}(\nabla f(\mathbf{x}), s)|\}$, implying that $I_1(\mathbf{x}) \subseteq T$ for some $T \in R_s(\nabla f(\mathbf{x}))$, in contradiction to our assumption. Thus, there exists some $l \in I_{>}(\nabla f(\mathbf{x}), s)$ such that $l \notin I_1(\mathbf{x})$.

Define \mathbf{z} as:

$$j = 1, \dots, n \quad z_j := \begin{cases} \operatorname{sgn}(\nabla_l f(\mathbf{x})) \cdot |x_h|, & j = l, \\ 0, & j = h, \\ x_j, & \text{otherwise.} \end{cases}$$

Clearly, $\mathbf{z} \in S_2(\mathbf{x})$. Furthermore, $\nabla f(\mathbf{x})^T(\mathbf{z} - \mathbf{x}) > 0$ since

$$\begin{aligned} \nabla f(\mathbf{x})^T(\mathbf{z} - \mathbf{x}) &= \nabla_l f(\mathbf{x}) \cdot \operatorname{sgn}(\nabla_l f(\mathbf{x})) \cdot |x_h| \\ &\quad - \nabla_h f(\mathbf{x}) \cdot x_h \\ &= |\nabla_l f(\mathbf{x})| \cdot |x_h| - |\nabla_h f(\mathbf{x})| \cdot |x_h| \quad (\nabla_h f(\mathbf{x}) \cdot x_h \geq 0) \\ &= (|\nabla_l f(\mathbf{x})| - |\nabla_h f(\mathbf{x})|) \cdot |x_h| > 0, \end{aligned}$$

where the last inequality holds since $x_h \neq 0$ and the indices l and h are such that $l \in I_{>}(\nabla f(\mathbf{x}), s)$ and $h \in I_{= }(\nabla f(\mathbf{x}), s)$, and thus according to Lemma 2.1 (part 3) $|\nabla_l f(\mathbf{x})| > |\nabla_h f(\mathbf{x})|$.

We have thus arrived at a contradiction, and the desired implication is established. \square

In order to show that the reverse implication is not valid, that is, that co-stationary points are not necessarily CW-maximal points, we present an example of a problem instance and a co-stationary point, that is not a CW-maximal point.

Example 3.1 For any $n > s > 0$, we consider problem (SPCA) with a diagonal matrix \mathbf{A} , whose entries on the main diagonal are given by the vector \mathbf{a} defined by

$$\mathbf{a} := \begin{pmatrix} 2 \cdot \mathbf{1}_{n-s} \\ 0.5 \cdot \mathbf{1}_s \end{pmatrix},$$

where for a given positive integer m , $\mathbf{1}_m$ and $\mathbf{0}_m$ are the vectors of size m with all entries equal to ones or zeros, respectively. We also define

$$\mathbf{x} := \begin{pmatrix} \mathbf{0}_{n-s} \\ s^{-0.5} \cdot \mathbf{1}_s \end{pmatrix} \quad \text{and} \quad \tilde{\mathbf{x}} := \begin{pmatrix} \mathbf{0}_{n-s-1} \\ s^{-0.5} \\ 0 \\ s^{-0.5} \cdot \mathbf{1}_{s-1} \end{pmatrix}.$$

It is easy to see that $\mathbf{x}, \tilde{\mathbf{x}} \in S$ and that $\mathbf{A} \succ \mathbf{0}$, since it is a diagonal matrix with positive diagonal elements. The gradient of f is given by:

$$\nabla f(\mathbf{x}) = 2\mathbf{A}\mathbf{x} = \begin{pmatrix} \mathbf{0}_{n-s} \\ s^{-0.5} \cdot \mathbf{1}_s \end{pmatrix}.$$

For any $\mathbf{v} \in S$:

$$\begin{aligned} \langle \nabla f(\mathbf{x}), \mathbf{v} - \mathbf{x} \rangle &= \sum_{i=n-s+1}^n s^{-0.5} (v_i - s^{-0.5}) \\ &= s^{-0.5} \left(\sum_{i=n-s+1}^n v_i - s^{0.5} \right) \leq s^{-0.5} (\|\mathbf{v}\|_1 - s^{0.5}) \leq 0, \end{aligned}$$

where the last inequality holds since $\|\mathbf{v}\|_1 \leq \sqrt{\|\mathbf{v}\|_0} \|\mathbf{v}\|_2 \leq \sqrt{s}$. Hence, \mathbf{x} is co-stationary. The vector $\tilde{\mathbf{x}}$ satisfies $\tilde{\mathbf{x}} \in S_2(\mathbf{x})$ and since:

$$\begin{aligned} f(\tilde{\mathbf{x}}) = \tilde{\mathbf{x}}^T \mathbf{A} \tilde{\mathbf{x}} &= (s-1) \cdot (2s)^{-1} + 2s^{-1} = (s+3) \cdot (2s)^{-1} \\ &> s \cdot (2s)^{-1} = \mathbf{x}^T \mathbf{A} \mathbf{x} = f(\mathbf{x}), \end{aligned}$$

it follows that \mathbf{x} is not a CW-maximum point.

3.3 Support Optimality

Theorem 3.1 establishes the relationship between the two stationarity conditions considered up to this point: co-stationarity and CW-maximality. A third condition, proposed in [4], that we will refer to as *support optimality* (SO), is given in the following definition.

Definition 3.1 (Support Optimality) A vector $\mathbf{x}^* \in S$ is called a **support optimal (SO) point** of (P) with respect to an index set $T \subseteq \{1, 2, \dots, n\}$ if and only if it is an optimal solution of the optimization problem

$$\max_{\mathbf{x} \in \mathbb{R}^n} \{f(\mathbf{x}) : \|\mathbf{x}\|_2 \leq 1, I_1(\mathbf{x}) \subseteq T\}. \quad (\text{SO})$$

It is clear that, if $\mathbf{x} \in S$ is an optimal solution of problem (P), then it must be an SO point of (P) with respect to any index set T satisfying $|T| \leq s$ and $I_1(\mathbf{x}) \subseteq T$. In that respect, support optimality is a necessary optimality condition for problem (P). It is a remarkably weak condition and cannot be used exclusively to derive a reasonable algorithm. Nevertheless, it is not totally futile. In order to enhance the performance, the CW-based algorithms that will be presented in Section 4 will produce a sequence of SO points, and in Section 5, we will adopt the variational re-normalization strategy suggested in [4], stating that for each sparse solution obtained by any technique, it is reasonable to replace this solution with the SO point that correspond to the same support.

We will conclude this section with an example that demonstrates the potential benefit of employing algorithms that produce a point that satisfies stronger necessary optimality conditions. Consider the pit-prop data, which consists of 13 variables measuring various physical properties of 180 pit-props. This data set was suggested originally in [21], and since then was extensively used as a benchmark example for sparse PCA; see, for example, [8, 15, 4]. The problem has 13 variables and we consider a sparsity level of $s = 4$. Note that we can list all the $\binom{13}{4} = 715$ SO points that correspond to index sets with exactly 4 indices, and the optimal solution must be one of these 715 points. Out of this set of points, 28 satisfy the co-stationarity condition and only 2 satisfy the CW-maximality condition. The following table presents the support sets of each of the co-stationarity points along with their function values.

Table 1 The supports of the co-stationary points for the pit prop data.

#	Support	CW-maximum	Value	#	Support	CW-maximum	Value
1	{1,2,9,10}	*	2.937	15	{5,6,7,10}		2.337
2	{1,2,7,10}		2.883	16	{7,8,10,12}		2.314
3	{1,2,7,9}		2.859	17	{7,8,10,13}		2.302
4	{1,2,8,9}		2.797	18	{5,6,7,13}		2.28
5	{1,2,8,10}		2.759	19	{3,4,6,7}		2.209
6	{1,2,6,7}		2.697	20	{4,5,6,7}		2.196
7	{2,7,9,10}		2.696	21	{7,10,12,13}		2.136
8	{2,6,7,10}		2.592	22	{3,4,8,12}		1.995
9	{1,6,7,10}		2.587	23	{3,4,10,12}		1.992
10	{1,2,3,4}		2.563	24	{3,10,11,12}		1.609
11	{7,8,9,10}		2.549	25	{3,5,12,13}		1.516
12	{6,7,9,10}		2.522	26	{1,5,12,13}		1.414
13	{6,7,10,13}		2.459	27	{2,5,12,13}		1.408
14	{6,7,8,10}		2.444	28	{3,5,11,13}		1.382

Since the number of CW-maximum points is significantly smaller than the number of co-stationary points, it is much more probable that the optimal solution will be found by an algorithm that produces CW-maximum points than an algorithm that produces co-stationary points.

4 Algorithms

In this section, we will present two CW-based algorithms – GCW and PCW – that are guaranteed to converge after a finite amount of iterations to a CW-maxima. Later on, in Section 5, we will demonstrate the superiority of these algorithm over methods which are based on the co-stationarity optimality condition such as the conditional gradient algorithm with unit step-size (ConGradU), that was suggested in [16], where it was also proven that limit points of the sequence generated by ConGradU are co-stationary point.

In [18] several algorithms that produce a CW-minimum point were considered. These block coordinate descent type algorithms perform at each iteration an optimization step with respect to one or two variables, while keeping the rest fixed. The coordinates that need to be altered are chosen to be the ones that produce the maximal decrease among all possible alternatives, or by applying an index selection strategy based on a local first order information. We adopt this approach and present similar algorithms for the sparse PCA problem.

At each iteration of a CW-based algorithm applied to (P), at most two variables will be updated. We can categorize each of the iterations according to whether the support is altered or not. Block coordinate algorithms suffer from a major drawback – a slow convergence rate. In order to reduce the effect of this displeasing characteristic, we will replace the point obtained at each step with an SO point that corresponds to the same support. This modification allows us to bypass the large amount of iterations that should have been devoted for optimizing the variables with respect to a fixed support.

Below we present the Greedy CW (GCW) algorithm. We denote by $\mathcal{O}(T)$ an oracle that produce an SO point with respect to a given support T by solving problem (SO). We will refer to this oracle as an *SO oracle*. In the specific case of the PCA problem, the SO oracle amounts to finding a normalized principal eigenvector of a submatrix of the covariance matrix. However, finding the maximum of a general convex function f over a unit ball is in principle a difficult task. We will assume that the solution produced by the oracle is uniquely defined by T . In addition, note that the oracle outputs an optimal solution of a problem consisting of maximizing a convex function over a compact and convex feasible set, and hence by [20, Corollary 32.3.2], there exists an optimal solution of the problem which is an extreme point. In particular, this means that we can assume without any loss of generality that the oracle outputs a vector with norm 1. This assumption will be made from now on.

The Greedy CW (GCW) Algorithm

Input: $f : \mathbb{R}^n \rightarrow \mathbb{R}$ – convex function; $\mathcal{O}(\cdot)$ – SO oracle; s – sparsity level.

Output: \mathbf{x} – a CW-maximum point of (SPCA).

Initialization: Take $T \in \{1, 2, \dots, n\}$ such that $1 \leq |T| \leq s$ and set $\mathbf{x}^0 = \mathcal{O}(T)$ and $k = 0$.

General step:

1. While $\|\mathbf{x}^k\|_0 < s$, compute

$$j_k \in \operatorname{argmax}_{j \in I_0(\mathbf{x}^k)} \{f(\mathbf{z}) : \mathbf{z} = \mathcal{O}(I_1(\mathbf{x}^k) \cup \{j\})\},$$

If $f(\mathcal{O}(I_1(\mathbf{x}^k) \cup \{j_k\})) > f(\mathbf{x}^k)$, then set

$$\begin{aligned} \mathbf{x}^{k+1} &= \mathcal{O}(I_1(\mathbf{x}^k) \cup \{j_k\}), \\ k &= k + 1, \end{aligned}$$

and return to **1**; otherwise, go to **2**.

2. For every $i \in I_1(\mathbf{x}^k)$ and $j \in I_0(\mathbf{x}^k)$ compute

$$f_{i,j} = \max_{\sigma \in \{-1,1\}} \{f(\mathbf{x}^k - x_i^k \mathbf{e}_i + \sigma |x_i^k| \mathbf{e}_j)\}.$$

Let $(i_k, j_k) = \operatorname{argmax} \{f_{i,j} : i \in I_1(\mathbf{x}^k), j \in I_0(\mathbf{x}^k)\}$. If $f_{i_k, j_k} > f(\mathbf{x}^k)$, then set

$$\begin{aligned} \mathbf{x}^{k+1} &= \mathcal{O}((I_1(\mathbf{x}^k) \setminus \{i_k\}) \cup \{j_k\}), \\ k &= k + 1, \end{aligned}$$

and return to **1**.

Otherwise, STOP and set $\mathbf{x} \leftarrow \mathbf{x}^{k+1}$.

Step **1** of the GCW algorithm is in fact the greedy forward selection method proposed in [4]. Hence, in some sense, the GCW method is a generalization of this method, that does not terminate at the moment that a solution with a full support is obtained. However, from a more practical point of view, this resemblance is irrelevant due to the fact that, if the initial support satisfies $|T| = s$, then the condition $\|\mathbf{x}^k\|_0 < s$ will probably be false for all k in any reasonable practical scenario.

The following theorem summarizes the key properties of the GCW algorithm.

Theorem 4.1 *Let $\{\mathbf{x}^k\}$ be the sequence generated by the GCW algorithm. Then, the following statements hold.*

- (i) *The sequence of function values $\{f(\mathbf{x}^k)\}$ is monotonically increasing.*
- (ii) *The algorithm terminates after a finite amount of iterations.*
- (iii) *At termination, the algorithm produces a CW maximum point.*

Proof Part (i) follows immediately from the description of the GCW algorithm. Part (ii) is a consequence of the monotonicity of the algorithm (part (i)) and the fact that it only passes through SO points, from which there is only a finite number under the standing assumption that the solution produced by the oracle $\mathcal{O}(T)$ is uniquely defined by T .

To prove (iii), consider the following partition of $S_2(\mathbf{x})$:

$$\begin{aligned} S_2(\mathbf{x}) &= \{\mathbf{z} : \|\mathbf{z} - \mathbf{x}\|_0 \leq 2, \mathbf{z} \in S\} \\ &= S_2^0(\mathbf{x}) \cup S_2^1(\mathbf{x}) \cup S_2^2(\mathbf{x}), \end{aligned}$$

where

$$\begin{aligned} S_2^0(\mathbf{x}) &= \{\mathbf{z} \in S : \|\mathbf{z} - \mathbf{x}\|_0 \leq 2, I_1(\mathbf{z}) \subseteq I_1(\mathbf{x})\} \\ S_2^1(\mathbf{x}) &= \{\mathbf{z} \in S : \|\mathbf{z} - \mathbf{x}\|_0 \leq 2, I_1(\mathbf{z}) = I_1(\mathbf{x}) \cup \{j\}, j \in I_0(\mathbf{x})\} \\ S_2^2(\mathbf{x}) &= \\ &\{\mathbf{z} \in S : \|\mathbf{z} - \mathbf{x}\|_0 \leq 2, I_1(\mathbf{z}) = (I_1(\mathbf{x}) \setminus \{i\}) \cup \{j\}, i \in I_1(\mathbf{x}), j \in I_0(\mathbf{x})\}, \end{aligned}$$

and assume that the algorithm produced the point $\bar{\mathbf{x}}$. Since $\bar{\mathbf{x}}$ is an SO point and $S_2^0(\bar{\mathbf{x}}) \subseteq \{\mathbf{x} : \|\mathbf{x}\|_2 \leq 1, I_1(\mathbf{x}) \subseteq I_1(\bar{\mathbf{x}})\}$, it follows that $f(\bar{\mathbf{x}}) \geq f(\mathbf{x})$ for any $\mathbf{x} \in S_2^0(\bar{\mathbf{x}})$. Now, note that the algorithm terminates only if after performing Step 2 we obtain that for any $i \in I_1(\bar{\mathbf{x}})$ and $j \in I_0(\bar{\mathbf{x}})$

$$\begin{aligned} f_{i,j} &= \max_{\sigma \in \{-1,1\}} \{f(\bar{\mathbf{x}} - \bar{x}_i \mathbf{e}_i + \sigma |\bar{x}_i| \mathbf{e}_j)\} \\ &= \max_{\alpha} \{f(\bar{\mathbf{x}} - \bar{x}_i \mathbf{e}_i + \alpha \mathbf{e}_j) : \alpha \in [-|\bar{x}_i|, |\bar{x}_i|]\} \\ &\leq f(\bar{\mathbf{x}}), \end{aligned}$$

where the first equality is due to the fact that the maximum of a convex function over a compact and convex set is attained at an extreme point, see [20, Corolalry 32.3.2]. Thus, $f(\bar{\mathbf{x}}) \geq f(\mathbf{x})$ for any $\mathbf{x} \in S_2^2(\bar{\mathbf{x}})$. This is enough for proving that $\bar{\mathbf{x}}$ is CW-maximal in the case when $\|\bar{\mathbf{x}}\|_0 = s$ since in this case $S_2^1(\bar{\mathbf{x}}) = \emptyset$. If $\|\bar{\mathbf{x}}\|_0 < s$, then prior to entering Step 2, Step 1 must be performed. This step is terminated only if $f(\bar{\mathbf{x}}) \geq f(\mathbf{x})$ for any

$$\mathbf{x} \in \{\mathbf{z} \in S : I_1(\mathbf{z}) = I_1(\bar{\mathbf{x}}) \cup \{j\}, j \in I_0(\bar{\mathbf{x}})\},$$

and since $S_2^1(\bar{\mathbf{x}}) \subseteq \{\mathbf{z} \in S : I_1(\mathbf{z}) = I_1(\bar{\mathbf{x}}) \cup \{j\}, j \in I_0(\bar{\mathbf{x}})\}$, it implies that $f(\bar{\mathbf{x}}) \geq f(\mathbf{x})$ for any $\mathbf{x} \in S_2^1(\bar{\mathbf{x}})$, concluding that $f(\bar{\mathbf{x}}) \geq f(\mathbf{x})$ for any $\mathbf{x} \in S_2(\bar{\mathbf{x}})$. \square

Practically, if the initial support T satisfies $|T| = s$, then most of the computation time in the GCW method is consumed in computing $f_{i,j}$ for each possible swap. This observation encourages us to consider the following variation of GCW, which we name *the Partial CW (PCW) algorithm*.

The Partial CW (PCW) Algorithm

Input: $f : \mathbb{R}^n \rightarrow \mathbb{R}$ – convex function; $\mathcal{O}(\cdot)$ – SO oracle; s – sparsity level.

Output: \mathbf{x} – a CW-maximum point of (SPCA).

Initialization: Take $T \in \{1, 2, \dots, n\}$ such that $1 \leq |T| \leq s$ and set $\mathbf{x}^0 = \mathcal{O}(T)$ and $k = 0$.

General step:

1. While $\|\mathbf{x}^k\|_0 < s$, compute

$$j_k \in \operatorname{argmax}_{j \in I_0(\mathbf{x}^k)} \{f(\mathbf{z}) : \mathbf{z} = \mathcal{O}(I_1(\mathbf{x}^k) \cup \{j\})\},$$

If $f(\mathcal{O}(I_1(\mathbf{x}^k) \cup \{j_k\})) > f(\mathbf{x}^k)$, then set

$$\begin{aligned} \mathbf{x}^{k+1} &= \mathcal{O}(I_1(\mathbf{x}^k) \cup \{j_k\}), \\ k &= k + 1, \end{aligned}$$

and return to 1; otherwise, go to 2.

2. Set $R = I_1(\mathbf{x}^k)$.

While $|R| > 0$

Set $i_k \in \operatorname{argmin} \{|x_i^k| : i \in R\}$ and for each $j \in I_0(\mathbf{x}^k)$ compute

$$f_{i_k, j} = \max_{\sigma \in \{-1, 1\}} \{f(\mathbf{x}^k - x_{i_k}^k \mathbf{e}_{i_k} + \sigma |x_{i_k}^k| \mathbf{e}_j)\}.$$

Let $j_k \in \operatorname{argmax} \{f_{i_k, j} : j \in I_0(\mathbf{x}^k)\}$.

If $f_{i_k, j_k} > f(\mathbf{x}^k)$, then set

$$\begin{aligned} \mathbf{x}^{k+1} &= \mathcal{O}((I_1(\mathbf{x}^k) \setminus \{i_k\}) \cup \{j_k\}), \\ k &= k + 1, \end{aligned}$$

and return to 1.

Otherwise, set $R = R \setminus \{i_k\}$.

STOP and set $\mathbf{x} \leftarrow \mathbf{x}^{k+1}$.

Before termination, PCW will perform the computation of all possible $f_{i,j}$, thus assuring the convergence to a CW-maximum point, given that the output is of a full support. For the general step, the amount of computation will significantly decrease on the expense of finding the indices that provide the maximal increase in the function value. Nevertheless, the empirical study suggests that PCW provides similar results as GCW with respect to function values in a fraction of the time, as demonstrated in Section 5.

5 Numerical Results

We will illustrate the effectiveness of the algorithms proposed in the previous section on simulated and a gene expression datasets. We compared the results with the following alternative algorithms: the novel l_0 -constrained version of ConGradU [16], the expectation maximization [11], approximate greedy [6] and thresholding [5]. The MATLAB implementation of ConGradU was kindly provided by the authors, for all the other alternative algorithms we used a MATLAB implementation available on the authors' web-pages. For

the thresholding algorithm and the algorithms proposed in this paper, we used a MATLAB implementation, which is available in the following URL:

http://tx.technion.ac.il/~yakovv/packages/CW_PCA.zip

Whenever an initialization is required, we set the initial point to be the solution of the thresholding method. Regarding the output, we adopt the variational renormalization strategy suggested in [4]. Hence, for each of the algorithms, we extracted the sparsity pattern (the set of indices of the nonzero elements). The actual output vector is determined to be equal to $\mathcal{O}(T)$, where T is the generated sparsity pattern. The experiments were conducted on a PC with a 3.40GHz processor with 16GB RAM.

5.1 Random Data

The covariance matrix \mathbf{A} is given by $\mathbf{A} = \mathbf{D}^T \mathbf{D}$, where \mathbf{D} is the so-called "data matrix". Each entry in the data matrix $\mathbf{D} \in \mathbb{R}^{m \times n}$ was randomly generated according to the Gaussian distribution with zero mean and variance $1/m$ ($D_{i,j} \sim \mathcal{N}(0, 1/m)$). We considered data matrices with $n = 2000, 5000, 10,000$ and $50,000$ variables. The number of observations is set to $m = 150$ for all matrices. The sparsity levels considered are $s = 5, 10, \dots, 250$, and for each sparsity level we generated 100 realizations. We will measure the effectiveness of the algorithms according to the average proportion of variability explained by the algorithm with respect to the largest eigenvalue of the data covariance matrix (i.e., $\mathbf{x}^T \mathbf{A} \mathbf{x} / \lambda_1(\mathbf{A})$, where \mathbf{x} is the solution and $\lambda_1(\mathbf{A})$ is the largest eigenvalue of \mathbf{A}).

5.1.1 GCW vs. PCW

First, we would like to compare the effectiveness and performance of the CW-based algorithms proposed in the previous section: GCW and PCW. We conducted the comparison based on data matrices with 2,000 variables and the results are given in Figure 1.

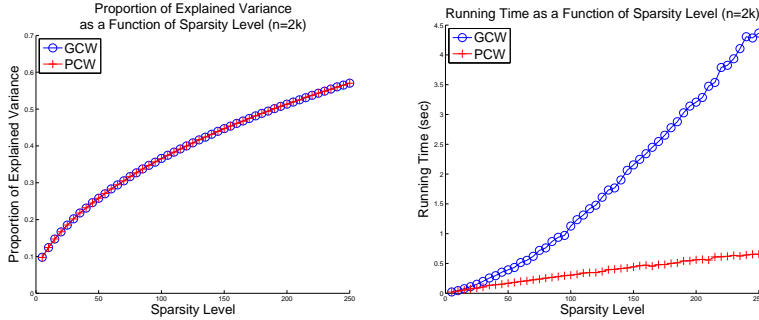


Fig. 1 GCW vs. PCW - The proportion of explained variability is given in the left figure and the computation time is given in the right one. The plot in both figures are given as a function of the sparsity level.

We can clearly see that both methods achieve similar results with respect to the function values, while PCW achieves these results in a fraction of the time. Thus, in the remaining numerical study we will omit GCW. Although the partial version remarkably reduces the computation time, it is still not competitive for very large-scale problems when a full path of solutions is required. Thus, for such cases, we will also examine the effect of initializing PCW with the solution of the previous run (with the smaller sparsity level), and we will refer to such a continuation scheme as PCW_{cont} .

5.1.2 PCW vs. Alternative Methods

We will now compare the effectiveness and performance of PCW with respect to the alternative algorithms mentioned earlier. The setting for this set of experiments is the same as the one described in the previous example, but with problems with $n = 5,000$, $10,000$ and $50,000$ variables. Figure 2 provides the proportion of explained variability as a function of the sparsity level.

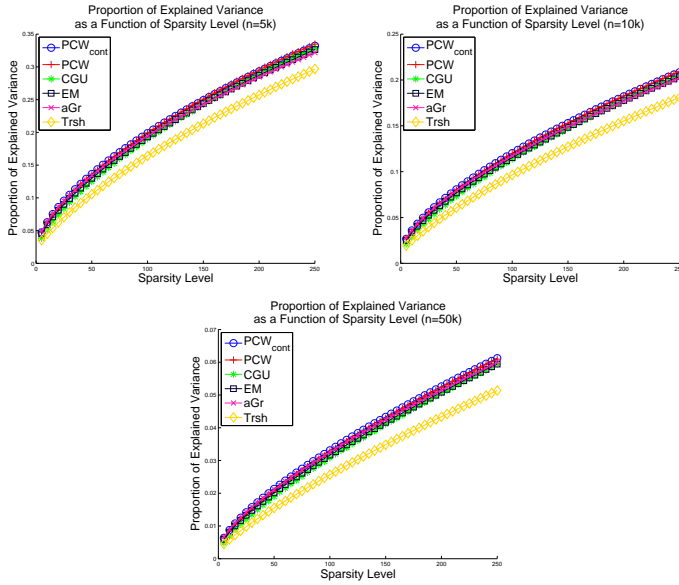


Fig. 2 PCW vs. Others - The proportion of explained variability as a function of the sparsity level for $n = 5000$, $10,000$ and $50,000$ are given in the upper left, upper right and bottom figures, respectively.

For small sparsity levels (< 50) most of the algorithms provide similar results, but as the sparsity level is increased, the CW algorithms becomes superior to all the other methods. This advantage is not achieved without a price. In Figure 3 we provide the cumulative computation time of the algorithms (the cumulative time is considered since the approximate greedy algorithm provides a full set of solutions).

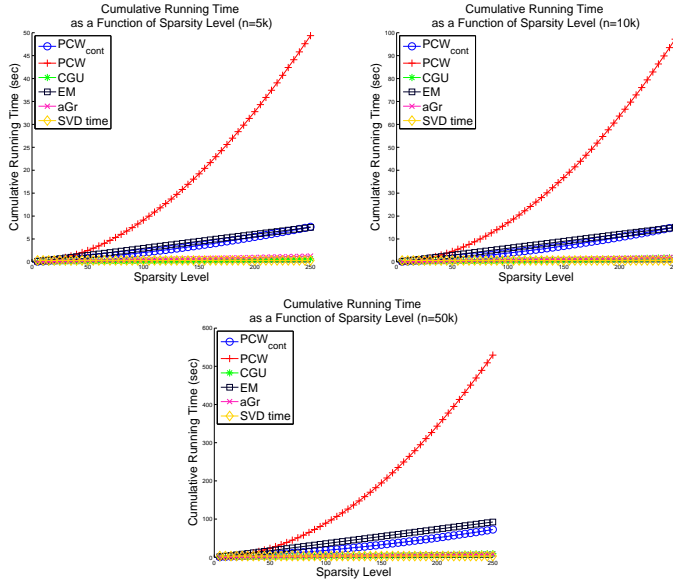


Fig. 3 PCW vs. Alternative Methods - The cumulative computation times as a function of the sparsity level for $n = 5,000, 10,000$ and $50,000$ are given in the upper left, upper right and bottom, respectively. The SVD time is the time required for computing the principal eigenvector of the covariance matrix that corresponds to the generated data, which is used in order to find the thresholding solution, and in order to initialize the CW and ConGradU algorithms.

Even though PCW has greatly decreased the computation time with respect to GCW, it still requires a notably higher amount of computation time with respect to the alternative algorithms. The scheme we referred as PCW_{cont} achieves similar results to PCW with respect to the function value. Regarding the running time, this scheme is competitive to the EM algorithm and requires somewhat more computational effort than the ConGradU and approximate greedy algorithms, thus providing a reasonable approach when a full set of solutions is required.

5.2 Gene Expression Dataset

Sparse PCA is extensively utilized in the identification of the genes that reflect the changes in the gene expression patterns during different biological states, thus contributing to the diagnosis and research of certain diseases such as cancer. Figure 4 illustrates the proportion of explained variability and the cumulative running time for a Leukemia data set [22]. This data set is composed from gene expression profiles of 72 patients with 12582 genes. The data set is normalized such that it has zero mean and unit variance.

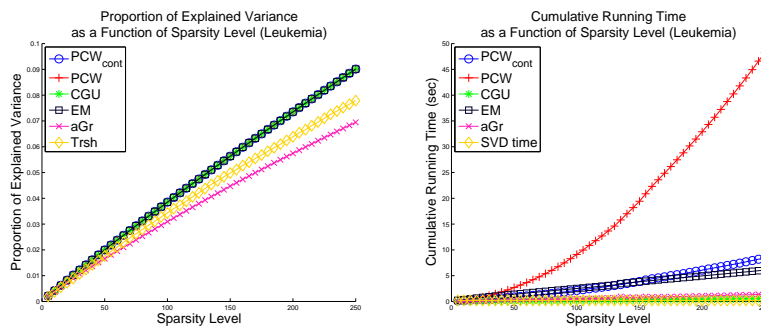


Fig. 4 Leukemia Gene Expression Data - The proportion of explained variability is given in the left figure and the cumulative computation time is given in the right one. The plot in both figures are given as a function of the sparsity level.

Most of the algorithms under consideration provide similar results with respect to the explained variability, which might indicate that this problem is, in a sense, rather easy to solve. We conducted similar experiments for additional 20 gene expression data sets from the GeneChip oncology database [23] that is publicly available in:

<http://compbio.dfci.harvard.edu/compbio/tools/gcod>

while commonly, all the algorithms provided similar results, we can still see in Figure 5 that PCW yields the best solution (with respect to the function value) more times than the alternative algorithms, and consequently it obtains the smallest mean error with respect to the best solution.

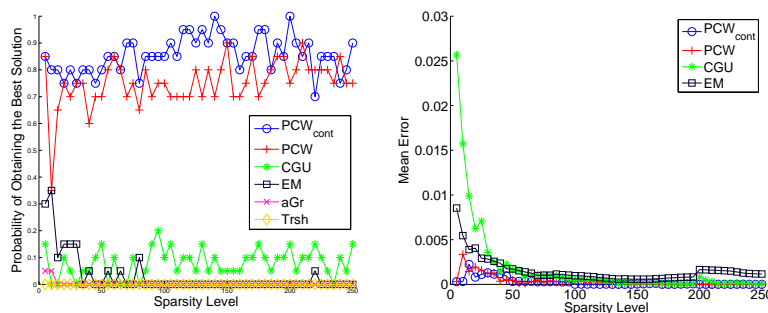


Fig. 5 Gene Expression Data - The left figure illustrates for each sparsity level the proportion of the number of data sets for which each algorithm obtained the best solution. The right figure illustrates for each sparsity level the mean error with respect to the best solution (the approximate greedy and thresholding algorithms were disregarded since both of them provide relative poor results).

6 Conclusions

In this paper, we considered the problem of maximizing a continuously differentiable convex function over the intersection of an l_2 unit ball and a sparsity constraint. We have shown that coordinate-wise maximality is a more restrictive condition than co-stationarity, which is the basis of many well-known methods for solving the sparse PCA problem. We introduced two algorithms (GCW and PCW) that are guaranteed to produce a CW-maximal solution, and demonstrated empirically the potential benefit of using this algorithms over some common algorithms proposed for this problem.

Acknowledgements We would like to thank two anonymous referees for their helpful remarks that helped to improve the presentation of the paper.

References

1. Jolliffe, I.T.: Principal Component Analysis, second edn. Springer, New York (2002)
2. Misra, J., Schmitt, W., Hwang, D., Hsiao, L.L., Gullans, S., Stephanopoulos, G., Stephanopoulos, G.: Interactive exploration of microarray gene expression patterns in a reduced dimensional space. *Genome research* **12**(7), 1112–1120 (2002)
3. d’Aspremont, A.: Identifying small mean-reverting portfolios. *Quant. Finance* **11**(3), 351–364 (2011)
4. Moghaddam, B., Weiss, Y., Avidan, S.: Spectral bounds for sparse pca: Exact and greedy algorithms. In: Y. Weiss, B. Schölkopf, J. Platt (eds.) *Adv. Neural. Inf. Process. Syst.* **18**, pp. 915–922. MIT Press, Cambridge, MA (2006)
5. Cadima, J., Jolliffe, I.T.: Loading and correlations in the interpretation of principle components. *J. Appl. Stat.* **22**(2), 203–214 (1995)
6. d’Aspremont, A., Bach, F., Ghaoui, L.E.: Optimal solutions for sparse principal component analysis. *J. Mach. Learn. Res.* **9**, 1269–1294 (2008)
7. d’Aspremont, A., El Ghaoui, L., Jordan, M., Lanckriet, G.: A direct formulation of sparse PCA using semidefinite programming. *SIAM Rev.* **49**(3) (2007)
8. Jolliffe, I.T., Trendafilov, N.T., Uddin, M.: A Modified Principal Component Technique Based on the LASSO. *J. Comput. Graph. Statist.* **12**(3), 531–547 (2003)
9. Tibshirani, R.: Regression shrinkage and selection via the lasso. *J. Roy. Statist. Soc. Ser. B* **58**, 267–288 (1996)
10. Witten, D.M., Hastie, T., Tibshirani, R.: A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis. *Biostatistics* **10**, 515–534 (2009)
11. Sigg, C.D., Buhmann, J.M.: Expectation-maximization for sparse and non-negative pca. In: *Proceedings of the 25th International Conference on Machine Learning, ICML ’08*, pp. 960–967. ACM, New York, NY, USA (2008)
12. Zou, H., Hastie, T., Tibshirani, R.: Sparse principal component analysis. *J. Comput. Graph. Statist.* **15**, 2006 (2004)
13. Zou, H., Hastie, T.: Regularization and variable selection via the elastic net. *J. Roy. Statist. Soc. Ser. B* **67**, 301–320 (2005)
14. Shen, H., Huang, J.Z.: Sparse principal component analysis via regularized low rank matrix approximation. *J. Multivariate Anal.* **99**(6), 1015 – 1034 (2008)
15. Journée, M., Nesterov, Y., Richtárik, P., Sepulchre, R.: Generalized power method for sparse principal component analysis. *J. Mach. Learn. Res.* **11**, 517–553 (2010)
16. Luss, R., Teboulle, M.: Conditional gradient algorithms for rank-one matrix approximations with a sparsity constraint. *SIAM Rev.* **55**(1), 65–98 (2013)
17. Sriperumbudur, B.K., Torres, D.A., Lanckriet, G.R.: A majorization-minimization approach to the sparse generalized eigenvalue problem. *Mach. Learn.* **85**(1), 3–39 (2011)

18. Beck, A., Eldar, Y.C.: Sparsity constrained nonlinear optimization: Optimality conditions and algorithms. *SIAM J. Optim.* **23**(3), 1480–1509 (2013)
19. Beck, A., Hallak, N.: On the minimization over sparse symmetric sets: Projections, optimality conditions, and algorithms. *Math. Oper. Res.* **41**(1), 196–223 (2016)
20. Rockafellar, R.: *Convex Analysis*. Princeton mathematical series. Princeton University Press (1970)
21. Jeffers, J.N.R.: Two case studies in the application of principal component analysis. *J. R. Stat. Soc. Ser. C. Appl. Stat.* **16**(3), pp. 225–236 (1967)
22. Armstrong, S.A., Staunton, J.E., Silverman, L.B., Pieters, R., den Boer, M.L., Minden, M.D., Sallan, S.E., Lander, E.S., Golub, T.R., Korsmeyer, S.J.: MLL translocations specify a distinct gene expression profile that distinguishes a unique leukemia. *Nature genetics* **30**(1), 41–47 (2002)
23. Liu, F., White, J., Antonescu, C., Gusenleitner, D., Quackenbush, J.: Gcod - genechip oncology database. *BMC Bioinformatics* **12**(1), 46 (2011)