

A majorization-minimization approach to the sparse generalized eigenvalue problem

Bharath K. Sriperumbudur · David A. Torres ·
Gert R.G. Lanckriet

Received: 1 March 2010 / Revised: 17 August 2010 / Accepted: 3 November 2010
© The Author(s) 2010. This article is published with open access at Springerlink.com

Abstract Generalized eigenvalue (GEV) problems have applications in many areas of science and engineering. For example, principal component analysis (PCA), canonical correlation analysis (CCA) and Fisher discriminant analysis (FDA) are specific instances of GEV problems, that are widely used in statistical data analysis. The main contribution of this work is to formulate a general, efficient algorithm to obtain sparse solutions to a GEV problem. Specific instances of sparse GEV problems can then be solved by specific instances of this algorithm. We achieve this by solving the GEV problem while constraining the cardinality of the solution. Instead of relaxing the cardinality constraint using a ℓ_1 -norm approximation, we consider a tighter approximation that is related to the negative log-likelihood of a Student's t-distribution. The problem is then framed as a d.c. (difference of convex functions) program and is solved as a sequence of convex programs by invoking the majorization-minimization method. The resulting algorithm is proved to exhibit *global convergence* behavior, i.e., for any random initialization, the sequence (subsequence) of iterates generated by the algorithm converges to a stationary point of the d.c. program. Finally, we illustrate the merits of this general sparse GEV algorithm with three specific examples of sparse GEV problems: sparse PCA, sparse CCA and sparse FDA. Empirical evidence for these examples suggests that the proposed sparse GEV algorithm, which offers a general framework

Editors: Süreyya Özgür-Akyüz, Devrim Ünay, and Alex Smola.

B.K. Sriperumbudur (✉) · G.R.G. Lanckriet
Department of Electrical and Computer Engineering, University of California, La Jolla, San Diego,
CA 92093-0407, USA
e-mail: bharathsv@ucsd.edu

G.R.G. Lanckriet
e-mail: gert@ece.ucsd.edu

D.A. Torres · G.R.G. Lanckriet
Department of Computer Science and Engineering, University of California, La Jolla, San Diego,
CA 92093-0407, USA

D.A. Torres
e-mail: datorres@cs.ucsd.edu

to solve any sparse GEV problem, will give rise to competitive algorithms for a variety of applications where specific instances of GEV problems arise.

Keywords Generalized eigenvalue problem · Principal component analysis · Canonical correlation analysis · Fisher discriminant analysis · Sparsity · D.c. program · Majorization-minimization · Zangwill's theory of global convergence · Music annotation · Cross-language document retrieval

1 Introduction

The generalized eigenvalue (GEV) problem for the matrix pair (A, B) is the problem of finding a pair (λ, x) such that

$$Ax = \lambda Bx, \quad (1)$$

where $A, B \in \mathbb{R}^{n \times n}$, $\mathbb{R}^n \ni x \neq 0$ and $\lambda \in \mathbb{R}$. When B is an identity matrix, the problem in (1) is simply referred to as an eigenvalue problem. Eigenvalue problems are so fundamental that they have applications in almost every area of science and engineering (Strang 1986).

In multivariate statistics, GEV problems are prominent and appear in problems dealing with high-dimensional data analysis, visualization and pattern recognition. In these applications, usually $x \in \mathbb{R}^n$, $A \in \mathbb{S}^n$ (the set of symmetric matrices of size $n \times n$ defined over \mathbb{R}) and $B \in \mathbb{S}_{++}^n$ (set of positive definite matrices of size $n \times n$ defined over \mathbb{R}). The variational formulation for the GEV problem in (1) is given by

$$\lambda_{\max}(A, B) = \max_x \{x^T A x : x^T B x = 1\}, \quad (\text{GEV-P})$$

where $\lambda_{\max}(A, B)$ is the maximum generalized eigenvalue associated with the matrix pair, (A, B) . The x that maximizes (GEV-P) is called the generalized eigenvector associated with $\lambda_{\max}(A, B)$. Some of the well-known and widely used data analysis techniques that are specific instances of (GEV-P) are:

- Principal component analysis (PCA) (Hotelling 1933; Jolliffe 1986), a classic tool for data analysis, data compression and visualization, finds the direction of maximal variance in a given multivariate data set. This technique is used in dimensionality reduction wherein the ambient space in which the data resides is approximated by a low-dimensional subspace without significant loss of information. The variational form of PCA is obtained by choosing A to be the covariance matrix (which is a positive semi-definite matrix defined over \mathbb{R}) associated with the multivariate data and B to be the identity matrix in (GEV-P).
- Canonical correlation analysis (CCA) (Hotelling 1936), similar to PCA, is also a data analysis and dimensionality reduction method. However, while PCA deals with only one data space \mathcal{X} (from which the multivariate data is obtained), CCA proposes a way for dimensionality reduction by taking into account relations between samples from two (or more) spaces \mathcal{X} and \mathcal{Y} . The assumption is that the data points from these two spaces contain some joint information that is reflected in correlations between them. Directions along which this correlation is high are thus assumed to be relevant directions when these relations are to be captured. The variational formulation for CCA is given by

$$\max_{w_x \neq 0, w_y \neq 0} \frac{w_x^T \Sigma_{xy} w_y}{\sqrt{w_x^T \Sigma_{xx} w_x} \sqrt{w_y^T \Sigma_{yy} w_y}}, \quad (2)$$

where \mathbf{w}_x and \mathbf{w}_y are the directions in \mathcal{X} and \mathcal{Y} along which the data is maximally correlated. Σ_{xx} and Σ_{yy} represent the covariance matrices for \mathcal{X} and \mathcal{Y} respectively and $\Sigma_{xy} = \Sigma_{yx}^T$ represents the cross-covariance matrix between \mathcal{X} and \mathcal{Y} . Equation (2) can be rewritten as

$$\max_{\mathbf{w}_x, \mathbf{w}_y} \{ \mathbf{w}_x^T \Sigma_{xy} \mathbf{w}_y : \mathbf{w}_x^T \Sigma_{xx} \mathbf{w}_x = 1, \mathbf{w}_y^T \Sigma_{yy} \mathbf{w}_y = 1 \}, \quad (3)$$

which in turn can be written in the form of (GEV-P) with:

$$A = \begin{pmatrix} \mathbf{0} & \Sigma_{xy} \\ \Sigma_{yx} & \mathbf{0} \end{pmatrix}, \quad B = \begin{pmatrix} \Sigma_{xx} & \mathbf{0} \\ \mathbf{0} & \Sigma_{yy} \end{pmatrix} \quad \text{and} \quad \mathbf{x} = \begin{pmatrix} \mathbf{w}_x \\ \mathbf{w}_y \end{pmatrix}.$$

- (c) In the binary classification setting, Fisher discriminant analysis (FDA) finds a one-dimensional subspace, $\mathbf{w} \in \mathbb{R}^n$, the projection of data onto which leads to maximal separation between the classes. Let μ_i and Σ_i denote the mean vector and covariance matrix associated with class i . The variational formulation of FDA is given by

$$\max_{\mathbf{w} \neq \mathbf{0}} \frac{(\mathbf{w}^T (\mu_1 - \mu_2))^2}{\mathbf{w}^T (\Sigma_1 + \Sigma_2) \mathbf{w}},$$

which can be rewritten as

$$\begin{aligned} \max_{\mathbf{w}} \quad & \mathbf{w}^T (\mu_1 - \mu_2) (\mu_1 - \mu_2)^T \mathbf{w} \\ \text{s.t.} \quad & \mathbf{w}^T (\Sigma_1 + \Sigma_2) \mathbf{w} = 1. \end{aligned} \quad (4)$$

Therefore, the FDA formulation is similar to (GEV-P) with $A = (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T$, called the *between-cluster variance* and $B = \Sigma_1 + \Sigma_2$, called the *within-cluster variance*. For multi-class problems, similar formulations lead to multiple-discriminant analysis.

Despite the simplicity and popularity of these data analysis and modeling methods, one potential limitation (depending on the application) is the lack of sparsity in their solution. Their solution vector, i.e., \mathbf{x} , is a linear combination of all input variables. Often, this makes it difficult to interpret the results and, therefore, limits their use for applications where interpretation of the solution is important. In the following, we motivate the need for sparse solutions to GEV problems by explaining how, for some specific instances of GEV problems, i.e., PCA, CCA and FDA, a variety of applications could benefit from sparsity for different reasons.

In many PCA applications, the coordinate axes have a physical interpretation; in biology, for example, each axis might correspond to a specific gene. In these cases, the interpretation of the principal components would be facilitated if they contained only few non-zero entries (or, loadings) while explaining most of the variance in the data. Moreover, in certain applications, e.g., financial asset trading strategies based on PCA techniques, the sparsity of the solution has important consequences, since fewer non-zero loadings imply fewer transaction costs. For CCA, consider a document translation application where two copies of a corpus of documents, one written in English and the other in German are given. The goal is to extract multiple low-dimensional representations of the documents, one in each language, each explaining most of the variation in the documents of a single language while maximizing the correlation between the representations to aid translation. Sparse representations, equivalent to representing the documents with a small set of words in each language, would allow to interpret the underlying translation mechanism and model it better. In music

annotation, CCA can be applied to model the correlation between semantic descriptions of songs (e.g., reviews) and their acoustic content. Sparsity in the semantic canonical components would allow to select the most meaningful words to describe musical content. This is expected to improve music annotation and retrieval systems. In a classification setting like FDA, feature selection aids generalization performance by promoting sparse solutions. To summarize, sparse representations are generally desirable as they aid human understanding, reduce computational and economic costs and promote better generalization.

In this paper, we propose an algorithm to find *sparse solutions to the generalized eigenvalue problem*. Specific instances of the sparse GEV problem, like, e.g., sparse CCA, can then be solved by specific instances of this algorithm. The sparse GEV problem can be written as

$$\max_{\mathbf{x}} \{ \mathbf{x}^T \mathbf{A} \mathbf{x} : \mathbf{x}^T \mathbf{B} \mathbf{x} = 1, \|\mathbf{x}\|_0 \leq k \}, \quad (\text{SGEV-P})$$

where $\mathbf{A} \in \mathbb{S}^n$, $\mathbf{B} \in \mathbb{S}_{++}^n$, $1 \leq k \leq n$ and $\|\mathbf{x}\|_0$ denotes the cardinality of \mathbf{x} , i.e., the number of non-zero elements of \mathbf{x} . Depending on the application and its sparsity requirements, k is chosen a priori. Equation (SGEV-P) can be solved either as a continuous optimization problem after relaxing the cardinality constraint or as a discrete optimization problem. In this paper, we follow the former approach, leading to our main contribution, a general algorithm for sparse GEV problems, which is presented in Sect. 3.

The first step in solving (SGEV-P) as a continuous optimization problem is to approximate the cardinality constraint. One usual heuristic is to approximate $\|\mathbf{x}\|_0$ by $\|\mathbf{x}\|_1$ (see Sect. 2 for the details on notation). Building on the earlier version of our work (Sriperumbudur et al. 2007), in Sect. 3.1, we approximate the cardinality constraint in (SGEV-P) as the negative log-likelihood of a Student's t-distribution, which is a tighter approximation than $\|\mathbf{x}\|_1$, and has been used earlier in many different contexts (Weston et al. 2003; Fazel et al. 2003; Candes et al. 2007). We then formulate this approximate problem as a d.c. (difference of convex functions) program and solve it in Sect. 3.3 as a sequence of quadratically constrained quadratic programs (QCQPs) (Boyd and Vandenberghe 2004), using the majorization-minimization (MM) method (Lange et al. 2000; Hunter and Lange 2004), which is briefly discussed in Sect. 3.2. As a special case, when $\mathbf{A} \in \mathbb{S}_+^n$ (the set of positive semidefinite matrices of size $n \times n$ defined over \mathbb{R}) and $\mathbf{B} \in \mathbb{S}_{++}^n$ is a diagonal matrix (as is the case, e.g., for PCA, where \mathbf{B} is an identity matrix), a simple iterative update rule can be obtained in closed form, which has a per iteration complexity of $O(n^2)$. Since the algorithm presented in this paper holds for any $\mathbf{A} \in \mathbb{S}^n$, it is more general than the ones in Sriperumbudur et al. (2007) and Moghaddam et al. (2007a), where $\mathbf{A} \in \mathbb{S}_+^n$. For example, the algorithm in Sriperumbudur et al. (2007) cannot handle the sparse CCA problem as $\mathbf{A} \in \mathbb{S}^n$ is indefinite. A preliminary version of the algorithm proposed in this paper, appeared in Torres et al. (2007b) and Sriperumbudur et al. (2009).

Since the proposed algorithm is an iterative procedure, using results from the global convergence theory of iterative algorithms (Zangwill 1969), we show in Sect. 3.4 that it is *globally convergent*. I.e., for any random initialization, the sequence (subsequence) of iterates generated by the algorithm converges to a stationary point of the d.c. program (see Sect. 3.4 for a detailed definition). A complete convergence analysis, including the rate of convergence and the quality of the obtained solution, is beyond the scope of this work and is the subject of future work.

In Sect. 4, we illustrate the merits of the proposed sparse GEV algorithm for some specific choices of \mathbf{A} and \mathbf{B} : algorithms for sparse PCA, sparse CCA and sparse FDA are obtained as special cases of the general sparse GEV algorithm. The empirical results observed for these special cases demonstrate the promise of the proposed algorithm, which is

more general than any of the special cases. More specifically, in Sect. 4.1, we compare the specific instance of our algorithm for $\mathbf{A} \in \mathbb{S}_+^n$ and \mathbf{B} the identity matrix (which we call DC-PCA) to other sparse PCA algorithms, SPCA (Zou et al. 2006), DSPCA (d’Aspremont et al. 2007), GSPCA (Moghaddam et al. 2007b) and GPower $_{\ell_0}$ (Journée et al. 2010), in terms of sparsity vs. explained variance on the “pit props” benchmark dataset, a random test dataset and three high-dimensional datasets, where the goal is to find relevant genes (as few as possible) while explaining the maximum possible variance. The results show that DC-PCA is competitive (with respect to sparsity, explained variance as well computational efficiency) with the state-of-the-art for sparse PCA. Given that sparse PCA is a well-studied problem, this indicates the potential of the general sparse GEV framework, proposed in this work, for instances of sparse GEV problems that have not been studied extensively. In Sect. 4.2, the proposed sparse GEV algorithm is used in two sparse CCA applications, one dealing with cross-language document retrieval and the other with vocabulary selection in music annotation. The cross-language document retrieval application involves a collection of documents with each document in different languages, say English and French. The goal is, given a query string in one language, retrieve the most relevant document(s) in the target language. We experimentally show that the proposed sparse CCA algorithm performs similar to the non-sparse version, however using only 10% of non-zero loadings in the canonical components. In the vocabulary selection application, we show that sparse CCA improves the performance of a statistical musical query system by selecting only those words (i.e., pruning the vocabulary) that are correlated to the underlying audio features. In Sect. 4.3, we consider the setting of sparse FDA. By exploiting the special structure of \mathbf{A} , i.e., $\mathbf{A} \in \mathbb{S}_+^n$ with $\text{rank}(\mathbf{A}) = 1$, we propose a sparse FDA algorithm that is more efficient than the general sparse GEV algorithm. We also discuss the relation of this algorithm to feature selection in least squares support vector machines (LS-SVMs) (Suykens et al. 2002, Chap. 3), which is a well-studied problem.

In summary, the main contribution of this paper is a globally convergent algorithm (see Algorithm 1)—a sequence of QCQPs—for the sparse GEV problem, which can handle any arbitrary symmetric matrix, \mathbf{A} and positive definite matrix, \mathbf{B} . When $\mathbf{A} \in \mathbb{S}_+^n$ and $\mathbf{B} \in \mathbb{S}_{++}^n$ is a diagonal matrix, the proposed algorithm has a simple, closed-form update rule. Encouraging empirical results obtained by applying the proposed, general algorithm to some special cases of sparse GEV problems suggest its potential for applications that rely on solving some instance of a sparse GEV problem.

In the following section, we introduce the notation that is used throughout the paper. Supplementary results and related discussions are collected in appendices.

2 Notation

\mathbb{S}^n (resp. \mathbb{S}_+^n , \mathbb{S}_{++}^n) denotes the set of symmetric (resp. positive semidefinite, positive definite) $n \times n$ matrices defined over \mathbb{R} . For $\mathbf{x} = (x_1, x_2, \dots, x_n)^T \in \mathbb{R}^n$, $\mathbf{x} \geq \mathbf{0}$ (resp. $\mathbf{x} > \mathbf{0}$) denotes $x_i \geq 0, \forall i$ (resp. $x_i > 0, \forall i$). $\mathbf{x} \geq \mathbf{y}$ (resp. $\mathbf{x} \leq \mathbf{y}$) denotes $x_i \geq y_i, \forall i$ (resp. $x_i \leq y_i, \forall i$). $\|\mathbf{x}\|_0$ denotes the number of non-zero elements of the vector \mathbf{x} , $\|\mathbf{x}\|_p := (\sum_{i=1}^n |x_i|^p)^{1/p}$, $1 \leq p < \infty$. \mathbf{I}_n denotes an $n \times n$ identity matrix. $\mathbf{D}(\mathbf{x})$ represents a diagonal matrix formed with \mathbf{x} as its principal diagonal.

3 Sparse generalized eigenvalue problem

In the following, we first discuss the intractability of (SGEV-P), which is then used in Sect. 3.1 to motivate and propose a d.c. (difference of convex functions) approximation

to (SGEV-P) that is based on a non-convex approximation to the cardinality constraint. This d.c. program is then solved as a sequence of QCQPs in Sect. 3.3 (see Algorithm 1) using the majorization-minimization method, which is briefly discussed in Sect. 3.2. Since the proposed algorithm is iterative, we present some convergence analysis in Sect. 3.4, which guarantees that the iterates generated by the algorithm converge to a stationary point of the d.c. program.

Intractability of (SGEV-P) Let us consider (SGEV-P), where $\mathbf{A} \in \mathbb{S}^n$ and $\mathbf{B} \in \mathbb{S}_{++}^n$. Suppose \mathbf{A} is not negative definite. Then (SGEV-P) is the maximization of a non-concave objective over the non-convex constraint set $\Phi := \{\mathbf{x} : \mathbf{x}^T \mathbf{B} \mathbf{x} = 1\} \cap \{\mathbf{x} : \|\mathbf{x}\|_0 \leq k\}$. Although Φ can be relaxed to a convex set $\tilde{\Phi} := \{\mathbf{x} : \mathbf{x}^T \mathbf{B} \mathbf{x} \leq 1\} \cap \{\mathbf{x} : \|\mathbf{x}\|_1 \leq k\}$, this does not simplify the problem as the maximization of a non-concave objective over a convex set is NP-hard [p. 342] (Rockafellar 1970).¹ Based on this discussion, it is evident that the intractability of (SGEV-P) is not only due to the cardinality constraint, but also due to the maximization of the non-concave objective function. Therefore, the ℓ_1 approximation that is popularly used in machine learning to obtain a convex approximation to an intractable problem involving a cardinality constraint, is not useful to obtain a tractable approximation to (SGEV-P).

In Appendix A, we discuss a convex approximation to (SGEV-P) using semidefinite programming (SDP) relaxation (Vandenberghe and Boyd 1996). We don't pursue such an approach since it is prohibitively expensive in computation for large n . Instead, in the following section, we trade convexity for good scalability and propose a non-convex approximation to (SGEV-P), resulting in a d.c. program, based on a non-convex approximation to the cardinality constraint.

3.1 Non-convex approximation to $\|\mathbf{x}\|_0$ and d.c. formulation

Let us consider the regularized (penalized) version of (SGEV-P) given by

$$\max_{\mathbf{x}} \left\{ \mathbf{x}^T \mathbf{A} \mathbf{x} - \rho \|\mathbf{x}\|_0 : \mathbf{x}^T \mathbf{B} \mathbf{x} \leq 1 \right\}, \quad (\text{SGEV-R})$$

where $\rho > 0$ is the regularization (penalization) parameter. Note that the quadratic equality constraint, $\mathbf{x}^T \mathbf{B} \mathbf{x} = 1$ is relaxed to the inequality constraint, $\mathbf{x}^T \mathbf{B} \mathbf{x} \leq 1$. Since

$$\|\mathbf{x}\|_0 = \sum_{i=1}^n \mathbb{1}_{\{|x_i| \neq 0\}} = \lim_{\varepsilon \rightarrow 0} \sum_{i=1}^n \frac{\log(1 + |x_i|/\varepsilon)}{\log(1 + 1/\varepsilon)}, \quad (5)$$

(SGEV-R) is equivalent (we define two programs to be *equivalent* if their optimizers, i.e., solutions are the same) to

$$\max_{\mathbf{x}} \left\{ \mathbf{x}^T \mathbf{A} \mathbf{x} - \rho \lim_{\varepsilon \rightarrow 0} \sum_{i=1}^n \frac{\log(1 + |x_i|/\varepsilon)}{\log(1 + 1/\varepsilon)} : \mathbf{x}^T \mathbf{B} \mathbf{x} \leq 1 \right\}. \quad (6)$$

¹Note that (GEV-P) also involves the maximization of a non-concave objective over a non-convex set. However, it is well-known that polynomial-time algorithms exist to solve (GEV-P), which is due to its *special* structure of a quadratic objective with a homogeneous quadratic constraint (Boyd and Vandenberghe 2004, p. 229).

The above program is approximated by the following *approximate sparse GEV program* by neglecting the limit in (6) and choosing $\varepsilon > 0$,

$$\max_{\mathbf{x}} \left\{ \mathbf{x}^T \mathbf{A} \mathbf{x} - \rho \sum_{i=1}^n \frac{\log(1 + |x_i|/\varepsilon)}{\log(1 + 1/\varepsilon)} : \mathbf{x}^T \mathbf{B} \mathbf{x} \leq 1 \right\}, \quad (7)$$

which is equivalent to

$$\max_{\mathbf{x}} \left\{ \mathbf{x}^T \mathbf{A} \mathbf{x} - \rho_\varepsilon \sum_{i=1}^n \log(|x_i| + \varepsilon) : \mathbf{x}^T \mathbf{B} \mathbf{x} \leq 1 \right\}, \quad (\text{SGEV-A})$$

where $\rho_\varepsilon := \rho / \log(1 + \varepsilon^{-1})$. Note that the approximate program in (SGEV-A) is a continuous optimization problem unlike the one in (SGEV-R), which has a combinatorial term. In addition, we show below that the objective of (SGEV-A) exhibits the d.c. structure and therefore can be written as a d.c. program. D.c. programs are well studied and many global optimization algorithms exist to solve them (Horst and Thoai 1999). Before we formulate (SGEV-A) as a d.c. program, we provide the definition of a d.c. program.

Definition 1 (D.c. program, Horst and Thoai 1999) Let Ω be a convex set in \mathbb{R}^n . A real valued function $f : \Omega \rightarrow \mathbb{R}$ is called a d.c. function on Ω , if there exist two *convex* functions $g, h : \Omega \rightarrow \mathbb{R}$ such that f can be expressed in the form $f(\mathbf{x}) = g(\mathbf{x}) - h(\mathbf{x})$, $\mathbf{x} \in \Omega$. Optimization problems of the form $\min\{f_0(\mathbf{x}) : \mathbf{x} \in \Omega, f_i(\mathbf{x}) \leq 0, i = 1, \dots, m\}$, where $f_i = g_i - h_i$, $i = 0, \dots, m$, are d.c. functions are called *d.c. programs*.

To formulate (SGEV-A) as a d.c. program, let us choose $\tau \in \mathbb{R}$ such that $\mathbf{A} + \tau \mathbf{I}_n \in \mathbb{S}_+^n$. If $\mathbf{A} \in \mathbb{S}_+^n$, such τ exists trivially (choose $\tau \geq 0$). If \mathbf{A} is indefinite, choosing $\tau \geq -\lambda_{\min}(\mathbf{A})$ ensures that $\mathbf{A} + \tau \mathbf{I}_n \in \mathbb{S}_+^n$. Therefore, choosing $\tau \geq \max(0, -\lambda_{\min}(\mathbf{A}))$ ensures that $\mathbf{A} + \tau \mathbf{I}_n \in \mathbb{S}_+^n$ for any $\mathbf{A} \in \mathbb{S}^n$. Equation (SGEV-A) is equivalently written as

$$\min_{\mathbf{x}} \left\{ \left[\tau \|\mathbf{x}\|_2^2 - \mathbf{x}^T (\mathbf{A} + \tau \mathbf{I}_n) \mathbf{x} \right] + \rho_\varepsilon \sum_{i=1}^n \log(|x_i| + \varepsilon) : \mathbf{x}^T \mathbf{B} \mathbf{x} \leq 1 \right\}. \quad (8)$$

Introducing the auxiliary variable, \mathbf{y} , yields the equivalent program

$$\min_{\mathbf{x}, \mathbf{y}} \left\{ \tau \|\mathbf{x}\|_2^2 - \left[\mathbf{x}^T (\mathbf{A} + \tau \mathbf{I}_n) \mathbf{x} - \rho_\varepsilon \sum_{i=1}^n \log(y_i + \varepsilon) \right] : \mathbf{x}^T \mathbf{B} \mathbf{x} \leq 1, -\mathbf{y} \preceq \mathbf{x} \preceq \mathbf{y} \right\}, \quad (9)$$

which is a d.c. program. Indeed, the term $\tau \|\mathbf{x}\|_2^2$ is convex in \mathbf{x} as $\tau \geq 0$ and, by construction, $\mathbf{x}^T (\mathbf{A} + \tau \mathbf{I}_n) \mathbf{x} - \rho_\varepsilon \sum_{i=1}^n \log(y_i + \varepsilon)$ is jointly convex in \mathbf{x} and \mathbf{y} . So, the above program is a minimization of the difference of two convex functions over a convex set.

As mentioned before, many global optimization algorithms exist to solve d.c. programs, e.g., branch and bound, cutting planes, etc. However, since these algorithms are not scalable for large n , in Sect. 3.3, we propose a local optimization algorithm based on majorization-minimization (see Sect. 3.2), which allows to exploit the d.c. structure and gives rise to a simple, efficient algorithm: a sequence of QCQPs. Before we discuss the MM method in Sect. 3.2 and present the sparse GEV algorithm in Sect. 3.3, we briefly discuss (a) the approximation to $\|\mathbf{x}\|_0$ that we consider in this paper and (b) the behavior of the solution to (SGEV-A)—or equivalently to (9)—as $\varepsilon \rightarrow 0$, assuming the global optimum to (SGEV-A) can be obtained.

Approximation to $\|\mathbf{x}\|_0$ The approximation (to $\|\mathbf{x}\|_0$) that we consider in this paper, i.e.,

$$\|\mathbf{x}\|_\varepsilon := \sum_{i=1}^n \frac{\log(1 + |x_i|\varepsilon^{-1})}{\log(1 + \varepsilon^{-1})},$$

has been used in many different contexts: feature selection using SVMs (Weston et al. 2003), sparse signal recovery (Candes et al. 2007), matrix rank minimization (Fazel et al. 2003), etc. This approximation is interesting because of its connection to sparse factorial priors that are studied in Bayesian inference, and can be interpreted as defining a Student's t-distribution prior over \mathbf{x} , an improper prior given by $\prod_{i=1}^n \frac{1}{|x_i| + \varepsilon}$. Tipping (2001) showed that this choice of prior leads to a sparse representation and demonstrated its validity for sparse kernel expansions in the Bayesian framework. Other approximations to $\|\mathbf{x}\|_0$ are possible, e.g., Bradley and Mangasarian (1998) used $\sum_{i=1}^n (1 - e^{-\alpha|x_i|})$ with $\alpha > 0$ ($\|\mathbf{x}\|_0 = \lim_{\alpha \rightarrow \infty} \sum_{i=1}^n (1 - e^{-\alpha|x_i|})$) as an approximation to $\|\mathbf{x}\|_0$ in the context of feature selection using SVMs.

In Appendix B, we show that the approximation (to $\|\mathbf{x}\|_0$) considered in this paper, i.e., $\|\mathbf{x}\|_\varepsilon$, is tighter than the ℓ_1 -norm approximation, for any $\varepsilon > 0$. Therefore, sparser solutions can be expected for (9)—obtained by replacing $\|\mathbf{x}\|_0$ in (SGEV-R) by $\|\mathbf{x}\|_\varepsilon$ —compared to replacing $\|\mathbf{x}\|_0$ by $\|\mathbf{x}\|_1$, in (SGEV-R), if the global solution(s) of the non-convex program in (9) can be found.

Behavior of the solution to (SGEV-A) as $\varepsilon \rightarrow 0$ As mentioned before, (SGEV-A) is an approximation to (SGEV-R), which is obtained by approximating $\|\mathbf{x}\|_0$ by $\|\mathbf{x}\|_\varepsilon$. Consider the objective functions of (SGEV-R) and (SGEV-A), given as

$$Q(\mathbf{x}) := \mathbf{x}^T \mathbf{A} \mathbf{x} - \rho \|\mathbf{x}\|_0,$$

$$Q_\varepsilon(\mathbf{x}) := \mathbf{x}^T \mathbf{A} \mathbf{x} - \rho \|\mathbf{x}\|_\varepsilon,$$

which are maximized over $\Omega := \{\mathbf{x} : \mathbf{x}^T \mathbf{B} \mathbf{x} \leq 1\}$. Suppose that $\hat{\mathbf{x}}$ and \mathbf{x}_ε are solutions to (SGEV-R) and (SGEV-A) respectively, i.e., maximizers of $Q(\mathbf{x})$ and $Q_\varepsilon(\mathbf{x})$ respectively over Ω . A natural question to answer is whether $\|\mathbf{x}_\varepsilon - \hat{\mathbf{x}}\| \rightarrow 0$ as $\varepsilon \rightarrow 0$? In general, this may not be true because $Q(\mathbf{x})$ may be quite flat near its maximum over Ω . If this is not true, at least, one would like to know whether $Q_\varepsilon(\mathbf{x}_\varepsilon) \rightarrow Q(\hat{\mathbf{x}})$ as $\varepsilon \rightarrow 0$, i.e.,

$$\lim_{\varepsilon \rightarrow 0} \max_{\mathbf{x} \in \Omega} Q_\varepsilon(\mathbf{x}) \stackrel{?}{=} \max_{\mathbf{x} \in \Omega} Q(\mathbf{x}) = \max_{\mathbf{x} \in \Omega} \lim_{\varepsilon \rightarrow 0} Q_\varepsilon(\mathbf{x}). \quad (10)$$

In other words, we would like to know whether the limit process and the maximization over Ω can be interchanged. It can be shown that if Q_ε converges uniformly over Ω to Q , then the equality in (10) holds. However, it is easy to see that Q_ε does not converge uniformly to Q over Ω , so nothing can be said about (10). Therefore, no theoretical guarantees can be provided about the behavior of $Q_\varepsilon(\mathbf{x}_\varepsilon)$ w.r.t. $Q(\hat{\mathbf{x}})$ as $\varepsilon \rightarrow 0$.

So, from a theoretical point of view, even if we could obtain a global optimum to (SGEV-A), it remains unclear how a global optimum of (SGEV-A) compares to that of (SGEV-R), as $\varepsilon \rightarrow 0$. From a practical perspective, as we pointed out, computing a global optimum to (SGEV-A) is computationally hard anyway. Therefore, instead of focusing on a global optimum for (SGEV-A), in Sect. 3.3, we propose a simple and efficient local optimization algorithm for (SGEV-A), based on the MM method. It remains an interesting open problem whether theoretical guarantees can be found, relating a global optimum of (SGEV-A) to that of (SGEV-R), which we will further investigate in future work.

3.2 Majorization-minimization method

The majorization-minimization (MM) method can be thought of as a generalization of the well-known expectation-maximization (EM) algorithm (Dempster et al. 1977). The general principle behind MM algorithms was first enunciated by the numerical analysts, Ortega and Rheinboldt (1970) in the context of line search methods. The MM principle appears in many places in statistical computation, including multidimensional scaling (de Leeuw 1977), robust regression (Huber 1981), correspondence analysis (Heiser 1987), variable selection (Hunter and Li 2005), signal/image processing (Daubechies et al. 2004; Bioucas-Dias et al. 2006; Figueiredo et al. 2007), etc. We refer the interested reader to a tutorial on MM algorithms (Hunter and Lange 2004) and the references therein.

The general idea of MM algorithms is as follows. Suppose we want to minimize f over $\Omega \subset \mathbb{R}^n$. The idea is to construct a *majorization function* g over $\Omega \times \Omega$ such that

$$f(x) \leq g(x, y), \quad \forall x, y \in \Omega \quad \text{and} \quad f(x) = g(x, x), \quad \forall x \in \Omega. \quad (11)$$

Thus, g as a function of x is an upper bound on f and coincides with f at y . The majorization-minimization algorithm corresponding to this majorization function g updates x at iteration l by

$$x^{(l+1)} \in \arg \min_{x \in \Omega} g(x, x^{(l)}), \quad (12)$$

unless we already have

$$x^{(l)} \in \arg \min_{x \in \Omega} g(x, x^{(l)}),$$

in which case the algorithm stops. $x^{(0)}$ is usually chosen randomly. The majorization function, g is usually constructed by using Jensen's inequality for convex functions, the first-order Taylor approximation or the quadratic upper bound principle (Böhning and Lindsay 1988). In fact, any other method can be used to construct g as long as it satisfies (11). It is easy to show that the above iterative scheme decreases the value of f monotonically in each iteration, i.e.,

$$f(x^{(l+1)}) \leq g(x^{(l+1)}, x^{(l)}) \leq g(x^{(l)}, x^{(l)}) = f(x^{(l)}), \quad (13)$$

where the first inequality and the last equality follow from (11) while the sandwiched inequality follows from (12).

Note that MM algorithms can be applied equally well to the maximization of f by simply reversing the inequality sign in (11) and changing the “min” to “max” in (12). In this case, the word MM refers to minorization-maximization, where the function g is called the *minorization function*. To put things in perspective, the EM algorithm can be obtained by constructing the minorization function g using Jensen's inequality for concave functions. The construction of such g is referred to as the E-step, while (12) with the “min” replaced by “max” is referred to as the M-step. The algorithm in (11) and (12) is used in machine learning, e.g., for non-negative matrix factorization (Lee and Seung 2001), under the name *auxiliary function method*. Lange et al. (2000) studied this algorithm under the name *optimization transfer* while Meng (2000) referred to it as the SM algorithm, where “S” stands for the surrogate step (same as the majorization/minorization step) and “M” stands for the minimization/maximization step depending on the problem at hand. g is called the surrogate function.

Having briefly discussed the idea behind MM algorithms, in the following example, we present an MM method to solve d.c. programs, which results in a sequence of convex programs. This will allow to derive an algorithm for the d.c. program in (9).

Example 1 (Linear majorization) Let us consider the optimization problem, $\min\{f(\mathbf{x}) : \mathbf{x} \in \Omega\}$ where $f = u - v$, with u and v both convex, and v continuously differentiable. Note that, by Definition 1, this is a d.c. program. Now, the goal is to construct an auxiliary function, g . Since v is convex, we have $v(\mathbf{x}) \geq v(\mathbf{y}) + (\mathbf{x} - \mathbf{y})^T \nabla v(\mathbf{y})$, $\forall \mathbf{x}, \mathbf{y} \in \Omega$. Therefore,

$$f(\mathbf{x}) = u(\mathbf{x}) - v(\mathbf{x}) \leq u(\mathbf{x}) - v(\mathbf{y}) - (\mathbf{x} - \mathbf{y})^T \nabla v(\mathbf{y}) =: g(\mathbf{x}, \mathbf{y}).$$

It is easy to verify that g is a majorization function of f . Therefore, by (12), we have

$$\mathbf{x}^{(l+1)} \in \arg \min_{\mathbf{x} \in \Omega} g(\mathbf{x}, \mathbf{x}^{(l)}) = \arg \min_{\mathbf{x} \in \Omega} u(\mathbf{x}) - \mathbf{x}^T \nabla v(\mathbf{x}^{(l)}). \quad (14)$$

If Ω is a convex set, then the above procedure solves a sequence of convex programs. Note that the same idea is used in the concave-convex procedure (CCCP) (Yuille and Rangarajan 2003).

In the following section, we apply the above example to (9) to obtain a sparse GEV algorithm, which is a sequence of QCQPs. We would like to mention that the d.c. program in (9) can be considered as an intermediate step in applying the MM method directly to (SGEV-A), to facilitate the construction of the majorization function. So, it is the general MM framework that essentially allows to derive a simple, efficient local optimization algorithm for the proposed sparse GEV problem formulation.

3.3 Sparse GEV algorithm

Let us return to the approximate sparse GEV program in (9). Let

$$f(\mathbf{x}, \mathbf{y}) = \tau \|\mathbf{x}\|_2^2 + \rho_\varepsilon \sum_{i=1}^n \log(\varepsilon + y_i) - \mathbf{x}^T (\mathbf{A} + \tau \mathbf{I}_n) \mathbf{x}, \quad (15)$$

where $\tau \geq \max(0, -\lambda_{\min}(\mathbf{A}))$ so that (9) can be written as $\min\{f(\mathbf{x}, \mathbf{y}) : (\mathbf{x}, \mathbf{y}) \in \Omega\}$ and $\Omega = \{(\mathbf{x}, \mathbf{y}) : \mathbf{x}^T \mathbf{B} \mathbf{x} \leq 1, -\mathbf{y} \leq \mathbf{x} \leq \mathbf{y}\}$. Defining

$$u(\mathbf{x}, \mathbf{y}) := \tau \|\mathbf{x}\|_2^2, \\ v(\mathbf{x}, \mathbf{y}) := \mathbf{x}^T (\mathbf{A} + \tau \mathbf{I}_n) \mathbf{x} - \rho_\varepsilon \sum_{i=1}^n \log(\varepsilon + y_i),$$

which are convex over $\mathbb{R}^n \times \mathbb{R}^n$ and invoking Example 1, yields the majorization function,

$$g((\mathbf{x}, \mathbf{y}), (\mathbf{z}, \mathbf{w})) := \tau \|\mathbf{x}\|_2^2 - \mathbf{z}^T (\mathbf{A} + \tau \mathbf{I}_n) \mathbf{z} + \rho_\varepsilon \sum_{i=1}^n \log(\varepsilon + w_i) \\ - 2(\mathbf{x} - \mathbf{z})^T (\mathbf{A} + \tau \mathbf{I}_n) \mathbf{z} + \rho_\varepsilon \sum_{i=1}^n \frac{y_i - w_i}{w_i + \varepsilon}, \quad (16)$$

and therefore the following algorithm based on (14):

$$(\mathbf{x}^{(l+1)}, \mathbf{y}^{(l+1)}) = \arg \min_{\mathbf{x}, \mathbf{y}} \tau \|\mathbf{x}\|_2^2 - 2\mathbf{x}^T (\mathbf{A} + \tau \mathbf{I}_n) \mathbf{x}^{(l)} + \rho_\varepsilon \sum_{i=1}^n \frac{y_i}{y_i^{(l)} + \varepsilon} \\ \text{s.t.} \quad \mathbf{x}^T \mathbf{B} \mathbf{x} \leq 1, \quad -\mathbf{y} \leq \mathbf{x} \leq \mathbf{y}. \quad (17)$$

Equation (17) is equivalent to the following *sparse GEV algorithm*,

$$\mathbf{x}^{(l+1)} = \arg \min_{\mathbf{x}} \left\{ \tau \|\mathbf{x}\|_2^2 - 2\mathbf{x}^T (\mathbf{A} + \tau \mathbf{I}_n) \mathbf{x}^{(l)} + \rho_\varepsilon \sum_{i=1}^n \frac{|x_i|}{|x_i^{(l)}| + \varepsilon} : \mathbf{x}^T \mathbf{B} \mathbf{x} \leq 1 \right\}, \quad (\text{ALG})$$

which is a sequence of quadratically constrained quadratic programs (QCQPs).² We refer to Boyd and Vandenberghe (2004) for details about QCQPs. For $\tau \neq 0$, the optimal solution of (ALG), $\mathbf{x}^{(l+1)}$, is unique, since the objective function in (ALG) is strictly convex in \mathbf{x} when $\tau \neq 0$. On the other hand, for $\tau = 0$, the objective function is no longer strictly convex. Below, in Special case 1, we show that for $\tau = 0$, (ALG) has a unique optimal solution, which can be found by solving a quadratic program rather than a QCQP.

Note that (ALG) requires the knowledge of τ , which depends on $\lambda_{\min}(\mathbf{A})$. Although for certain choices of \mathbf{A} , $\lambda_{\min}(\mathbf{A})$ can be computed efficiently, it is generally an expensive task, which can significantly raise the cost of solving (ALG), if the computation of $\lambda_{\min}(\mathbf{A})$ is included as part of solving (ALG). However, if $\mathbf{A} \in \mathbb{S}_+^n$ (e.g., for PCA and FDA), τ can be chosen to be zero, to not raise the cost of solving (ALG).

Special case 1 ($\mathbf{A} \in \mathbb{S}_+^n$, $\mathbf{B} \in \mathbb{S}_+^n$, $\tau = 0$) Define $w_i^{(l)} := (|x_i^{(l)}| + \varepsilon)^{-1}$, $\mathbf{w}^{(l)} := (w_1^{(l)}, \dots, w_n^{(l)})$ and $\mathbf{W}^{(l)} := \mathbf{D}(\mathbf{w}^{(l)})$, a diagonal matrix with $\mathbf{w}^{(l)}$ as its principal diagonal. Choosing $\tau = 0$ in (ALG), we have that

$$\mathbf{x}^{(l+1)} = \arg \max_{\mathbf{x}} \left\{ \mathbf{x}^T \mathbf{A} \mathbf{x}^{(l)} - \frac{\rho_\varepsilon}{2} \|\mathbf{W}^{(l)} \mathbf{x}\|_1 : \mathbf{x}^T \mathbf{B} \mathbf{x} \leq 1 \right\}.$$

We show in Appendix C that when $\rho_\varepsilon > 2 \max_i |(\mathbf{A} \mathbf{x}^{(l)})_i| (w_i^{(l)})^{-1}$, the unique optimal solution of (ALG) is $\mathbf{x}^{(l+1)} = \mathbf{0}$. When $\rho_\varepsilon < 2 \max_i |(\mathbf{A} \mathbf{x}^{(l)})_i| (w_i^{(l)})^{-1}$, we show that

$$\mathbf{x}^{(l+1)} = \frac{\mathbf{S}(\mathbf{SBS})^+(\boldsymbol{\gamma} + \boldsymbol{\lambda}^{(l+1)})}{\sqrt{(\boldsymbol{\gamma} + \boldsymbol{\lambda}^{(l+1)})^T (\mathbf{SBS})^+ (\boldsymbol{\gamma} + \boldsymbol{\lambda}^{(l+1)})}}, \quad (\text{ALG-R})$$

where \mathbf{C}^+ denotes the Moore-Penrose pseudoinverse of \mathbf{C} , $\mathbf{S} := \mathbf{D}(\mathbf{s})$, $(\mathbf{s})_i := \text{sign}((\mathbf{A} \mathbf{x}^{(l)})_i)$, $(\boldsymbol{\gamma})_i := |(\mathbf{A} \mathbf{x}^{(l)})_i| - \frac{\rho_\varepsilon}{2} w_i^{(l)}$, and

$$\boldsymbol{\lambda}^{(l+1)} \in \arg \min_{\boldsymbol{\lambda}} \{ (\boldsymbol{\gamma} + \boldsymbol{\lambda})^T (\mathbf{SBS})^+ (\boldsymbol{\gamma} + \boldsymbol{\lambda}) : \boldsymbol{\lambda} \geq \mathbf{0} \}. \quad (18)$$

The solution to (18) is unique for $(\boldsymbol{\lambda}^{(l+1)})_i$, $i \in \mathcal{I} := \{j : (\mathbf{s})_j \neq 0\}$, while it is not unique for $(\boldsymbol{\lambda}^{(l+1)})_i$, $i \notin \mathcal{I}$. However, because of the pre-multiplication by \mathbf{S} in (ALG-R), $(\mathbf{x}^{(l+1)})_i = 0$, $i \notin \mathcal{I}$ and, therefore, $\mathbf{x}^{(l+1)}$ in (ALG-R) is the unique optimal solution of (ALG). Note that in this case of $\tau = 0$, we just need to solve a quadratic program in (18) compared to solving a QCQP in (ALG) when $\tau \neq 0$.

²Although (ALG) is not a QCQP in standard form, its equivalent program in (17) can be easily written in standard form as $\mathbf{z}^{(l+1)} = \arg \min_{\mathbf{z}} \{ \mathbf{z}^T \mathbf{C} \mathbf{z} + \mathbf{b}^T \mathbf{z} : \mathbf{z}^T \mathbf{E} \mathbf{z} \leq 1, \mathbf{z}^T \mathbf{G} \mathbf{z} + \mathbf{c}_i^T \mathbf{z} \leq 0, \mathbf{z}^T \mathbf{G} \mathbf{z} + \mathbf{d}_i^T \mathbf{z} \leq 0, 1 \leq i \leq n \}$, where $\mathbf{z} := [\mathbf{x}^T \ \mathbf{y}^T]^T$, $\mathbf{b} = [-2(\mathbf{A} + \tau \mathbf{I}_n) \mathbf{x}^{(l)}]$, $\mathbf{C} = [\tau \mathbf{I}_n \ \mathbf{0}]$, $\mathbf{E} = [\mathbf{B} \ \mathbf{0}]$, $\mathbf{F} = \mathbf{D}(\mathbf{y}^{(l)} + \varepsilon)$ and \mathbf{G} is a $2n \times 2n$ matrix of zeroes. \mathbf{c}_i and \mathbf{d}_i are defined as: for $1 \leq i \leq n$, $(\mathbf{c}_i)_j = \delta_{ij}$, $1 \leq j \leq n$ and $(\mathbf{c}_i)_j = -\delta_{(n+i)j}$, $(n+1) \leq j \leq 2n$; for $1 \leq i \leq n$, $(\mathbf{d}_i)_j = -\delta_{ij}$, $1 \leq j \leq n$ and $(\mathbf{d}_i)_j = -\delta_{(n+i)j}$, $(n+1) \leq j \leq 2n$, where δ represents the Kronecker delta. Therefore, for simplicity, we refer to (ALG) as a QCQP.

In case $\rho_\varepsilon = 2 \max_i |(\mathbf{A}\mathbf{x}^{(l)})_i| (w_i^{(l)})^{-1}$ (rather hypothetical since ρ_ε is a positive *real* number chosen by the user), $\mathbf{x}^{(l+1)}$ is no longer unique. However, $\|\mathbf{x}^{(l+1)}\|_0$ is usually very small in this case since $(\mathbf{x}^{(l+1)})_j = 0$ for $j \notin \arg \max_i |(\mathbf{A}\mathbf{x}^{(l)})_i| (w_i^{(l)})^{-1}$, as discussed in Appendix C. Also, $\mathbf{0}$ is a solution of (ALG) in this case. Therefore, we set $\mathbf{x}^{(l+1)} = \mathbf{0}$ as the optimal solution of (ALG) when $\rho_\varepsilon \geq 2 \max_i |(\mathbf{A}\mathbf{x}^{(l)})_i| (w_i^{(l)})^{-1}$, to implement practical algorithms.

Special case 2 ($\mathbf{A} \in \mathbb{S}_+^n$, $\mathbf{B} \in \mathbb{S}_{++}^n$ diagonal, $\tau = 0$) Let \mathbf{B}_{++}^n be a diagonal matrix, i.e., $\mathbf{B} = \mathbf{D}(\mathbf{b})$, where $\mathbf{b} = (b_1, \dots, b_n) > \mathbf{0}$. Then, it can be shown that (ALG-R) reduces to a simple update rule, given by the closed form expression

$$x_i^{(l+1)} = \frac{[(\mathbf{A}\mathbf{x}^{(l)})_i] - \frac{\rho_\varepsilon}{2} w_i^{(l)}]_+ \text{sign}((\mathbf{A}\mathbf{x}^{(l)})_i)}{b_i \sqrt{\sum_{i=1}^n b_i^{-1} [(\mathbf{A}\mathbf{x}^{(l)})_i] - \frac{\rho_\varepsilon}{2} w_i^{(l)}]^2_+}}, \quad \forall i, \quad (\text{ALG-S})$$

where $[a]_+ := \max(0, a)$. See Appendix C for the derivation of (ALG-S). Note that (ALG-S) has a per-iteration complexity of $O(n^2)$ compared to the worst-case complexity of $O(n^3)$ for (ALG) and (ALG-R).

Interpretation of (ALG) Assuming $\tau \neq 0$, (ALG) reduces to

$$\mathbf{x}^{(l+1)} = \arg \min_{\mathbf{x}} \left\{ \|\mathbf{x} - (\tau^{-1} \mathbf{A} + \mathbf{I}_n) \mathbf{x}^{(l)}\|_2^2 + \frac{\rho_\varepsilon}{\tau} \|\mathbf{W}^{(l)} \mathbf{x}\|_1 : \mathbf{x}^T \mathbf{B} \mathbf{x} \leq 1 \right\}. \quad (19)$$

Equation (19) is very similar to LASSO (Tibshirani 1996) except for the *weighted* ℓ_1 -penalty and the quadratic constraint. When $\mathbf{x}^{(0)}$ is chosen such that $\mathbf{x}^{(0)} = \mathbf{a}\mathbf{1}$, then the first iteration of (19) is a LASSO minimization problem except for the quadratic constraint. Let us analyze (19) to get an intuitive interpretation.

(a) $\rho_\varepsilon = \rho = 0$: (19) reduces to

$$\min_{\mathbf{x}} \left\{ \|\mathbf{x} - \mathbf{t}^{(l)}\|_2^2 : \mathbf{x}^T \mathbf{B} \mathbf{x} \leq 1 \right\},$$

where $\mathbf{t}^{(l)} = (\tau^{-1} \mathbf{A} + \mathbf{I}_n) \mathbf{x}^{(l)}$, i.e., the first term in the objective of (19) computes the best approximation to $\mathbf{t}^{(l)}$ in the ℓ_2 -norm so that the approximation lies in the ellipsoid $\mathbf{x}^T \mathbf{B} \mathbf{x} \leq 1$, and therefore, the solution \mathbf{x} is non-sparse.

(b) $\rho_\varepsilon = \rho = \infty$: In this case, (19) reduces to

$$\min_{\mathbf{x}} \left\{ \|\mathbf{W}^{(l)} \mathbf{x}\|_1 : \mathbf{x}^T \mathbf{B} \mathbf{x} \leq 1 \right\},$$

which is a weighted ℓ_1 -norm minimization problem. Intuitively, it is clear that if $x_i^{(l)}$ is small, its weighting factor, $w_i^{(l)} = (|x_i^{(l)}| + \varepsilon)^{-1}$ in the next minimization step is large, which therefore pushes $x_i^{(l+1)}$ to be small. This way the small entries in \mathbf{x} are generally pushed toward zero as far as the constraints on \mathbf{x} allow, therefore yielding a sparse solution.

From the above discussion, it is clear that (19) is a trade-off between the solution to the least-squares problem and the solution to the weighted ℓ_1 -norm problem. From now on, we refer to (ALG) as the *Sparse GEV algorithm*, which is detailed in Algorithm 1.

Algorithm 1 Sparse generalized eigenvalue algorithm

Require: $A \in \mathbb{S}^n$, $B \in \mathbb{S}_{++}^n$, $\varepsilon > 0$ and $\rho > 0$

1: Choose $\tau \geq \max(0, -\lambda_{\min}(A))$

2: Choose $\mathbf{x}^{(0)} \in \{\mathbf{x} : \mathbf{x}^T B \mathbf{x} \leq 1\}$

3: $\rho_\varepsilon = \frac{\rho}{\log(1+\varepsilon^{-1})}$

4: **if** $\tau = 0$ **then**

5: **if** $B = D(\mathbf{b})$ **then**

6: **repeat**

7: $w_i^{(l)} = (|x_i^{(l)}| + \varepsilon)^{-1}$

8: **if** $\rho_\varepsilon < 2 \max_i |(A\mathbf{x}^{(l)})_i| (w_i^{(l)})^{-1}$ **then**

9:

$$x_i^{(l+1)} = \frac{[|(A\mathbf{x}^{(l)})_i| - \frac{\rho_\varepsilon}{2} w_i^{(l)}]_+ \text{sign}((A\mathbf{x}^{(l)})_i)}{b_i \sqrt{\sum_{i=1}^n b_i^{-1} [|(A\mathbf{x}^{(l)})_i| - \frac{\rho_\varepsilon}{2} w_i^{(l)}]_+^2}}, \quad \forall i$$

10: **else**

11:

$$\mathbf{x}^{(l+1)} = \mathbf{0}$$

12: **end if**

13: **until** convergence

14: **else**

15: **repeat**

16: $w_i^{(l)} = (|x_i^{(l)}| + \varepsilon)^{-1}$

17: **if** $\rho_\varepsilon < 2 \max_i |(A\mathbf{x}^{(l)})_i| (w_i^{(l)})^{-1}$ **then**

18: $(\boldsymbol{\gamma})_i = |(A\mathbf{x}^{(l)})_i| - \frac{\rho_\varepsilon}{2} w_i^{(l)}$

19: $(\mathbf{s})_i = \text{sign}((A\mathbf{x}^{(l)})_i)$, $S = D(\mathbf{s})$

20: $\boldsymbol{\lambda}^{(l+1)} \in \arg \min_{\boldsymbol{\lambda}} \{(\boldsymbol{\gamma} + \boldsymbol{\lambda})^T (SBS)^+ (\boldsymbol{\gamma} + \boldsymbol{\lambda}) : \boldsymbol{\lambda} \geq \mathbf{0}\}$

21:

$$\mathbf{x}^{(l+1)} = \frac{S(SBS)^+ (\boldsymbol{\gamma} + \boldsymbol{\lambda}^{(l+1)})}{\sqrt{(\boldsymbol{\gamma} + \boldsymbol{\lambda}^{(l+1)})^T (SBS)^+ (\boldsymbol{\gamma} + \boldsymbol{\lambda}^{(l+1)})}}$$

22: **else**

23:

$$\mathbf{x}^{(l+1)} = \mathbf{0}$$

24: **end if**

25: **until** convergence

26: **end if**

27: **else**

28: **repeat**

29: $w_i^{(l)} = (|x_i^{(l)}| + \varepsilon)^{-1}$

30: $\mathbf{W}^{(l)} = D(\mathbf{w}^{(l)})$

31:

$$\mathbf{x}^{(l+1)} = \arg \min_{\mathbf{x}} \left\{ \|\mathbf{x} - (\tau^{-1} A + I_n) \mathbf{x}^{(l)}\|_2^2 + \frac{\rho_\varepsilon}{\tau} \|\mathbf{W}^{(l)} \mathbf{x}\|_1 : \mathbf{x}^T B \mathbf{x} \leq 1 \right\}$$

32: **until** convergence

33: **end if**

34: **return** $\mathbf{x}^{(l)}$

Choice of ρ , τ and ε To run Algorithm 1, ρ , τ and ε need to be chosen. Based on the above discussion, although it is clear that ρ controls the sparsity of the solution output by Algorithm 1, it is not possible to know a priori the value of ρ that will provide some desired sparsity level. Therefore, in practice, Algorithm 1 is solved for a fixed set of ρ values that are chosen a priori and the solution with desired cardinality is selected. Depending on the specific sparse GEV under consideration, ρ may sometimes be selected in a different way. For example, for supervised learning with FDA, the algorithm can be cross-validated over the set of ρ values that was fixed a priori and the ρ that provides the smallest cross-validation error may be selected. Since ρ is a free parameter, τ and ε can be set to any value (that satisfies the constraints in Algorithm 1) and ρ can be tuned to obtain the desired sparsity as mentioned above. However, it has to be noted that for a fixed value of ρ , increasing τ or ε reduces sparsity.³ So, in practice τ is chosen to be $\max(0, -\lambda_{\min}(\mathbf{A}))$, ε to be close to zero and ρ is set by searching for an appropriate value as discussed above.

Post-processing Suppose that Algorithm 1 outputs a solution, \mathbf{x}^* such that $\|\mathbf{x}^*\|_0 = k$. Can we say that \mathbf{x}^* is the optimal solution of (SGEV-P) among all \mathbf{x} with cardinality k ? The following proposition provides a condition to check for the non-optimality of \mathbf{x}^* . In addition, it also presents a post-processing step (called *variational renormalization*) that improves the performance of Algorithm 1. See Moghaddam et al. (2007b, Proposition 2) for a similar result in the case of $\mathbf{A} \in \mathbb{S}_+^n$ and $\mathbf{B} = \mathbf{I}_n$.

Proposition 1 Suppose Algorithm 1 converges to a solution \mathbf{x}^* such that $\|\mathbf{x}^*\|_0 = k$. Let \mathbf{z} be the sub-vector of \mathbf{x}^* obtained by removing the zero entries of \mathbf{x}^* and

$$\mathbf{u}_k = \arg \max_{\mathbf{x}} \{\mathbf{x}^T \mathbf{A}_k \mathbf{x} : \mathbf{x}^T \mathbf{B}_k \mathbf{x} = 1\},$$

where \mathbf{A}_k and \mathbf{B}_k are submatrices of \mathbf{A} and \mathbf{B} defined by the same non-zero indices of \mathbf{x}^* . If $\mathbf{z} \neq \mathbf{u}_k$, then \mathbf{x}^* is not the optimal solution of (SGEV-P) among all \mathbf{x} with the same sparsity pattern as \mathbf{x}^* (and therefore, is not the optimal solution of (SGEV-P) among all \mathbf{x} with $\|\mathbf{x}\|_0 = k$). Nevertheless, by replacing the non-zero entries of \mathbf{x}^* with those of \mathbf{u}_k , the value of the objective function in (SGEV-P) increases from $[\mathbf{x}^*]^T \mathbf{A} \mathbf{x}^*$ to $\lambda(\mathbf{A}_k, \mathbf{B}_k)$, its optimal value among all \mathbf{x} with the same sparsity pattern as \mathbf{x}^* .

Proof Assume that \mathbf{x}^* , the solution output by Algorithm 1, is the optimal solution of (SGEV-P). Define \mathbf{v} such that $v_i = \mathbb{1}_{\{|\mathbf{x}_i^*| \neq 0\}}$. Since \mathbf{x}^* is the optimal solution of (SGEV-P), we have

$$\mathbf{x}^* = \arg \max_{\mathbf{y}} \{ \mathbf{y}^T \mathbf{D}(\mathbf{v}) \mathbf{A} \mathbf{D}(\mathbf{v}) \mathbf{y} : \mathbf{y}^T \mathbf{D}(\mathbf{v}) \mathbf{B} \mathbf{D}(\mathbf{v}) \mathbf{y} = 1 \},$$

which is equivalent to $\mathbf{z} = \arg \max_{\mathbf{w}} \{ \mathbf{w}^T \mathbf{A}_k \mathbf{w} : \mathbf{w}^T \mathbf{B}_k \mathbf{w} = 1 \} = \mathbf{u}_k$ and the result follows. Note that $\lambda(\mathbf{A}_k, \mathbf{B}_k)$ is the optimal value of (SGEV-P) among all \mathbf{x} with the same sparsity pattern as \mathbf{x}^* . Therefore, if $\mathbf{z} = \mathbf{u}_k$, then $[\mathbf{x}^*]^T \mathbf{A} \mathbf{x}^* = \lambda(\mathbf{A}_k, \mathbf{B}_k)$. \square

The variational renormalization suggests that given a solution (in our case, \mathbf{x}^* at the termination of Algorithm 1), it is almost certainly better to discard the loadings, keep only the

³Increasing ε increases the approximation error between $\|\mathbf{x}\|_0$ and $\|\mathbf{x}\|_\varepsilon$ and therefore reduces sparsity. From (19), it is clear that increasing τ reduces the weight on the term $\|\mathbf{W}^{(l)} \mathbf{x}\|_1$, which means more importance is given to reducing the approximation error, $\|\mathbf{x} - (\tau^{-1} \mathbf{A} + \mathbf{I}_n) \mathbf{x}^{(l)}\|_2^2$, leading to a less sparse solution.

sparsity pattern and solve the smaller unconstrained subproblem to obtain the final loadings, given the sparsity pattern. This procedure surely improves any continuous algorithm's performance.

In Algorithm 1, we mention that the iterative scheme is continued until convergence. What does convergence mean here? Does the algorithm really converge? If it converges, what does it converge to? Does it converge to an optimal solution of (SGEV-A)? To address these questions, in the following section, we provide some convergence analysis of Algorithm 1, using tools from global convergence theory (Zangwill 1969).

3.4 Convergence analysis

For an iterative procedure like Algorithm 1 to be useful, it must converge to point solutions from all or at least a significant number of initialization states and not exhibit other nonlinear system behaviors, such as divergence or oscillation. *Global convergence theory of iterative algorithms* (Zangwill 1969) can be used to investigate this behavior. We mention up front that this *does not* deal with proving convergence to a global optimum. Using this theory, recently, Sriperumbudur and Lanckriet (2009) analyzed the convergence behavior of the iterative linear majorization algorithm in (14) and showed that under certain conditions on u and v , the algorithm in (14) is globally convergent, i.e., for any random initialization, $\mathbf{x}^{(0)}$, the sequence of iterates $\{\mathbf{x}^{(l)}\}_{l=0}^{\infty}$ converges to some stationary point— \mathbf{x}_* is said to be stationary point of a constrained optimization problem if it satisfies the corresponding Karush-Kuhn-Tucker (KKT) conditions (Bonnans et al. 2006, Sect. 13.3), which assuming constraint qualification are necessary for the local optimality of \mathbf{x}_* —of the d.c. program, $\min\{u(\mathbf{x}) - v(\mathbf{x}) : \mathbf{x} \in \Omega\}$. Since (ALG) is obtained by applying linear majorization to (9), as shown in Sect. 3.3, the convergence analysis of (ALG) can be carried out by invoking the results in Sriperumbudur and Lanckriet (2009). This results in Theorem 1, which states that Algorithm 1 is globally convergent. This expresses, in a sense, the certainty that the algorithm works. The rate of convergence remains an open problem and is the subject of future work. It is important to stress the fact that global convergence does not imply (contrary to what the term might suggest) convergence to a global optimum for all initial points $\mathbf{x}^{(0)}$. Based on Theorem 1, we also present the convergence analysis for some special cases of Algorithm 1 in Corollaries 1–3. The proofs of all these results are provided in Appendix D.

Theorem 1 (Global convergence of sparse GEV algorithm) *Let $\{\mathbf{x}^{(l)}\}_{l=0}^{\infty}$ be any sequence generated by the sparse GEV algorithm in Algorithm 1. Then, all the limit points of $\{\mathbf{x}^{(l)}\}_{l=0}^{\infty}$ are fixed points of Algorithm 1, which are stationary points of the program in (SGEV-A),*

$$\rho_{\varepsilon} \sum_{i=1}^n \log(\varepsilon + |x_i^{(l)}|) - [\mathbf{x}^{(l)}]^T \mathbf{A} \mathbf{x}^{(l)} \rightarrow \rho_{\varepsilon} \sum_{i=1}^n \log(\varepsilon + |x_i^*|) - [\mathbf{x}^*]^T \mathbf{A} \mathbf{x}^* := L^*, \quad (20)$$

for some fixed point \mathbf{x}^* , $\|\mathbf{x}^{(l+1)} - \mathbf{x}^{(l)}\| \rightarrow 0$, and either $\{\mathbf{x}^{(l)}\}_{l=0}^{\infty}$ converges or the set of limit points of $\{\mathbf{x}^{(l)}\}_{l=0}^{\infty}$ is a connected and compact subset of $\mathcal{S}(L^*)$, where $\mathcal{S}(a) := \{\mathbf{x} \in \mathcal{S} : \mathbf{x}^T \mathbf{A} \mathbf{x} - \rho_{\varepsilon} \sum_{i=1}^n \log(\varepsilon + |x_i|) = -a\}$ and \mathcal{S} is the set of fixed points of Algorithm 1. If $\mathcal{S}(L^*)$ is finite, then any sequence $\{\mathbf{x}^{(l)}\}_{l=0}^{\infty}$ generated by Algorithm 1 converges to some \mathbf{x}^* in $\mathcal{S}(L^*)$.

Having considered the convergence of Algorithm 1, we now consider the convergence of some of its special cases. The following result shows that a simple iterative algorithm can be obtained to compute the generalized eigenvector associated with $\lambda_{\max}(\mathbf{A}, \mathbf{B})$, when $\mathbf{A} \in \mathcal{S}_+^n$.

Corollary 1 Let $\mathbf{A} \in \mathbb{S}_+^n$, $\tau = 0$ and $\rho = 0$. Then, any sequence $\{\mathbf{x}^{(l)}\}_{l=0}^\infty$ generated by the following algorithm

$$\mathbf{x}^{(l+1)} = \frac{\mathbf{B}^{-1} \mathbf{A} \mathbf{x}^{(l)}}{\sqrt{[\mathbf{x}^{(l)}]^T \mathbf{A} \mathbf{B}^{-1} \mathbf{A} \mathbf{x}^{(l)}}}, \quad (21)$$

converges to some \mathbf{x}^* such that $\lambda_{\max}(\mathbf{A}, \mathbf{B}) = [\mathbf{x}^*]^T \mathbf{A} \mathbf{x}^*$ and $[\mathbf{x}^*]^T \mathbf{B} \mathbf{x}^* = 1$.

Assuming $\mathbf{B} = \mathbf{I}_n$, $\rho = 0$ and $\tau = 0$, as a corollary to Corollary 1, the following result shows that Algorithm 1 reduces to the power method for computing $\lambda_{\max}(\mathbf{A})$, and the sequence of iterates generated by Algorithm 1 converges to the eigenvector associated with $\lambda_{\max}(\mathbf{A})$.

Corollary 2 (Power method) Let $\mathbf{A} \in \mathbb{S}_+^n$, $\mathbf{B} = \mathbf{I}_n$, $\tau = 0$ and $\rho = 0$. Then Algorithm 1 is the power method for computing $\lambda_{\max}(\mathbf{A})$, wherein any sequence $\{\mathbf{x}^{(l)}\}_{l=0}^\infty$ generated by Algorithm 1 converges to some $\mathbf{x}^* \in \{\mathbf{x} : \|\mathbf{x}\|_2^2 = 1\}$ such that $\lambda_{\max}(\mathbf{A}) = [\mathbf{x}^*]^T \mathbf{A} \mathbf{x}^*$.

Finally, Corollary 3 generalizes the above result for any $\mathbf{A} \in \mathbb{S}^n$ with $\lambda_{\max}(\mathbf{A}) > 0$, i.e., the sequence of iterates generated by Algorithm 1 converges to the eigenvector associated with $\lambda_{\max}(\mathbf{A})$ for any $\mathbf{A} \in \mathbb{S}^n$ with $\lambda_{\max}(\mathbf{A}) > 0$, and not just for $\mathbf{A} \in \mathbb{S}_+^n$.

Corollary 3 Let $\mathbf{A} \in \mathbb{S}^n$ such that $\lambda_{\max}(\mathbf{A}) > 0$. Assume $\mathbf{B} = \mathbf{I}_n$ and $\rho = 0$. Then, any sequence $\{\mathbf{x}^{(l)}\}_{l=0}^\infty$ generated by Algorithm 1 converges to some $\mathbf{x}^* \in \{\mathbf{x} : \|\mathbf{x}\|_2^2 = 1\}$ such that $\lambda_{\max}(\mathbf{A}) = [\mathbf{x}^*]^T \mathbf{A} \mathbf{x}^*$.

In summary, in this section, we have presented our main contribution: a globally convergent sparse GEV algorithm (Algorithm 1) which can handle any $\mathbf{A} \in \mathbb{S}^n$ and $\mathbf{B} \in \mathbb{S}_{++}^n$. In the following section, we illustrate the potential of this general algorithm for some specific instances of sparse GEV problems, corresponding to specific choices for \mathbf{A} and \mathbf{B} : sparse PCA, sparse CCA and sparse FDA.

4 Experimental results

To illustrate the potential of the sparse GEV algorithm (Algorithm 1), suited for any application that involves a specific instance of a sparse GEV problem, we present empirical results for some of its special cases in this section. To this end, in Sect. 4.1, we choose $\mathbf{A} \in \mathbb{S}_+^n$, $\mathbf{B} = \mathbf{I}_n$ and $\tau = 0$ —a special case of (ALG-S)—to obtain a sparse PCA algorithm, called DC-PCA, which is then compared to various other sparse PCA algorithms (that are proposed in literature) on both benchmark and three high-dimensional gene datasets. Experiments show that the performance of DC-PCA in terms of explained variance vs. cardinality, as well as its computational efficiency compare well to the state-of-the-art amongst these other algorithms. In Sect. 4.2, we choose $\mathbf{A} \in \mathbb{S}^n$ (with \mathbf{A} being indefinite as shown beneath (3)) and $\mathbf{B} \in \mathbb{S}_{++}^n$ to obtain a sparse CCA algorithm, which is evaluated for two sparse CCA applications: cross-language document retrieval and vocabulary selection in music annotation. In Sect. 4.3, we choose $\mathbf{A} \in \mathbb{S}_+^n$ with $\text{rank}(\mathbf{A}) = 1$ and $\mathbf{B} \in \mathbb{S}_{++}^n$, which is the setting of FDA (see (4)). Exploiting the special structure of \mathbf{A} for FDA allows to propose a sparse FDA algorithm that is more efficient than the general sparse GEV algorithm in Algorithm 1. We also discuss its relation to feature selection in least squares SVMs (Suykens et al. 2002, Chap. 3).

Algorithm 2 Sparse PCA algorithm (DC-PCA)**Require:** $\mathbf{A} \in \mathbb{S}_+^n$, $\varepsilon > 0$ and $\rho > 0$ 1: Choose $\mathbf{x}^{(0)} \in \{\mathbf{x} : \mathbf{x}^T \mathbf{x} \leq 1\}$ 2: $\rho_\varepsilon = \frac{\rho}{\log(1+\varepsilon^{-1})}$ 3: **repeat**4: $w_i^{(l)} = (|x_i^{(l)}| + \varepsilon)^{-1}$ 5: **if** $\rho_\varepsilon < 2 \max_i |(\mathbf{A}\mathbf{x}^{(l)})_i| (w_i^{(l)})^{-1}$ **then**

6:

$$x_i^{(l+1)} = \frac{[|(\mathbf{A}\mathbf{x}^{(l)})_i| - \frac{\rho_\varepsilon}{2} w_i^{(l)}]_+ \text{sign}((\mathbf{A}\mathbf{x}^{(l)})_i)}{\sqrt{\sum_{i=1}^n [|(\mathbf{A}\mathbf{x}^{(l)})_i| - \frac{\rho_\varepsilon}{2} w_i^{(l)}]_+^2}}, \quad \forall i$$

7: **else**

8:

$$\mathbf{x}^{(l+1)} = \mathbf{0}$$

9: **end if**10: **until** convergence11: **return** $\mathbf{x}^{(l)}$

4.1 Sparse principal component analysis

In this section, we consider sparse PCA as a special case of the sparse GEV algorithm, presented in Sect. 3.3, for $\mathbf{A} \in \mathbb{S}_+^n$ being a covariance matrix, $\mathbf{B} = \mathbf{I}_n$ and $\tau = 0$, and call the resulting algorithm DC-PCA (see Algorithm 2). Note that this is a special case of (ALG-S) with $b_i = 1$, $\forall i$. Since the computation of $\mathbf{x}^{(l+1)}$, from $\mathbf{x}^{(l)}$, involves computing $\mathbf{A}\mathbf{x}^{(l)}$, which has a complexity of $O(n^2)$, the DC-PCA algorithm has a per iteration complexity of $O(n^2)$. Being a special case of Algorithm 1, which is a globally convergent algorithm, DC-PCA is also globally convergent.

Before empirically comparing the performance of DC-PCA to other sparse PCA algorithms that have been proposed in literature, we briefly discuss the prior work on sparse PCA algorithms. The earliest attempts at “sparsifying” PCA consisted of simple axis rotations and component thresholding (Cadima and Jolliffe 1995) for subset selection, often based on the identification of principal variables (McCabe 1984). The first true computational technique, called SCoTLASS (Jolliffe et al. 2003), provided an optimization framework using LASSO (Tibshirani 1996) by enforcing a sparsity constraint on the PCA solution by bounding its ℓ_1 -norm, leading to a non-convex procedure. Zou et al. (2006) proposed a ℓ_1 -penalized regression algorithm for PCA (called SPCA) using an *elastic net* (Zou and Hastie 2005) and solved it efficiently using least angle regression (Efron et al. 2004). Subsequently, d’Aspremont et al. (2007) proposed a convex relaxation to the non-convex cardinality constraint for PCA (called DSPCA) leading to a SDP with a complexity of $O(n^4 \sqrt{\log n})$. Although this method shows performance comparable to SPCA on a small-scale benchmark data set, it is not scalable for high-dimensional data sets, even possibly with Nesterov’s first-order method (Nesterov 2005). Moghaddam et al. (2007b) proposed a combinatorial optimization algorithm (called GSPCA) using greedy search and branch-and-bound methods to solve the sparse PCA problem, leading to a total complexity of $O(n^4)$ for a full set of solutions (one for each target sparsity between 1 and n). d’Aspremont et al. (2008) formulated a new SDP relaxation to the sparse PCA problem and derived a more efficient greedy algorithm (compared to GSPCA) for computing a full set of solutions at a total numerical

complexity of $O(n^3)$, which is based on the convexity of the largest eigenvalue of a symmetric matrix. Recently, Journée et al. (2010) proposed a simple, iterative sparse PCA algorithm (GPower $_{\ell_0}$) with a per iteration complexity of $O(n^2)$ and showed that it performs similar or better than many of the above mentioned algorithms.

We now illustrate the effectiveness of DC-PCA in terms of sparsity and scalability on various datasets by comparing to different approaches. On small datasets, the performance of DC-PCA is compared against SPCA, DSPCA, GSPCA and GPower $_{\ell_0}$, while on large datasets, DC-PCA is compared to all these algorithms except DSPCA and GSPCA due to scalability issues. Since GPower $_{\ell_0}$ has been compared to the greedy algorithm of d'Aspremont et al. (2008) by Journée et al. (2010), wherein it is shown that these two algorithms perform similarly except for the greedy algorithm being computationally more complex, we do not include the greedy algorithm in our comparison. The results show that the performance of DC-PCA is comparable to the performance of many of these algorithms, with state-of-the-art *scalability*, i.e., comparable to GPower $_{\ell_0}$. This competitive comparison to the state-of-the-art for the special case of sparse PCA, which is a well-studied problem, illustrates the potential of the general sparse GEV algorithm, proposed in this work, for other instances of sparse GEV problems that have not been studied extensively. The experiments in this paper are carried out on a Linux 3 GHz, 4 GB RAM workstation. On the implementation side, we fix ε to be the machine precision in all our experiments, which is motivated from the discussion in Sect. 3.1.

4.1.1 Pit props data

The pit props dataset (Jeffers 1967) has become a standard benchmark example to test sparse PCA algorithms. The first 6 principal components (PCs) capture 87% of the total variance. Therefore, the explanatory power of sparse PCA methods is often compared on the first 6 sparse PCs.⁴ Table 1 shows the first 3 sparse PCs and their loadings for SPCA, DSPCA, GSPCA, GPower $_{\ell_0}$ and DC-PCA. Using the first 6 sparse PCs, SPCA captures 75.8% of the variance with a cardinality pattern of (7, 4, 4, 1, 1, 1), which indicates the number of non-zero loadings for the first to the sixth sparse PC, respectively. This results in a total of 18 non-zero loadings for SPCA, while DSPCA captures 75.5% of the variance with a sparsity pattern of (6, 2, 3, 1, 1, 1), totaling 14 non-zero loadings. With a sparsity pattern of (6, 2, 2, 1, 1, 1) (total of *only* 13 non-zero loadings), DC-PCA, GSPCA and GPower $_{\ell_0}$ can capture 77.1% of the total variance. Comparing the cumulative variance and cumulative cardinality, Figs. 1(a–b) show that DC-PCA explains more variance with fewer non-zero loadings than SPCA and DSPCA. In addition, its performance is similar to that of GSPCA and GPower $_{\ell_0}$. For the first sparse PC, Fig. 1(c) shows that DC-PCA consistently explains more variance with better sparsity than SPCA, while performing similar to other algorithms. Figure 1(d) shows the variation of sparsity and explained variance with respect to ρ for the first sparse PC computed with DC-PCA. This plot summarizes the method for setting ρ : the algorithm is run for various ρ and the value of ρ that achieves the desired sparsity is selected.

⁴The discussion so far dealt with computing the first sparse eigenvector (which we also call the “first sparse PC”) of \mathbf{A} . To compute the second sparse eigenvector (or “second sparse PC”), the matrix \mathbf{A} is deflated with the first sparse eigenvector (see Mackey 2009 for details) and the sparse PCA algorithm is applied again. In general, subsequent sparse eigenvectors (or “sparse PCs”) are obtained by applying the sparse PCA algorithm to a sequence of deflated matrices. In this paper, we use the orthogonalized Hotelling’s deflation, as mentioned in Mackey (2009). More details can also be found in Sriperumbudur et al. (2007).

Table 1 Loadings for first three sparse principal components (PCs) of the pit props data. The SPCA and DSPCA loadings are taken from Zou et al. (2006) and d'Aspremont et al. (2007) respectively

PC	x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9	x_{10}	x_{11}	x_{12}	x_{13}
SPCA	1	-.48	-.48	0	.18	0	-.25	-.34	-.42	-.40	0	0	0
	2	0	0	.79	.62	0	0	-.02	0	0	0	.01	0
	3	0	0	0	0	.64	.49	0	0	0	0	0	-.02
DSPCA	1	-.56	-.58	0	0	0	-.26	-.10	-.37	-.36	0	0	0
	2	0	0	.71	.71	0	0	0	0	0	0	0	0
	3	0	0	0	0	-.79	-.61	0	0	0	0	0	.01
GSPCA	1	.44	.45	0	0	0	.38	.34	.40	.42	0	0	0
	2	0	0	.71	.71	0	0	0	0	0	0	0	0
	3	0	0	0	0	.82	.58	0	0	0	0	0	0
GPower _{t₀}	1	.44	.45	0	0	0	.38	.34	.40	.42	0	0	0
	2	0	0	.71	.71	0	0	0	0	0	0	0	0
	3	0	0	0	0	.82	.58	0	0	0	0	0	0
DC-PCA	1	.45	.46	0	0	0	.37	.33	.40	.42	0	0	0
	2	0	0	.71	.71	0	0	0	0	0	0	0	0
	3	0	0	0	0	.82	.58	0	0	0	0	0	0

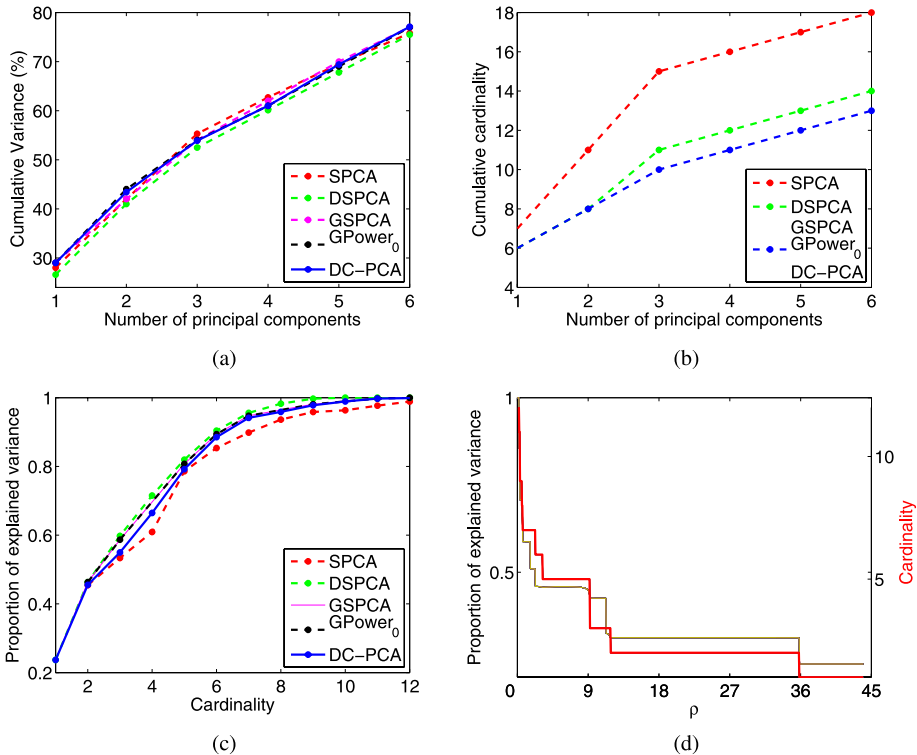


Fig. 1 (Color online) Pit props: (a) cumulative variance and (b) cumulative cardinality for the first 6 sparse principal components (PCs); (c) proportion of explained variance (PEV) vs. cardinality for the first sparse PC (obtained by varying the sparsity parameter and computing the cardinality and explained variance for the solution vector); (d) dependence of sparsity (shown in red) and PEV (shown in brown) on ρ for the first sparse PC computed with DC-PCA

4.1.2 Random test problems

In this section, we follow the experimental setup that is considered in Journée et al. (2010). Throughout this section, we assume $A = C^T C$, where C is a $p \times n$ random matrix whose entries are generated according to a Gaussian distribution, with zero mean and unit variance. In the following, we present the trade-off curves (proportion of explained variance vs. cardinality for the first sparse PC associated with A), computational complexity vs. cardinality and computational complexity vs. problem size for various sparse PCA algorithms.

Trade-off curves Figure 2(a) shows the trade-off between the proportion of explained variance and the cardinality for the first sparse PC associated with A for various sparse PCA algorithms. For each algorithm, the sparsity inducing parameter (k in the case of DSPCA and GSPCA, and the regularization parameter in the case of SPCA, GPower_{ℓ_0} and DC-PCA) is incrementally increased to obtain the solution vector with cardinality that decreases from n to 1. The results displayed in Fig. 2(a) are averages of computations on 100 random matrices with dimensions $p = 100$ and $n = 300$. It can be seen from Fig. 2(a) that DC-PCA performs similar to DSPCA, GSPCA and GPower_{ℓ_0} , while performing better than SPCA.

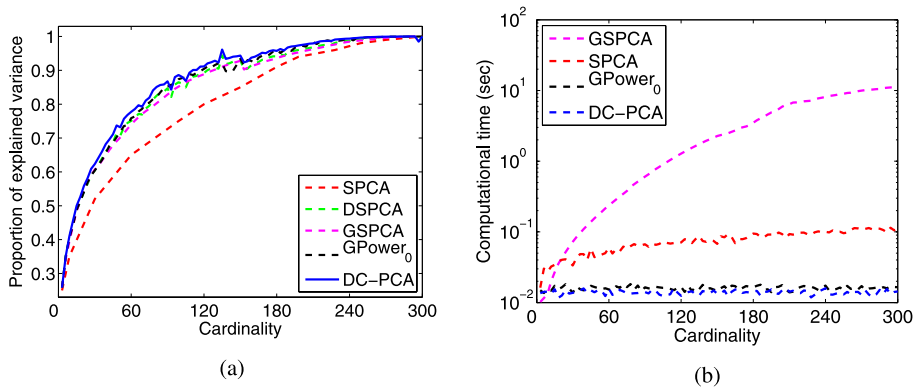


Fig. 2 Random test data: (a) (average) proportion of explained variance vs. cardinality for the first sparse PC of A ; (b) (average) computation time vs. cardinality. In (a), all the sparse PCA algorithms perform similarly and better than SPCA. In (b), the complexity of GSPCA grows significantly with increasing cardinality of the solution vector, while the speed of the other methods is almost independent of the cardinality

Table 2 Average computation time (in seconds) for the first sparse PC associated with A for a fixed regularization parameter

$p \times n$	100×1000	250×2500	500×5000	750×7500	1000×10000
SPCA	0.135	1.895	10.256	34.367	87.459
GPower _{ℓ₀}	0.027	0.159	0.310	1.224	1.904
DC-PCA	0.034	0.151	0.301	1.202	1.913

Computational complexity vs. cardinality Figure 2(b) shows the average time required by the sparse PCA algorithms to extract the first sparse PC of A with $p = 100$ and $n = 300$, for varying cardinality. It is obvious from Fig. 2(b) that as the cardinality increases, GSPCA tends to get slower while the speed of SPCA, GPower_{ℓ₀} and DC-PCA is not really affected by the cardinality. We did not show the results of DSPCA in Fig. 2(b) as its computational complexity is an order of magnitude (around 100 times) more than the scale on the vertical axis of Fig. 2(b). Journée et al. (2010) have demonstrated that the greedy method proposed by d’Aspremont et al. (2008) also exhibits the behavior of increasing computational complexity with the increase in cardinality.

Computational complexity vs. problem size Figure 3 shows the average computation time in seconds, required by various sparse PCA algorithms, to extract the first sparse PC of A , for various problem sizes, n , where n is increased exponentially and p is fixed to 500. The times shown are averages over 100 random instances of A for each problem size, where the sparsity inducing parameters are chosen such that the solution vectors of these algorithms exhibit comparable cardinality. It is clear from Fig. 3 that DC-PCA and GPower_{ℓ₀} scale better to large-dimensional problems than the other algorithms. Since, on average, GSPCA and DSPCA are much slower than the other methods, even for low cardinalities (see Fig. 2(b)), we discard them from all the following numerical experiments that deal with large n .

For the remaining algorithms, SPCA, GPower_{ℓ₀} and DC-PCA, we run another round of experiments, now examining the computational complexity with varying n and p but with a fixed aspect ratio $n/p = 10$. The results are depicted in Table 2. Again, the corresponding

Fig. 3 Average computation time (seconds) for the first sparse PC of A vs. problem size, n , over 100 randomly generated matrices A

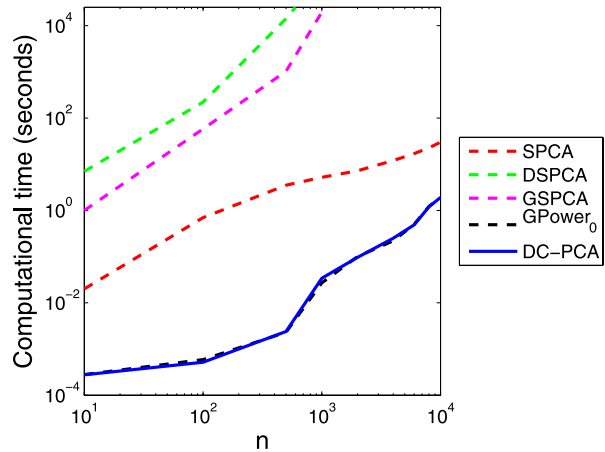


Table 3 Gene expression datasets

Dataset	Samples (p)	Genes (n)	Reference
Colon cancer	62	2000	Alon et al. (1999)
Leukemia	38	7129	Golub et al. (1999)
Ramaswamy	127	16063	Ramaswamy et al. (2001)

regularization parameters are set in such a way that the solution vectors of these algorithms exhibit comparable cardinality. The values displayed in Table 2 correspond to the average running times of the algorithms on 100 random instances of A for each problem size. It can be seen that our proposed method, DC-PCA, is comparable to GPower_{ℓ_0} and faster than SPCA.

4.1.3 Gene expression data

Gene expression data from DNA microarrays provides the expression level of thousands of genes across several hundreds or thousands of experiments. To enhance the interpretation of these large data sets, sparse PCA algorithms can be applied, to extract sparse principal components that involve only a few genes.

Datasets Usually, gene expression data is specified by a $p \times n$ matrix (say C) of p samples and n genes. The covariance matrix, A is therefore computed as $C^T C$. In our experiments, we consider three gene expression datasets which are tabulated in Table 3. The colon cancer dataset (Alon et al. 1999) consists of 62 tissue samples (22 normal and 40 cancerous) with the gene expression profiles of $n = 2000$ genes extracted from DNA microarray data. Its first principal component explains 44.96% of the total variance. The leukemia dataset (Golub et al. 1999) consists of a training set of 38 samples (27 ALL and 11 AML, two variants of leukemia) from bone marrow specimens and a test set of 34 samples (20 ALL and 14 AML). This dataset has been used widely in classification settings where the goal is to distinguish between two variants of leukemia. All samples have 7129 features, each of which corresponds to a normalized gene expression value extracted from the microarray image. The

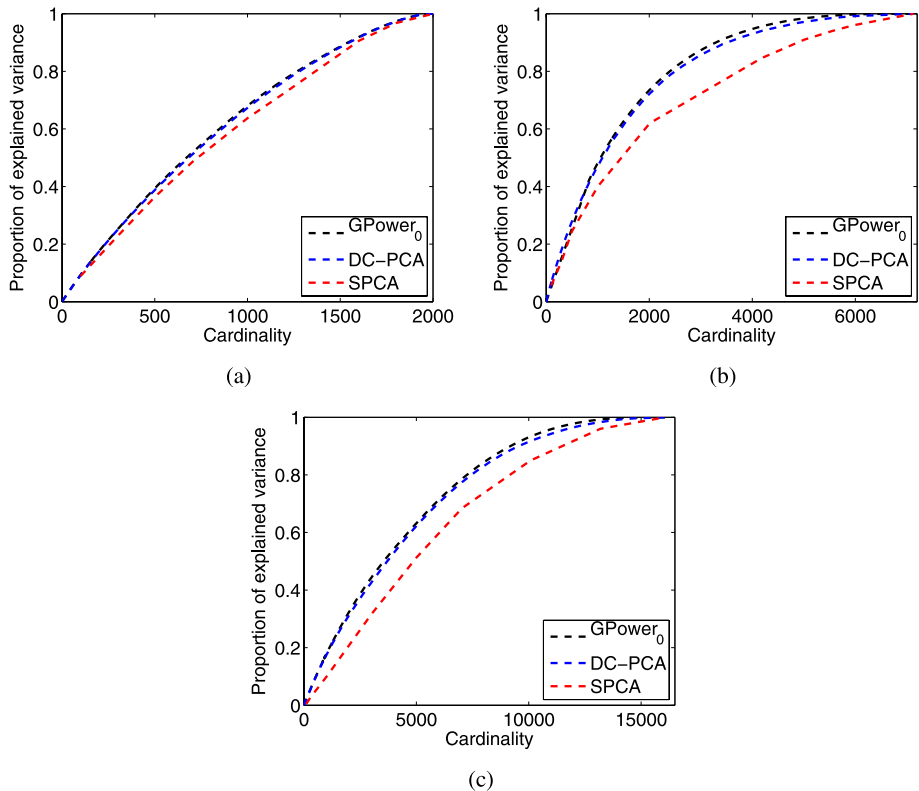


Fig. 4 Trade-off curves between explained variance and cardinality for (a) colon cancer, (b) leukemia and (c) Ramaswamy datasets. The proportion of variance explained is computed on the first sparse principal component. (a–c) show that DC-PCA performs similar to GPower_{ℓ_0} , while explaining more variance (for a fixed cardinality) than SPCA

first principal component explains 87.64% of the total variance. The Ramaswamy dataset has 16063 genes and 127 samples, its first principal component explaining 76.5% of the total variance.

The high dimensionality of these datasets makes them suitable candidates for studying the performance of sparse PCA algorithms, by investigating their ability to explain variance in the data based on a small number of genes, to obtain interpretable results. Since DSPCA and GSPCA are not scalable for these large datasets, in our study, we compare DC-PCA to SPCA and GPower_{ℓ_0} .

Trade-off curves Figures 4(a–c) show the proportion of explained variance versus the cardinality for the first sparse PC for the datasets shown in Table 3. It can be seen that DC-PCA performs similar to GPower_{ℓ_0} and performs better than SPCA.

Computational complexity The average computation time required by the sparse PCA algorithms on each dataset is shown in Table 4. The indicated times are averages over n computations, one for each cardinality ranging from n down to 1. The results show that DC-PCA and GPower_{ℓ_0} are significantly faster than SPCA, which, for a long time, was widely accepted as the algorithm that can handle large datasets.

Table 4 Computation time (in seconds) to obtain the first sparse PC, averaged over cardinalities ranging from 1 to n , for the colon cancer, leukemia and Ramaswamy datasets

n	Colon cancer 2000	Leukemia 7129	Ramaswamy 16063
SPCA	2.057	3.548	38.731
GPower $_{\ell_0}$	0.182	0.223	2.337
DC-PCA	0.034	0.156	0.547

Overall, the results in this section demonstrate that DC-PCA performs similar to or better than various sparse PCA algorithms proposed in literature, both in terms of scalability and proportion of variance explained vs. cardinality. As mentioned before, DC-PCA (see Algorithm 2) is a special instance of a more general framework (see Algorithm 1), that can be used to address other sparse generalized eigenvalue problems as well, e.g., sparse CCA, sparse FDA, etc., whereas most other sparse PCA algorithms cannot be readily extended to these other settings.

4.2 Sparse canonical correlation analysis

In this section, we consider sparse CCA as a special case of the sparse GEV algorithm and present two CCA applications where sparsity is helpful. We call our sparse CCA algorithm DC-CCA, where \mathbf{A} and \mathbf{B} are determined from the covariance and cross-covariance matrices as explained right below (3). Note that \mathbf{A} is indefinite and, therefore, in our experiments, we choose $\tau = -\lambda_{\min}(\mathbf{A})$ in Algorithm 1. In the following, we present two sparse CCA applications, one related to the task of cross-language document retrieval and the other dealing with semantic annotation and retrieval of music (Torres et al. 2007a, 2007b).

4.2.1 Cross-language document retrieval

The problem of cross-language document retrieval involves a collection, $\{D_i\}_{i=1}^N$ of documents, with each document being represented in different languages, say English and French. The goal of the task is, given a query string in one language, retrieve the most relevant document(s) in the target language. The first step is to obtain a semantic representation of the documents in both languages, which models the correlation between translated versions, so we can detect similarities in content between the two document spaces (one for English and the other for French). This is exactly what CCA does by finding a low-dimensional representation in both languages, with maximal correlation between them. Vinokourov et al. (2003) used CCA to address this problem and showed that the CCA approach performs better than the latent semantic indexing approach used by Littman et al. (1998). CCA provides an *efficient* basis representation (that captures the maximal correlation) for the two document spaces.

Using a bag-of-words representation for the documents, sparse CCA would allow to find a low-dimensional model based on a small subset of words in both languages. This would improve the interpretability of the model and could identify small subsets of words that are used in similar contexts in both languages and, possibly, are translations of one another. Representing documents by their similarity to all other documents (e.g., by taking inner products

of bag-of-word vectors, as explained below), sparse CCA would create a low-dimensional model that only requires to measure the similarity for a small subset of the training documents. This would immediately improve storage requirements and the efficiency of retrieval computations. In this study, we follow the second approach, representing documents by their similarity to all other training documents by applying a linear kernel function to a binary bag-of-words representation of the documents, as proposed in Vinokourov et al. (2003). This will illustrate how we can achieve significant sparsity without significant loss of retrieval performance.

More specifically, each version of a document (English or French) is modeled using a bag-of-words feature vector. Within a feature vector, we associate an element in $\{0, 1\}$ with each word w_i in its language vocabulary. A value of 1 indicates that w_i is found in the document. We collect the feature vectors into the $n \times p$ matrix \mathbf{E} , where we collect the English feature vectors, and the $n \times q$ matrix \mathbf{F} , where we collect the French feature vectors. n is the number of documents and p and q are the vocabulary sizes of \mathbf{E} and \mathbf{F} respectively. Computing the similarity between English documents as the inner product between their binary bag-of-words vectors (i.e., the rows of \mathbf{E}) results in computing an $n \times n$ data matrix $\mathbf{E}\mathbf{E}^T$. Similarly, we compute an $n \times n$ data matrix $\mathbf{F}\mathbf{F}^T$ and obtain two feature spaces which are both n -dimensional.

By applying sparse CCA, we effectively perform simultaneous feature selection across two vector spaces and characterize the content of and correlation between English and French documents in an efficient manner. We use the DC-CCA algorithm, using the covariance and cross-variance matrices associated with the document matrices $\mathbf{E}\mathbf{E}^T$ and $\mathbf{F}\mathbf{F}^T$ and obtain successive pairs of sparse canonical components which we stack into the columns of \mathbf{V}_E and \mathbf{V}_F . Subsequent pairs of these sparse canonical components are obtained by deflating $\mathbf{E}\mathbf{E}^T$ and $\mathbf{F}\mathbf{F}^T$ with respect to previous canonical components. For a detailed review on deflation, we refer the reader to Shawe-Taylor and Cristianini (2004).

Then, given a query document in an input language, say English, we convert the query into the appropriate feature vector, \mathbf{q}_E . We project \mathbf{q}_E onto the subspace spanned by the sparse canonical components in the English language space by computing $\mathbf{V}_E^T \mathbf{q}_E$.⁵ Similarly, we project all the French documents onto the subspace spanned by the sparse canonical components, \mathbf{V}_F associated with the French language. Finally, we perform document retrieval by selecting those French documents whose projections are closest to the projected query, where we measure distance in a nearest neighbor sense.

Experimental details The data set used was the Aligned Hansards of the 36th Parliament of Canada (Germann 2001), which is a collection of 1.3 million pairs of text chunks (sentences or smaller fragments) aligned into English and French translations. The text chunks are split into documents based on *** delimiters. After removing stop words and rare words (those that occur less than 3 times), we are left with an 1800×26328 English document-by-term matrix and a 1800×30167 French matrix. Computing $\mathbf{E}\mathbf{E}^T$ and $\mathbf{F}\mathbf{F}^T$ results in matrices of size 1800×1800 .

To generate a query, we select English test documents from a test set not used for training. The appropriate retrieval result is the corresponding French language version of the query document. To perform retrieval, the query and the French test documents are projected onto the sparse canonical components and retrieval is performed as described before. Table 5 shows the performance of DC-CCA (sparse CCA) against CCA. We measure our

⁵Notice how this projection, onto the sparse canonical components, only requires to compute a few elements of \mathbf{q}_E , i.e., the ones corresponding to the non-zero loadings of the sparse canonical components; differently said, we only need to compute the similarity of the query document to a small subset of all training documents.

Table 5 Average area under the ROC curve (in %) using CCA and sparse CCA (DC-CCA) in a cross-language document retrieval task. d represents the number of canonical components and *sparsity* represents the percentage of zero loadings in the canonical components

d	100	200	300	400	500
CCA	99.92	99.93	99.96	99.95	99.93
DC-CCA	95.72	97.57	98.45	98.75	99.04
Sparsity	87.15	87.56	87.95	88.21	88.44

results using the average area under the ROC curve (average AROC). The results in Table 5 are shown in percentages. To go into detail, for each test query we generate an ROC curve from the ranked retrieval results. Results are ranked according to their projected feature vector's Euclidean distance from the query. The area under this ROC curve is used to measure performance. For example, if the first returned document was the most relevant (i.e., the corresponding French language version of the query document) this would result in an ROC with area under the curve (AROC) of 1. If the most relevant document came in above the 75th percentile of all documents, this would lead to an AROC of 0.75, and so on. So, we're basically measuring how highly the corresponding French language document ranks in the retrieval results. For a collection of queries we take the simple average of each query's AROC to obtain the average AROC. An average AROC of 1 is best, a value of 0.5 is as good as chance.

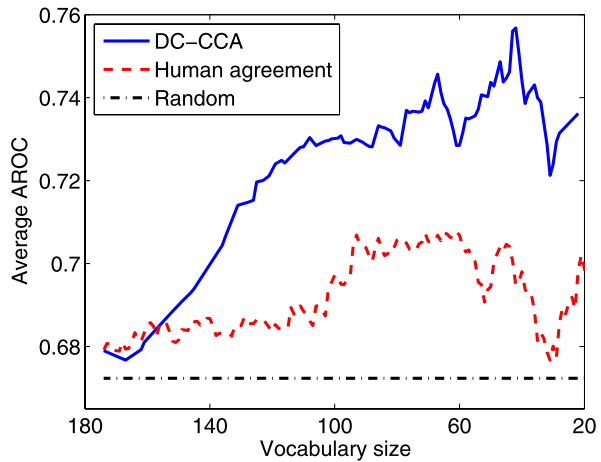
In Table 5, we compare retrieval using sparse CCA to regular CCA. For sparse CCA, we use a sparsity parameter that leads to loadings that are approximately 10% of the full dimensionality, i.e., the canonical components are approximately 90% sparse. We note that sparse CCA is able to achieve good retrieval rates, only slightly sacrificing performance compared to regular CCA. This is the key result of this section: we can achieve performance close to regular CCA, by using only about 12% of the number of loadings (i.e., documents) required by regular CCA. This shows that sparse CCA can narrow in on the most informative dimensions exhibited by data and can be used as an effective dimensionality reduction technique.

In summary, Torres et al. (2007a) illustrates that vocabulary selection using sparse CCA significantly improves the retrieval performance of a computer audition system (by effectively removing noisy words), outperforming a random baseline and a human agreement heuristic.

4.2.2 Vocabulary selection for music information retrieval

In this section we provide a short summary of the results in Torres et al. (2007a), which nicely illustrate how sparse CCA can be used to improve the performance of a statistical musical query application, by identifying problematic query words and eliminating them from the model. The application involves a computer audition system (Turnbull et al. 2008) that can annotate songs with semantically meaningful words or *tags* (such as, e.g., *rock* or *mellow*), or retrieve songs from a database, based on a semantic query. This system is based on a joint probabilistic model between words and acoustic signals, learned from a training data set of songs and song tags. "Noisy" words, that are not or only weakly related to the musical content, will decrease the system's performance and waste computational resources. Sparse CCA is employed to prune away those noisy words and improve the system's performance.

Fig. 5 Comparison of vocabulary selection techniques for music retrieval



The details of this experiment are beyond the scope of this work and can be found in Torres et al. (2007a). In short, each song from the CAL-500 dataset⁶ is represented in two different spaces: in a semantic space, based on a bag-of-words representation of a song's semantic tags, and in an audio space, based on Mel-frequency cepstral coefficients (McKinney 2003) extracted from a song's audio content. This representation allows sparse CCA to identify a small subset of words spanning a semantic subspace that is highly correlated with audio content. In Fig. 5, we use sparse CCA to generate a sequence of vocabularies of progressively smaller size, ranging from full size (containing about 180 words) to very sparse (containing about 20 words), depicted on the horizontal axis. For each vocabulary size, the computer audition system is trained and the average area under the receiver operating characteristic curve (AROC) is shown on the vertical axis, measuring its retrieval performance on an independent test set. The AROC (ranging between 0.5 for randomly ranked retrieval results and 1.0 for a perfect ranking) initially clearly improves, as sparse CCA (DC-CCA) generates vocabularies of smaller size: it is effectively removing noisy words that are detrimental for the system's performance. Also shown in Fig. 5 are the results of training the music retrieval system based on two alternative vocabulary selection techniques: random selection (offering no improvement) and a heuristic that eliminates words exhibiting less agreement amongst the human subjects that were surveyed to collect the CAL-500 dataset (only offering a slight improvement, initially).

4.3 Sparse fisher discriminant analysis

In this section, we show that the FDA problem is an *interesting* special case of the GEV problem and that the special structure of \mathbf{A} allows the sparse FDA problem to be solved more efficiently than the general sparse GEV problem.

Let us consider the GEV problem in (GEV-P) with $\mathbf{A} \in \mathbb{S}_{++}^n$, $\mathbf{B} \in \mathbb{S}_{++}^n$ and $\text{rank}(\mathbf{A}) = 1$. This is exactly the FDA problem as shown in (4) where \mathbf{A} is of the form $\mathbf{A} = \mathbf{a}\mathbf{a}^T$, with $\mathbf{a} = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \in \mathbb{R}^n$. The corresponding GEV problem, written as $\lambda_{\max}(\mathbf{A}, \mathbf{B}) = \max\{(\mathbf{a}^T \mathbf{x})^2 : \mathbf{x}^T \mathbf{B} \mathbf{x} = 1\}$, can be posed as a minimization problem as shown by the following result.

⁶The CAL-500 data set consists of a set of songs, annotated with semantic tags, obtained by conducting human surveys. More details can be found in Turnbull et al. (2008).

Proposition 2 Suppose \mathbf{x}_1 is the solution to

$$\max_{\mathbf{x}} \{(\mathbf{a}^T \mathbf{x})^2 : \mathbf{x}^T \mathbf{B} \mathbf{x} = 1\}, \quad (22)$$

and \mathbf{x}_2 is the solution to

$$\min_{\mathbf{x}} \{\mathbf{x}^T \mathbf{B} \mathbf{x} : \mathbf{a}^T \mathbf{x} = 1\}. \quad (23)$$

Then

$$\mathbf{x}_1 = \mathbf{x}_2 \sqrt{\mathbf{a}^T \mathbf{B}^{-1} \mathbf{a}}. \quad (24)$$

Proof Consider $\max\{(\mathbf{a}^T \mathbf{x})^2 : \mathbf{x}^T \mathbf{B} \mathbf{x} = 1\}$, whose Lagrangian is given as $L_1(\mathbf{x}, \lambda_1) = (\mathbf{a}^T \mathbf{x})^2 - \lambda_1(\mathbf{x}^T \mathbf{B} \mathbf{x} - 1)$, where $\lambda_1 > 0$ is the Lagrange multiplier. Differentiating L_1 w.r.t. \mathbf{x} and setting it to zero gives $\mathbf{a} \mathbf{a}^T \mathbf{x} = \lambda_1 \mathbf{B} \mathbf{x}$, i.e., $\mathbf{x} = \lambda_1^{-1} (\mathbf{a}^T \mathbf{x}) \mathbf{B}^{-1} \mathbf{a}$. Note that $\mathbf{a}^T \mathbf{x}$ determines the scale of \mathbf{x} and does not determine its direction. Therefore, let $\mathbf{x} = c \mathbf{B}^{-1} \mathbf{a}$, where $c \in \mathbb{R}$ is chosen such that $\mathbf{x}^T \mathbf{B} \mathbf{x} = 1$, which gives $\mathbf{x}_1 = (\mathbf{a}^T \mathbf{B}^{-1} \mathbf{a})^{-\frac{1}{2}} (\mathbf{B}^{-1} \mathbf{a})$.

On the other hand, the Lagrangian of $\min\{\mathbf{x}^T \mathbf{B} \mathbf{x} : \mathbf{a}^T \mathbf{x} = 1\}$ is given by $L_2(\mathbf{x}, \lambda_2) = \mathbf{x}^T \mathbf{B} \mathbf{x} - \lambda_2(\mathbf{a}^T \mathbf{x} - 1)$, where $\lambda_2 \in \mathbb{R}$ is the Lagrangian multiplier. Differentiating L_2 w.r.t. \mathbf{x} and setting it to zero gives $\mathbf{x} = \frac{\lambda_2}{2} \mathbf{B}^{-1} \mathbf{a}$, where λ_2 is chosen by setting $\mathbf{a}^T \mathbf{x} = 1$, therefore yielding $\mathbf{x}_2 = (\mathbf{a}^T \mathbf{B}^{-1} \mathbf{a})^{-1} \mathbf{B}^{-1} \mathbf{a}$. The result in (24) follows. \square

The above result shows that the solutions to (22) and (23) are the same except for scale. Therefore, instead of considering (22), we will consider its equivalent minimization formulation in (23), the advantage of which will become clear when we consider its sparse version, i.e., after introducing the constraint $\{\mathbf{x} : \|\mathbf{x}\|_0 \leq k\}$ in (23). Clearly, introducing the sparsity constraint in (23) makes it intractable. However, introducing an ℓ_1 -norm relaxation in this formulation gives rise to a *convex program*,

$$\min_{\mathbf{x}} \{\mathbf{x}^T \mathbf{B} \mathbf{x} : \mathbf{a}^T \mathbf{x} = 1, \|\mathbf{x}\|_1 \leq k\}, \quad (25)$$

more specifically a quadratic program (QP). On the other hand, as discussed in Sect. 3, the ℓ_1 -norm relaxation to the cardinality constraint, when used in (22), does not convexify the problem (and a convex approximation requires additional relaxation, e.g., as discussed in Appendix A). Therefore, the formulation in (23) is better suited to introduce sparsity than the one in (22).

Suppose that one would like to use a better approximation to $\|\mathbf{x}\|_0$ than $\|\mathbf{x}\|_1$, for sparse FDA. Using the approximation to $\|\mathbf{x}\|_0$ we consider in this work, the sparse version of (23), given by

$$\min_{\mathbf{x}} \{\mathbf{x}^T \mathbf{B} \mathbf{x} + v \|\mathbf{x}\|_0 : \mathbf{a}^T \mathbf{x} = 1\}, \quad (26)$$

is approximated as

$$\min_{\mathbf{x}} \left\{ \mathbf{x}^T \mathbf{B} \mathbf{x} + v_{\varepsilon} \sum_{i=1}^n \log(\varepsilon + |x_i|) : \mathbf{a}^T \mathbf{x} = 1 \right\}, \quad (27)$$

where (26) is the regularized version of $\min\{\mathbf{x}^T \mathbf{B} \mathbf{x} : \mathbf{a}^T \mathbf{x} = 1, \|\mathbf{x}\|_0 \leq k\}$, $v > 0$ is the regularization parameter and $v_{\varepsilon} := v / \log(1 + \varepsilon^{-1})$. Note that (27) can be written as a d.c. program and therefore, applying the MM method (see Example 1) to (27) results in the following iterative scheme,

$$\mathbf{x}^{(l+1)} = \arg \min_{\mathbf{x}} \left\{ \mathbf{x}^T \mathbf{B} \mathbf{x} + v_\varepsilon \sum_{i=1}^n \frac{|x_i|}{|x_i^{(l)}| + \varepsilon} : \mathbf{a}^T \mathbf{x} = 1 \right\}, \quad (28)$$

which is a sequence of QPs.⁷ In this case, Algorithm 1 would solve a sequence of QPs as well (since $\mathbf{A} \in \mathbb{S}_+^n$, we can have $\tau = 0$).

Suykens et al. (2002, Sect. 3.3) and Mika et al. (2001, Proposition 1) have shown connections between the FDA formulation in (23) with $\mathbf{a} = \boldsymbol{\mu}_1 - \boldsymbol{\mu}_2$ and $\mathbf{B} = \boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_2$ (see paragraph below (4) for details) and least squares SVMs⁸ (classifiers that minimize the squared loss, see Chap. 3 in Suykens et al. 2002). Therefore, sparse FDA is equivalent to feature selection with least squares SVMs, i.e., (25) is equivalent to LASSO, while the formulation in (27) is similar to the one considered in Weston et al. (2003). Since these are well studied problems, we do not pursue further showing the numerical performance of sparse FDA.

5 Conclusion and discussion

This work presents a general, efficient algorithm that allows to obtain sparse solutions to a generalized eigenvalue problem. After proposing a non-convex but tight approximation to the cardinality constraint, we formulate the resulting optimization problem as a d.c. program and derive an iterative algorithm, based on the majorization-minimization method. This results in solving a sequence of quadratically constrained quadratic programs, an algorithm which exhibits global convergence behavior, as we show. To empirically demonstrate the merit of this general framework, specific versions of the proposed algorithm are derived for some special cases, like, e.g., sparse PCA (DC-PCA) and sparse CCA (DC-CCA). For the special case of sparse PCA, we experimentally demonstrate on both benchmark and real-life datasets of varying dimensionality that DC-PCA has performance and scalability similar to the state-of-the-art for sparse PCA, while being derived from a more general framework that can readily be extended to other sparse generalized eigenvalue problems, unlike other sparse PCA algorithms. For the setting of sparse CCA, we illustrate the benefits of DC-CCA in two applications: cross-language document retrieval and vocabulary selection for music information retrieval.

The proposed algorithm does not allow to set the regularization parameter a priori, to guarantee a given sparsity level. SDP-based relaxation methods (see Appendix A), on the other hand, are better suited to achieve a given sparsity level in one shot, by incorporating an explicit constraint on the sparsity of the solution (although, eventually, through relaxation, an approximation of the original problem is solved). Since the algorithm we propose solves a LASSO problem in each step but with a quadratic constraint, in future work, we will explore a modified version of path following techniques like least angle regression (Efron et al. 2004) to learn the entire regularization path.

Another topic for future work is to study the behavior of the solution obtained for (SGEV-A), compared to that of (SGEV-R) as $\varepsilon \rightarrow 0$. As mentioned in Sect. 3.1, at present, we do not have any theoretical guarantees about the behavior of \mathbf{x}_ε or $Q_\varepsilon(\mathbf{x}_\varepsilon)$ as $\varepsilon \rightarrow 0$. In

⁷Although (28) is not a QP in standard form, it can be reformulated as the following standard QP: $\mathbf{z}^{(l+1)} = \arg \min_{\mathbf{z}} \{ \mathbf{z}^T \mathbf{C} \mathbf{z} + \mathbf{b}^T \mathbf{z} : \mathbf{e}^T \mathbf{z} = 1, \mathbf{F} \mathbf{z} \leq \mathbf{0} \}$, where $\mathbf{z} := [\mathbf{x}^T \ \mathbf{y}^T]^T$, $\mathbf{C} = \begin{bmatrix} \mathbf{B} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$, $\mathbf{b} = \begin{bmatrix} \mathbf{0} \\ v_\varepsilon \mathbf{G}^{-1} \mathbf{1} \end{bmatrix}$, $\mathbf{F} = \begin{bmatrix} \mathbf{I}_n & -\mathbf{I}_n \\ -\mathbf{I}_n & -\mathbf{I}_n \end{bmatrix}$, $\mathbf{G} = \mathbf{D}(\mathbf{y}^{(l)} + \varepsilon)$ and $\mathbf{e} = \begin{bmatrix} \mathbf{a} \\ \mathbf{0} \end{bmatrix}$. Here, y_i are the auxiliary variables introduced for $|x_i|$ (see Boyd and Vandenberghe 2004, Chap. 4). Therefore, for simplicity, we refer to (28) as a QP.

⁸Here, we mean the linear classifier case and not the kernelized version.

future work, we would like to study in detail the conditions under which $Q_\varepsilon(\mathbf{x}_\varepsilon) \rightarrow Q(\hat{\mathbf{x}})$ as $\varepsilon \rightarrow 0$ (see Sect. 3.1 for details on the notation).

Acknowledgements The authors thank the editors and reviewers for their constructive comments. B.K.S. thanks Suvrit Sra for helpful discussions while the former was an intern at the Max Planck Institute for Biological Cybernetics, Tübingen. The authors wish to acknowledge support from the National Science Foundation (grant DMS-MSPA 0625409), the Fair Isaac Corporation and the University of California MICRO program.

Open Access This article is distributed under the terms of the Creative Commons Attribution Noncommercial License which permits any noncommercial use, distribution, and reproduction in any medium, provided the original author(s) and source are credited.

Appendix A: Semidefinite programming relaxation for (SGEV-P)

Let us consider the following approximate program that is obtained by relaxing the non-convex constraint set $\{\mathbf{x} : \mathbf{x}^T \mathbf{B} \mathbf{x} = 1\} \cap \{\|\mathbf{x}\|_0 \leq k\}$ in (SGEV-P):

$$\max_{\mathbf{x}} \{\mathbf{x}^T \mathbf{A} \mathbf{x} : \mathbf{x}^T \mathbf{B} \mathbf{x} \leq 1, \|\mathbf{x}\|_1 \leq k\}. \quad (29)$$

As mentioned before, this program is still intractable due to the maximization of the non-concave objective. Had the objective function been linear, (29) would have been a canonical convex program, which could then be solved efficiently. One approach to linearize the objective function is by using the *lifting* technique (Lemaréchal and Oustry 1999, Sect. 4.4), which was considered by d’Aspremont et al. (2005) when $\mathbf{A} \in \mathbb{S}_+^n$ and $\mathbf{B} = \mathbf{I}_n$. The lifted version of (29) is given by,

$$\max_{\mathbf{X}, \mathbf{x}} \{\text{tr}(\mathbf{X} \mathbf{A}) : \text{tr}(\mathbf{X} \mathbf{B}) \leq 1, \|\mathbf{x}\|_1 \leq k, \mathbf{X} = \mathbf{x} \mathbf{x}^T\}.$$

Note that in the above program, the objective function is linear in \mathbf{X} , and the constraints are convex except for the non-convex constraint, $\mathbf{X} = \mathbf{x} \mathbf{x}^T$.⁹ Relaxing $\mathbf{X} = \mathbf{x} \mathbf{x}^T$ to $\mathbf{X} - \mathbf{x} \mathbf{x}^T \in \mathbb{S}_+^n$ results in the following program

$$\max_{\mathbf{X}, \mathbf{x}} \{\text{tr}(\mathbf{X} \mathbf{A}) : \text{tr}(\mathbf{X} \mathbf{B}) \leq 1, \|\mathbf{x}\|_1 \leq k, \mathbf{X} - \mathbf{x} \mathbf{x}^T \in \mathbb{S}_+^n\}, \quad (30)$$

which is a semidefinite program (SDP). The ℓ_1 -norm constraint in (30) can be relaxed as $\|\mathbf{x}\|_1^2 \leq k^2 \Rightarrow \mathbf{1}^T \mathbf{X} \mathbf{1} \leq k^2$ so that the problem reduces to solving only for \mathbf{X} . Here $[\mathbf{X}]_{ij} = |[\mathbf{X}]_{ij}|$. Therefore, we have obtained a tractable convex approximation to (SGEV-P).

Although (30) is a *convex* approximation to (SGEV-P), it is computationally very intensive as general purpose interior-point methods for SDP scale as $O(n^6 \log \epsilon^{-1})$, where ϵ is the required accuracy on the optimal value. For large-scale problems, first-order methods (Nesterov 2005; d’Aspremont et al. 2007) can be used which scale as $O(\epsilon^{-1} n^4 \sqrt{\log n})$. Therefore, the SDP-based convex relaxation to (SGEV-P) is prohibitively expensive in computation for large n .

⁹ $\mathbf{X} = \mathbf{x} \mathbf{x}^T \Leftrightarrow \mathbf{X} \in \mathbb{S}_+^n$, $\text{rank}(\mathbf{X}) = 1$, where $\text{rank}(\mathbf{X}) = 1$ is a non-convex constraint and therefore $\mathbf{X} = \mathbf{x} \mathbf{x}^T$ is a non-convex constraint.

Appendix B: Approximation to $\|\mathbf{x}\|_0$

In Sect. 3.1, we approximated $\|\mathbf{x}\|_0$ by

$$\|\mathbf{x}\|_\varepsilon := \sum_{i=1}^n \frac{\log(1 + |x_i|\varepsilon^{-1})}{\log(1 + \varepsilon^{-1})}.$$

We now show that $\|\mathbf{x}\|_\varepsilon$ is a tighter approximation to $\|\mathbf{x}\|_0$ than $\|\mathbf{x}\|_1$, for any $\varepsilon > 0$. To this end, let us define $a_\varepsilon := \frac{\log(1+a\varepsilon^{-1})}{\log(1+\varepsilon^{-1})}$, where $a \geq 0$, so that $\|\mathbf{x}\|_\varepsilon = \sum_{i=1}^n |x_i| a_\varepsilon$. Note that $\|\mathbf{x}\|_0 = \lim_{\varepsilon \rightarrow 0} \|\mathbf{x}\|_\varepsilon$ and

$$\lim_{\varepsilon \rightarrow \infty} \|\mathbf{x}\|_\varepsilon = \lim_{h \rightarrow 0} \sum_{i=1}^n \frac{\log(1 + |x_i|h)}{\log(1 + h)} \stackrel{(a)}{=} \sum_{i=1}^n \lim_{h \rightarrow 0} \frac{|x_i|(1+h)}{(1+h|x_i|)} = \sum_{i=1}^n |x_i| = \|\mathbf{x}\|_1,$$

where we have invoked L'Hôpital's rule in (a). In addition, we have $a > a_{\varepsilon_1} > a_{\varepsilon_2} > \dots > 1$ for $a > 1$ and $1 > \dots > a_{\varepsilon_2} > a_{\varepsilon_1} > a$ for $0 < a < 1$, if $\varepsilon_1 > \varepsilon_2 > \dots$, i.e., for any $a > 0$ and any $0 < \varepsilon < \infty$, the value a_ε is closer to 1 than a is to 1. This means a_ε for any $0 < \varepsilon < \infty$ is a better approximation to $\mathbb{1}_{\{a \neq 0\}}$ than a is to $\mathbb{1}_{\{a \neq 0\}}$. Therefore, $\|\mathbf{x}\|_\varepsilon$ for any $0 < \varepsilon < \infty$ is a better approximation to $\|\mathbf{x}\|_0$ than $\|\mathbf{x}\|_1$ is to $\|\mathbf{x}\|_0$.

Appendix C: Derivation of (ALG-R) and (ALG-S)

Let $\mathbf{A} \in \mathbb{S}_+^n$ and $\mathbf{B} \in \mathbb{S}_{++}^n$. Using $\tau = 0$ in (ALG), we have that $\mathbf{x}^{(l+1)}$ is a maximizer of the following program:

$$\begin{aligned} \max_{\mathbf{x}^T \mathbf{B} \mathbf{x} \leq 1} \mathbf{x}^T \mathbf{A} \mathbf{x}^{(l)} - \frac{\rho_\varepsilon}{2} \|\mathbf{W}^{(l)} \mathbf{x}\|_1 &= \max_{\mathbf{x}^T \mathbf{B} \mathbf{x} \leq 1} \sum_{i=1}^n \left[x_i (\mathbf{A} \mathbf{x}^{(l)})_i - \frac{\rho_\varepsilon}{2} w_i^{(l)} |x_i| \right] \\ &= \max_{\mathbf{x}^T \mathbf{B} \mathbf{x} \leq 1} \sum_{i=1}^n |x_i| \left[(\mathbf{A} \mathbf{x}^{(l)})_i \text{sign}(x_i) - \frac{\rho_\varepsilon}{2} w_i^{(l)} \right] \\ &\stackrel{(a)}{=} \max_{\mathbf{z}^T \mathbf{S} \mathbf{B} \mathbf{S} \mathbf{z} \leq 1, \mathbf{z} \geq \mathbf{0}} \sum_{i=1}^n z_i \left(|(\mathbf{A} \mathbf{x}^{(l)})_i| - \frac{\rho_\varepsilon}{2} w_i^{(l)} \right), \quad (31) \end{aligned}$$

where $x_i = s_i z_i$, $z_i \geq 0$ is used in (a) with $s_i := (\mathbf{s})_i = \text{sign}((\mathbf{A} \mathbf{x}^{(l)})_i)$ and $\mathbf{S} := \mathbf{D}(\mathbf{s})$. Note that $\mathbf{S} \mathbf{B} \mathbf{S} \in \mathbb{S}_+^n$. Define $\gamma_i := (\boldsymbol{\gamma})_i = |(\mathbf{A} \mathbf{x}^{(l)})_i| - \frac{\rho_\varepsilon}{2} w_i^{(l)}$. Then, we have

$$\mathbf{x}^{(l+1)} = \mathbf{S} \mathbf{z}^{(l+1)}, \quad (32)$$

where $\mathbf{z}^{(l+1)}$ is an optimizer of

$$\max_{\mathbf{z}} \{ \boldsymbol{\gamma}^T \mathbf{z} : \mathbf{z}^T \mathbf{S} \mathbf{B} \mathbf{S} \mathbf{z} \leq 1, \mathbf{z} \geq \mathbf{0} \}. \quad (33)$$

If $\rho_\varepsilon > 2 \max_i |(\mathbf{A} \mathbf{x}^{(l)})_i| (w_i^{(l)})^{-1}$, then $\boldsymbol{\gamma} < \mathbf{0}$ and, therefore, $\mathbf{z}^{(l+1)} = \mathbf{0}$ is the unique optimizer of (33). This implies $\mathbf{x}^{(l+1)} = \mathbf{0}$ is the unique optimal solution of (ALG). For $\rho_\varepsilon = 2 \max_i |(\mathbf{A} \mathbf{x}^{(l)})_i| (w_i^{(l)})^{-1}$, all the elements of $\boldsymbol{\gamma}$ are negative except for one or more that are zero, i.e., $\gamma_j = 0$ for $j \in \arg \max_i |(\mathbf{A} \mathbf{x}^{(l)})_i| (w_i^{(l)})^{-1}$ and $\gamma_i < 0$, $i \neq j$. So, we have

$(\mathbf{z}^{(l+1)})_i = 0$, $i \neq j$, while any $(\mathbf{z}^{(l+1)})_j \geq 0$ is optimal as long as $\mathbf{z}^T \mathbf{SBSz} \leq 1$. Apart from the uninteresting case $\rho_\varepsilon = 0$ and $\mathbf{x}^{(l)} = \mathbf{0}$, having $\rho_\varepsilon = 2 \max_i |(\mathbf{Ax}^{(l)})_i| (w_i^{(l)})^{-1}$ is clearly rather hypothetical, as ρ_ε is a positive *real* number chosen by the user, and having $\gamma_j = 0$ (and thus many $(\mathbf{z}^{(l+1)})_j \geq 0$ optimal) for multiple j even more hypothetical. Therefore, and since $\mathbf{0}$ is a solution of (ALG) when $\rho_\varepsilon = 2 \max_i |(\mathbf{Ax}^{(l)})_i| (w_i^{(l)})^{-1}$, we set $\mathbf{x}^{(l+1)} = \mathbf{0}$ as the optimal solution of (ALG) for $\rho_\varepsilon \geq 2 \max_i |(\mathbf{Ax}^{(l)})_i| (w_i^{(l)})^{-1}$, when implementing a practical algorithm (e.g., Algorithms 1 and 2).

Let $\rho_\varepsilon < 2 \max_i |(\mathbf{Ax}^{(l)})_i| (w_i^{(l)})^{-1}$. Consider the Lagrangian of (33):

$$L(\mathbf{z}, b, \boldsymbol{\lambda}) = -\boldsymbol{\gamma}^T \mathbf{z} + b(\mathbf{z}^T \mathbf{SBSz} - 1) - \boldsymbol{\lambda}^T \mathbf{z},$$

where $b \geq 0$ and $\boldsymbol{\lambda} \geq \mathbf{0}$. Differentiating L w.r.t. \mathbf{z} and setting it zero gives

$$\mathbf{z}^{(l+1)} = \frac{(\mathbf{SBS})^+(\boldsymbol{\gamma} + \boldsymbol{\lambda})}{2b}, \quad (34)$$

where \mathbf{C}^+ represents the Moore-Penrose pseudoinverse of \mathbf{C} . Substituting for $\mathbf{z}^{(l+1)}$ in L gives the following dual program:

$$\min_{b, \boldsymbol{\lambda}} \left\{ b + \frac{1}{4b} (\boldsymbol{\lambda} + \boldsymbol{\gamma})^T (\mathbf{SBS})^+ (\boldsymbol{\lambda} + \boldsymbol{\gamma}) : b \geq 0, \boldsymbol{\lambda} \geq \mathbf{0} \right\}. \quad (35)$$

Solving (35) w.r.t. b gives

$$b = \frac{1}{2} \sqrt{(\boldsymbol{\gamma} + \boldsymbol{\lambda})^T (\mathbf{SBS})^+ (\boldsymbol{\gamma} + \boldsymbol{\lambda})}, \quad (36)$$

which satisfies the constraint, $b \geq 0$, as $(\mathbf{SBS})^+ \in \mathbb{S}_+^n$. We show that $b > 0$. For this, let $\mathcal{I} := \{i : s_i \neq 0\}$. Note that for $i \notin \mathcal{I}$, $((\mathbf{SBS})^+)_{ij} = 0$, $\forall j$, which means $(\boldsymbol{\gamma} + \boldsymbol{\lambda})^T (\mathbf{SBS})^+ (\boldsymbol{\gamma} + \boldsymbol{\lambda}) = \sum_{i, j \in \mathcal{I}} (\gamma_i + (\boldsymbol{\lambda})_i)(\gamma_j + (\boldsymbol{\lambda})_j)((\mathbf{SBS})^+)_{ij}$. Let $\eta : \mathcal{I} \rightarrow \{1, \dots, |\mathcal{I}|\}$ be a bijective map. By defining $(\tilde{\boldsymbol{\gamma}})_{\eta(i)} := (\boldsymbol{\gamma})_i$, $i \in \mathcal{I}$, $(\tilde{\boldsymbol{\lambda}})_{\eta(i)} := (\boldsymbol{\lambda})_i$, $i \in \mathcal{I}$, $(\tilde{s})_{\eta(i)} = s_i$, $i \in \mathcal{I}$, $\tilde{\mathbf{S}} = \mathbf{D}(\tilde{s})$ and $(\tilde{\mathbf{B}})_{\eta(i), \eta(j)} = \mathbf{B}_{ij}$, $i, j \in \mathcal{I}$, we have that $(\boldsymbol{\gamma} + \boldsymbol{\lambda})^T (\mathbf{SBS})^+ (\boldsymbol{\gamma} + \boldsymbol{\lambda}) = (\tilde{\boldsymbol{\gamma}} + \tilde{\boldsymbol{\lambda}})^T (\tilde{\mathbf{S}} \tilde{\mathbf{B}} \tilde{\mathbf{S}})^{-1} (\tilde{\boldsymbol{\gamma}} + \tilde{\boldsymbol{\lambda}}) = (\tilde{\boldsymbol{\gamma}} + \tilde{\boldsymbol{\lambda}})^T \tilde{\mathbf{S}} \tilde{\mathbf{B}}^{-1} \tilde{\mathbf{S}} (\tilde{\boldsymbol{\gamma}} + \tilde{\boldsymbol{\lambda}})$. Suppose $b = 0$. This means $\tilde{\boldsymbol{\gamma}} + \tilde{\boldsymbol{\lambda}} = \mathbf{0}$ as $\tilde{\mathbf{S}} \tilde{\mathbf{B}}^{-1} \tilde{\mathbf{S}} \in \mathbb{S}_{++}^{|\mathcal{I}|}$ (since $\mathbf{B} \in \mathbb{S}_{++}^n$ and thus $\tilde{\mathbf{B}} \in \mathbb{S}_{++}^{|\mathcal{I}|}$), and, therefore, $\tilde{\boldsymbol{\lambda}} = -\tilde{\boldsymbol{\gamma}}$. However, this violates the constraint $\tilde{\boldsymbol{\lambda}} \geq \mathbf{0}$ because our assumption of $\rho_\varepsilon < 2 \max_i |(\mathbf{Ax}^{(l)})_i| (w_i^{(l)})^{-1}$ implies that there exists $i \in \{1, \dots, |\mathcal{I}|\}$ such that $(\tilde{\boldsymbol{\gamma}})_i > 0$ (since $(\boldsymbol{\gamma})_j \leq 0$ for $j \notin \mathcal{I}$). Therefore, $b > 0$. Substituting (36) in (35) yields

$$\boldsymbol{\lambda}^{(l+1)} \in \arg \min_{\boldsymbol{\lambda}} \left\{ (\boldsymbol{\gamma} + \boldsymbol{\lambda})^T (\mathbf{SBS})^+ (\boldsymbol{\gamma} + \boldsymbol{\lambda}) : \boldsymbol{\lambda} \geq \mathbf{0} \right\}. \quad (37)$$

From the above analysis, it is clear that $\tilde{\boldsymbol{\lambda}}^{(l+1)}$ is uniquely computed as

$$\tilde{\boldsymbol{\lambda}}^{(l+1)} = \arg \min_{\tilde{\boldsymbol{\lambda}}} \left\{ (\tilde{\boldsymbol{\gamma}} + \tilde{\boldsymbol{\lambda}})^T \tilde{\mathbf{S}} \tilde{\mathbf{B}}^{-1} \tilde{\mathbf{S}} (\tilde{\boldsymbol{\gamma}} + \tilde{\boldsymbol{\lambda}}) : \tilde{\boldsymbol{\lambda}} \geq \mathbf{0} \right\},$$

because $\tilde{\mathbf{S}} \tilde{\mathbf{B}}^{-1} \tilde{\mathbf{S}} \in \mathbb{S}_{++}^{|\mathcal{I}|}$ and, therefore, the objective is strictly convex in $\tilde{\boldsymbol{\lambda}}$, while $(\boldsymbol{\lambda}^{(l+1)})_i$, $i \notin \mathcal{I}$ cannot be uniquely computed. Combining (32), (34), (36) and (37), we have

$$\mathbf{x}^{(l+1)} = \frac{\mathbf{S}(\mathbf{SBS})^+(\boldsymbol{\gamma} + \boldsymbol{\lambda}^{(l+1)})}{\sqrt{(\boldsymbol{\gamma} + \boldsymbol{\lambda}^{(l+1)})^T (\mathbf{SBS})^+ (\boldsymbol{\gamma} + \boldsymbol{\lambda}^{(l+1)})}}. \quad (38)$$

Although $(\lambda^{(l+1)})_i$, $i \notin \mathcal{I}$ cannot be uniquely computed, $(x^{(l+1)})_i$, $i \notin \mathcal{I}$ is uniquely determined as zero because of the pre-multiplication by S . Therefore, $x^{(l+1)}$ is the unique optimizer of (31), as shown in (ALG-R).

If B is a diagonal matrix, i.e., $B = D(b)$ with $b = (b_1, \dots, b_n) \succ 0$, then (37) reduces to

$$\lambda^{(l+1)} = \arg \min_{\lambda} \left\{ \sum_{i \in \mathcal{I}} \frac{(\gamma_i + (\lambda)_i)^2}{b_i} : \lambda \geq 0 \right\}.$$

It is easy to see that $(\lambda^{(l+1)})_i = [-\gamma_i]_+$, $\forall i \in \mathcal{I}$, where $[a]_+ := \max(0, a)$. Substituting this in (38) yields

$$x_i^{(l+1)} = \frac{s_i [\gamma_i]_+}{b_i \sqrt{\sum_{i=1}^n b_i^{-1} [\gamma_i]_+^2}},$$

as shown in (ALG-S).

Appendix D: Proofs of results in Sect. 3.4

In this section, we present the proofs of Theorem 1 and Corollaries 1–3.

The proof of Theorem 1 is based on a result due to Sriperumbudur and Lanckriet (2009), which is mentioned as Theorem 2 below. In order to understand this result, we introduce some notation and terminology as follows.

The convergence analysis of an iterative procedure like Algorithm 1 uses the notion of a *set-valued mapping*, or *point-to-set mapping*, which is central to the theory of global convergence. A point-to-set map Ψ from a set X into a set Y is defined as $\Psi : X \rightarrow \mathcal{P}(Y)$, which assigns a subset of Y with each point of X , where $\mathcal{P}(Y)$ denotes the power set of Y . Ψ is said to be *uniformly compact* on X if there exists a compact set H independent of x such that $\Psi(x) \subset H$ for all $x \in X$. Note that if X is compact, then Ψ is uniformly compact on X . A *fixed point* of the map $\Psi : X \rightarrow \mathcal{P}(X)$ is a point x for which $\{x\} = \Psi(x)$.

Many iterative algorithms in mathematical programming can be described using the notion of point-to-set maps. Let X be a set and $x_0 \in X$ a given point. Then an *algorithm*, \mathcal{A} , with initial point x_0 is a point-to-set map $\mathcal{A} : X \rightarrow \mathcal{P}(X)$ which generates a sequence $\{x_k\}_{k=1}^\infty$ via the rule $x_{k+1} \in \mathcal{A}(x_k)$, $k = 0, 1, \dots$. \mathcal{A} is said to be *globally convergent* if for any chosen initial point x_0 , the sequence $\{x_k\}_{k=0}^\infty$ generated by $x_{k+1} \in \mathcal{A}(x_k)$ (or a subsequence) converges to a point for which a necessary condition of optimality holds: the Karush-Kuhn-Tucker (KKT) conditions in the case of constrained optimization and stationarity in the case of unconstrained optimization.

We now state the convergence result for (14) by Sriperumbudur and Lanckriet (2009), using which we provide the proof of Theorem 1.

Theorem 2 (Sriperumbudur and Lanckriet 2009) *Consider the program,*

$$\min\{u(x) - v(x) : x \in \Omega\}, \quad (\text{DC})$$

where $\Omega = \{x : c_i(x) \leq 0, i \in [m], d_j(x) = 0, j \in [p]\}$ and $[m] := \{1, \dots, m\}$. Let u and v be strictly convex, differentiable functions defined on \mathbb{R}^n . Also assume ∇v is continuous. Let $\{c_i\}$ be differentiable convex functions and $\{d_j\}$ be affine functions on \mathbb{R}^n . Suppose (DC) is solved iteratively as $x^{(l+1)} \in \mathcal{A}_{dc}(x^{(l)})$, where \mathcal{A}_{dc} is the point-to-set map defined as

$$\mathcal{A}_{dc}(y) = \arg \min_{x \in \Omega} u(x) - x^T \nabla v(y). \quad (\text{DC-ALG})$$

Let $\{x^{(l)}\}_{l=0}^\infty$ be any sequence generated by \mathcal{A}_{dc} defined by (DC-ALG). Suppose \mathcal{A}_{dc} is uniformly compact on Ω and $\mathcal{A}_{dc}(x)$ is nonempty for any $x \in \Omega$. Then, assuming suitable constraint qualification, all the limit points of $\{x^{(l)}\}_{l=0}^\infty$ are fixed points of (DC-ALG), which are stationary points of the d.c. program in (DC), $u(x^{(l)}) - v(x^{(l)}) \rightarrow u(x_*) - v(x_*) =: f^*$ as $l \rightarrow \infty$, for some fixed point x_* , $\|x^{(l+1)} - x^{(l)}\| \rightarrow 0$, and either $\{x^{(l)}\}_{l=0}^\infty$ converges or the set of limit points of $\{x^{(l)}\}_{l=0}^\infty$ is a connected and compact subset of $\mathcal{S}(f^*)$, where $\mathcal{S}(a) := \{x \in \mathcal{S} : u(x) - v(x) = a\}$ and \mathcal{S} is the set of fixed points of (DC-ALG). If $\mathcal{S}(f^*)$ is finite, then any sequence $\{x^{(l)}\}_{l=0}^\infty$ generated by \mathcal{A}_{dc} converges to some x_* in $\mathcal{S}(f^*)$.

Proof (Theorem 1) Since Algorithm 1 and (ALG) are equivalent, let \mathcal{A} correspond to the point-to-set map in (ALG). As noted before, (ALG) is obtained by applying linear majorization to (9), which is equivalent to (SGEV-A), with $\Omega = \{(x, y) : x^T Bx \leq 1, -y \leq x \leq y\}$, $u(x, y) = \tau \|x\|_2^2$ and $v(x, y) = x^T (A + \tau I_n)x - \rho_\varepsilon \sum_{i=1}^n \log(y_i + \varepsilon)$. Let $\tau \neq 0$. It is easy to check that u and v satisfy the conditions of Theorem 2. Define $z := (x^T, y^T)^T$. Note that $\Omega = \{z : z^T \tilde{A}z \leq 1, \tilde{B}z \leq 0, \tilde{C}z \geq 0\}$, where $\tilde{A} = \begin{bmatrix} B & 0 \\ 0 & 0 \end{bmatrix}$, $\tilde{B} = [I_n \ -I_n]$ and $\tilde{C} = [I_n \ I_n]$. $\Omega_1 := \{z : \tilde{B}z \leq 0\}$ and $\Omega_2 := \{z : \tilde{C}z \geq 0\}$ are closed sets in \mathbb{R}^{2n} , while $\Omega_3 := \{z : z^T \tilde{A}z \leq 1\}$ is compact in \mathbb{R}^{2n} . Note that $\Omega = \Omega_1 \cap \Omega_2 \cap \Omega_3$. Since Ω_1 and Ω_2 are closed, $\Omega_1 \cap \Omega_2$ is closed. Since Ω_3 is compact, $(\Omega_1 \cap \Omega_2) \cap \Omega_3$ is also compact, i.e., Ω is compact and therefore \mathcal{A} is uniformly compact. By the Weierstrass theorem¹⁰ (Minoux 1986), it is clear that $\mathcal{A}(z)$ is nonempty for any $z \in \Omega$. The result therefore follows from Theorem 2. When $\tau = 0$, all the conditions of Theorem 2 hold except for u being strictly convex. The strict convexity of u in Theorem 2 is mainly needed to ensure that (DC-ALG) has a unique optimizer (for details, see the proof of Theorem 2). Since we showed in Appendix C how Algorithm 1 obtains a unique optimizer in every iteration, when $\tau = 0$, the convergence result also holds for the case of $\tau = 0$. \square

Proof (Corollary 1) The proof idea is as follows. First, we consider Algorithm 1 with $\tau = 0$, $\rho_\varepsilon = \rho = 0$ and $A \in \mathbb{S}_+^n$ to show that it reduces to (21). Second, we consider the fixed points of (21) and show that the fixed point of (21), which is a stationary point of (SGEV-A) with $\rho = 0$, is indeed the generalized eigenvector associated with $\lambda_{\max}(A, B)$. Therefore, by Theorem 1, the sequence $\{x^{(l)}\}_{l=0}^\infty$ generated by (21) converges to the generalized eigenvector associated with $\lambda_{\max}(A, B)$.

Consider Algorithm 1 with $\tau = 0$, $\rho_\varepsilon = \rho = 0$ and $A \in \mathbb{S}_+^n$, which reduces to

$$x^{(l+1)} = \arg \min_x \{-x^T A x^{(l)} : x^T B x \leq 1\}. \quad (39)$$

Since the objective in (39) is linear in x , the minimum occurs at the boundary of the constraint set, i.e., $\{x : x^T B x = 1\}$. The corresponding Lagrangian is given by

$$L(x, \mu) = -2x^T A x^{(l)} + \mu(x^T B x - 1),$$

where $\mu > 0$ is the Lagrange multiplier. Differentiating L w.r.t. x and setting it to zero gives

$$x = \mu^{-1} B^{-1} A x^{(l)},$$

¹⁰The Weierstrass theorem states: If f is a real continuous function on a compact set $K \subset \mathbb{R}^n$, then the problem $\min\{f(x) : x \in K\}$ has an optimal solution $x^* \in K$.

where solving for μ using $\mathbf{x}^T \mathbf{B} \mathbf{x} = 1$ yields (21). So far, we have shown that (39) reduces to (21). Next, let us consider the fixed point, \mathbf{x}^* of (21). At \mathbf{x}^* , (21) reduces to $\mathbf{x}^* = ([\mathbf{x}^*]^T \mathbf{A} \mathbf{B}^{-1} \mathbf{A} \mathbf{x}^*)^{-\frac{1}{2}} \mathbf{B}^{-1} \mathbf{A} \mathbf{x}^*$, which can be written as $\mathbf{A} \mathbf{x}^* = \mu^* \mathbf{B} \mathbf{x}^*$, where $\mu^* = [\mathbf{x}^*]^T \mathbf{A} \mathbf{x}^*$. Note that \mathbf{x}^* is the stationary point of (SGEV-A), when $\rho = 0$. The objective of (39) at \mathbf{x}^* is $-[\mathbf{x}^*]^T \mathbf{A} \mathbf{x}^*$, which is minimized when \mathbf{x}^* is the generalized eigenvector associated with $\lambda_{\max}(\mathbf{A}, \mathbf{B})$. Since $\mathcal{S}(-\lambda_{\max}(\mathbf{A}, \mathbf{B})) = \{\pm \mathbf{x}^*\}$, the result follows from Theorem 1. \square

Proof (Corollary 2) When $\mathbf{B} = \mathbf{I}_n$, (21) reduces to $\mathbf{x}^{(l+1)} = \frac{\mathbf{A} \mathbf{x}^{(l)}}{\|\mathbf{A} \mathbf{x}^{(l)}\|_2}$, which is the power method for computing $\lambda_{\max}(\mathbf{A})$. The rest of the result simply follows from Corollary 1. \square

Proof (Corollary 3) It is easy to see that when $\mathbf{B} = \mathbf{I}_n$ and $\rho = 0$, Algorithm 1 reduces to

$$\mathbf{x}^{(l+1)} = \arg \min_{\mathbf{x}} \left\{ \tau \|\mathbf{x}\|_2^2 - 2\mathbf{x}^T (\mathbf{A} + \tau \mathbf{I}_n) \mathbf{x}^{(l)} : \mathbf{x}^T \mathbf{x} \leq 1 \right\}. \quad (40)$$

The Lagrangian of (40) is given by

$$L(\mathbf{x}, \mu) = \tau \|\mathbf{x}\|_2^2 - 2\mathbf{x}^T (\mathbf{A} + \tau \mathbf{I}_n) \mathbf{x}^{(l)} + \mu (\mathbf{x}^T \mathbf{x} - 1),$$

where $\mu \geq 0$ is the Lagrange multiplier. Differentiating L w.r.t. \mathbf{x} and setting it zero gives $\mathbf{x} = (\mu + \tau)^{-1} (\mathbf{A} + \tau \mathbf{I}_n) \mathbf{x}^{(l)}$, i.e.,

$$\mathbf{x}^{(l+1)} = (\mu^{(l+1)} + \tau)^{-1} (\mathbf{A} + \tau \mathbf{I}_n) \mathbf{x}^{(l)},$$

where

$$\mu^{(l+1)} ([\mathbf{x}^{(l+1)}]^T \mathbf{x}^{(l+1)} - 1) = 0 \quad \text{and} \quad [\mathbf{x}^{(l+1)}]^T \mathbf{x}^{(l+1)} \leq 1.$$

At the fixed point, \mathbf{x}^* , we have $\mathbf{A} \mathbf{x}^* = \mu^* \mathbf{x}^*$, where $[\mathbf{x}^*]^T \mathbf{x}^* \leq 1$ and $\mu^* ([\mathbf{x}^*]^T \mathbf{x}^* - 1) = 0$. Therefore, $[\mathbf{x}^*]^T \mathbf{A} \mathbf{x}^* = \mu^*$. In addition, the objective of (40) at \mathbf{x}^* is given by

$$\psi(\mathbf{x}^*) = -2[\mathbf{x}^*]^T \mathbf{A} \mathbf{x}^* - \tau \|\mathbf{x}^*\|_2^2 = -2\mu^* - \tau \|\mathbf{x}^*\|_2^2,$$

which is minimized when $\mu^* = \lambda_{\max}(\mathbf{A})$, i.e., \mathbf{x}^* is the eigenvector associated with $\lambda_{\max}(\mathbf{A})$. The result therefore follows from Theorem 1 as $\mathcal{S}(-\lambda_{\max}(\mathbf{A})) = \{\pm \mathbf{x}^*\}$. \square

References

- Alon, U., Barkai, N., Notterman, D. A., Gish, K., Ybarra, S., Mack, D., & Levine, A. J. (1999). Broad patterns of gene expression revealed by clustering analysis of tumor and normal colon cancer tissues. *Cell Biology*, 96, 6745–6750.
- Bioucas-Dias, J., Figueiredo, M., & Oliveira, J. (2006). Total-variation image deconvolution: a majorization-minimization approach. In: *Proc. IEEE international conference on acoustics, speech, and signal processing*, Toulouse, France.
- Böhning, D., & Lindsay, B. G. (1988). Monotonicity of quadratic-approximation algorithms. *Annals of the Institute of Statistical Mathematics*, 40(4), 641–663.
- Bonnans, J. F., Gilbert, J. C., Lemaréchal, C., & Sagastizábal, C. A. (2006). *Numerical optimization: theoretical and practical aspects*. Berlin: Springer.
- Boyd, S. P., & Vandenberghe, L. (2004). *Convex optimization*. Cambridge: Cambridge University Press.
- Bradley, P. S., & Mangasarian, O. L. (1998). Feature selection via concave minimization and support vector machines. In *Proc. 15th international conf. on machine learning* (pp. 82–90). San Francisco: Kaufmann.
- Cadima, J., & Jolliffe, I. (1995). Loadings and correlations in the interpretation of principal components. *Applied Statistics*, 22, 203–214.

- Candes, E. J., Wakin, M., & Boyd, S. (2007). Enhancing sparsity by reweighted ℓ_1 minimization. *The Journal of Fourier Analysis and Applications*, 14, 877–905.
- d'Aspremont, A., El Ghaoui, L., Jordan, M. I., & Lanckriet, G. R. G. (2005). A direct formulation for sparse PCA using semidefinite programming. In L. K. Saul, Y. Weiss, & L. Bottou (Eds.), *Advances in neural information processing systems* (Vol. 17, pp. 41–48). Cambridge: MIT Press.
- d'Aspremont, A., El Ghaoui, L., Jordan, M. I., & Lanckriet, G. R. G. (2007). A direct formulation for sparse PCA using semidefinite programming. *SIAM Review*, 49(3), 434–448.
- d'Aspremont, A., Bach, F. R., & El Ghaoui, L. (2008). Optimal solutions for sparse principal component analysis. *Journal of Machine Learning Research*, 9, 1269–1294.
- Daubechies, I., Defrise, M., & Mol, C. D. (2004). An iterative thresholding algorithm for linear inverse problems with a sparsity constraint. *Communications on Pure and Applied Mathematics*, 57, 1413–1457.
- de Leeuw, J. (1977). Applications of convex analysis to multidimensional scaling. In J. R. Barra, F. Brodeau, G. Romier, & B. V. Cutsem (Eds.), *Recent advantages in statistics* (pp. 133–146). Amsterdam: North Holland.
- Dempster, A. P., Laird, N. M., & Rubin, D. B. (1977). Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society, Series B*, 39, 1–38.
- Efron, B., Hastie, T., Johnstone, I., & Tibshirani, R. (2004). Least angle regression. *Annals of Statistics*, 32(2), 407–499.
- Fazel, M., Hindi, H., & Boyd, S. (2003). Log-det heuristic for matrix rank minimization with applications to Hankel and Euclidean distance matrices. In: *Proc. American control conference*, Denver, Colorado.
- Figueiredo, M. A. T., Bioucas-Dias, J. M., & Nowak, R. D. (2007). Majorization-minimization algorithms for wavelet-based image restoration. *IEEE Transactions on Image Processing*, 16, 2980–2991.
- Germann, U. (2001) *Aligned Hansards of the 36th parliament of Canada*. <http://www.isi.edu/natural-language/download/hansard/>.
- Golub, T. R., Slonim, D. K., Tamayo, P., Huard, C., Gaasenbeek, M., Mesirov, J., Collier, H., Loh, M. K., Downing, J. R., Caligiuri, M. A., Bloomfield, C. D., & Lander, E. (1999). Molecular classification of cancer: class discovery and class prediction by gene expression monitoring. *Science*, 286, 531–537.
- Heiser, W. J. (1987). Correspondence analysis with least absolute residuals. *Computational Statistics and Data Analysis*, 5, 337–356.
- Horst, R., & Thoai, N. V. (1999). D.c. programming: overview. *Journal of Optimization Theory and Applications*, 103, 1–43.
- Hotelling, H. (1933). Analysis of a complex of statistical variables into principal components. *Journal of Educational Psychology*, 24, 417–441.
- Hotelling, H. (1936). Relations between two sets of variates. *Biometrika*, 28, 321–377.
- Huber, P. J. (1981). *Robust statistics*. New York: Wiley.
- Hunter, D. R., & Lange, K. (2004). A tutorial on MM algorithms. *The American Statistician*, 58, 30–37.
- Hunter, D. R., & Li, R. (2005). Variable selection using MM algorithms. *The Annals of Statistics*, 33, 1617–1642.
- Jeffers, J. (1967). Two case studies in the application of principal components. *Applied Statistics*, 16, 225–236.
- Jolliffe, I. (1986). *Principal component analysis*. New York: Springer.
- Jolliffe, I. T., Trendafilov, N. T., & Uddin, M. (2003). A modified principal component technique based on the LASSO. *Journal of Computational and Graphical Statistics*, 12, 531–547.
- Journée, M., Nesterov, Y., Richtárik, P., & Sepulchre, R. (2010). Generalized power method for sparse principal component analysis. *Journal of Machine Learning Research*, 11, 517–553.
- Lange, K., Hunter, D. R., & Yang, I. (2000). Optimization transfer using surrogate objective functions with discussion. *Journal of Computational and Graphical Statistics*, 9(1), 1–59.
- Lee, D. D., & Seung, H. S. (2001). Algorithms for non-negative matrix factorization. In T. Leen, T. Dietterich, & V. Tresp (Eds.), *Advances in neural information processing systems* (Vol. 13, pp. 556–562). Cambridge: MIT Press.
- Lemaréchal, C., & Oustry, F. (1999). *Semidefinite relaxations and Lagrangian duality with application to combinatorial optimization* (Tech. Rep. RR3710). INRIA.
- Littman, M. L., Dumais, S. T., & Landauer, T. K. (1998). Automatic cross-language information retrieval using latent semantic indexing. In G. Grefenstette (Ed.), *Cross-language information retrieval* (pp. 51–62). Norwell: Kluwer Academic.
- Mackey, L. (2009). Deflation methods for sparse PCA. In D. Koller, D. Schuurmans, Y. Bengio, & L. Bottou (Eds.), *Advances in neural information processing systems* (Vol. 21, pp. 1017–1024). Cambridge: MIT Press.
- McCabe, G. (1984). Principal variables. *Technometrics*, 26, 137–144.

- McKinney, M. F. (2003). Features for audio and music classification. In *Proc. of the international symposium on music information retrieval* (pp. 151–158).
- Meng, X. L. (2000). Discussion on “optimization transfer using surrogate objective functions”. *Journal of Computational and Graphical Statistics*, 9(1), 35–43.
- Mika, S., Rätsch, G., & Müller, K. R. (2001). A mathematical programming approach to the kernel Fisher algorithm. In T. Leen, T. Dietterich, & V. Tresp (Eds.), *Advances in neural information processing systems* (Vol. 13). Cambridge: MIT Press.
- Minoux, M. (1986). *Mathematical programming: theory and algorithms*. New York: Wiley.
- Moghaddam, B., Weiss, Y., & Avidan, S. (2007a). Generalized spectral bounds for sparse LDA. In: *Proc. of international conference on machine learning*.
- Moghaddam, B., Weiss, Y., & Avidan, S. (2007b). Spectral bounds for sparse PCA: exact and greedy algorithms. In B. Schölkopf, J. Platt, & T. Hoffman (Eds.), *Advances in neural information processing systems* (Vol. 19). Cambridge: MIT Press.
- Nesterov, Y. (2005). Smooth minimization of non-smooth functions. *Mathematical Programming, Series A*, 103, 127–152.
- Ortega, J. M., & Rheinboldt, W. C. (1970). *Iterative solution of nonlinear equations in several variables*. New York: Academic Press.
- Ramaswamy, S., Tamayo, P., Rifkin, R., Mukherjee, S., Yeang, C., Angelo, M., Ladd, C., Reich, M., Latulippe, E., Mesirov, J., Poggio, T., Gerald, W., Loda, M., Lander, E., & Golub, T. (2001). Multiclass cancer diagnosis using tumor gene expression signature. *Proceedings of the National Academy of Sciences*, 98, 15,149–15,154.
- Rockafellar, R. T. (1970). *Convex analysis*. Princeton: Princeton University Press.
- Shawe-Taylor, J., & Cristianini, N. (2004). *Kernel methods for pattern analysis*. Cambridge: Cambridge University Press.
- Sriperumbudur, B. K., & Lanckriet, G. R. G. (2009). On the convergence of the concave-convex procedure. In Y. Bengio, D. Schuurmans, J. Lafferty, C. K. I. Williams, & A. Culotta (Eds.), *Advances in neural information processing systems* (Vol. 22, pp. 1759–1767). Cambridge: MIT Press.
- Sriperumbudur, B. K., Torres, D. A., & Lanckriet, G. R. G. (2007). Sparse eigen methods by d.c. programming. In: *Proc. of the 24th annual international conference on machine learning*.
- Sriperumbudur, B. K., Torres, D. A., & Lanckriet, G. R. G. (2009). A d.c. programming approach to the sparse generalized eigenvalue problem. In: *Optimization for machine learning workshop, NIPS*.
- Strang, G. (1986). *Introduction to applied mathematics*. Wellesley: Cambridge Press.
- Suykens, J. A. K., Gestel, T. V., Brabanter, J. D., Moor, B. D., & Vandewalle, J. (2002). *Least squares support vector machines*. Singapore: World Scientific.
- Tibshirani, R. (1996). Regression shrinkage and selection via the LASSO. *Journal of Royal Statistical Society, Series B*, 58(1), 267–288.
- Tipping, M. E. (2001). Sparse Bayesian learning and the relevance vector machine. *Journal of Machine Learning Research*, 1, 211–244.
- Torres, D., Turnbull, D., Barrington, L., & Lanckriet, G. R. G. (2007a). Identifying words that are musically meaningful. In: *Proc. of international symposium on music information and retrieval*.
- Torres, D. A., Turnbull, D., Sriperumbudur, B. K., Barrington, L., & Lanckriet, G. R. G. (2007b). Finding musically meaningful words using sparse CCA. In: *Music brain & cognition workshop, NIPS*.
- Turnbull, D., Barrington, L., Torres, D., & Lanckriet, G. R. G. (2008). Semantic annotation and retrieval of music and sound effects. *IEEE Transactions on Audio, Speech and Language Processing*, 16, 467–476.
- Vandenbergh, L., & Boyd, S. (1996). Semidefinite programming. *SIAM Review*, 38, 49–95.
- Vinokourov, A., Shawe-Taylor, J., & Cristianini, N. (2003). Inferring a semantic representation of text via cross-language correlation analysis. In S. Becker, S. Thrun, & K. Obermayer (Eds.), *Advances in neural information processing systems* (Vol. 15, pp. 1473–1480). Cambridge: MIT Press.
- Weston, J., Elisseeff, A., Schölkopf, B., & Tipping, M. (2003). Use of the zero-norm with linear models and kernel methods. *Journal of Machine Learning Research*, 3, 1439–1461.
- Yuille, A. L., & Rangarajan, A. (2003). The concave-convex procedure. *Neural Computation*, 15, 915–936.
- Zangwill, W. I. (1969). *Nonlinear programming: a unified approach*. Englewood Cliffs: Prentice-Hall.
- Zou, H., & Hastie, T. (2005). Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society, Series B*, 67, 301–320.
- Zou, H., Hastie, T., & Tibshirani, R. (2006). Sparse principal component analysis. *Journal of Computational and Graphical Statistics*, 15, 265–286.