



CHRISTOPH HELMBERG  
FRANZ RENDL

## A Spectral Bundle Method for Semidefinite Programming

# A Spectral Bundle Method for Semidefinite Programming

C. Helmberg\*      F. Rendl†

August 1997, revised October 1997

## Abstract

A central drawback of primal-dual interior point methods for semidefinite programs is their lack of ability to exploit problem structure in cost and coefficient matrices. This restricts applicability to problems of small dimension. Typically semidefinite relaxations arising in combinatorial applications have sparse and well structured cost and coefficient matrices of huge order. We present a method that allows to compute acceptable approximations to the optimal solution of large problems within reasonable time.

Semidefinite programming problems with constant trace on the primal feasible set are equivalent to eigenvalue optimization problems. These are convex nonsmooth programming problems and can be solved by bundle methods. We propose to replace the traditional polyhedral cutting plane model constructed by means of subgradient information by a semidefinite model that is tailored for eigenvalue problems. Convergence follows from the traditional approach but a proof is included for completeness. We present numerical examples demonstrating the efficacy of the approach on combinatorial examples.

**Key Words.** Eigenvalue optimization, convex optimization, semidefinite programming, proximal bundle method, large-scale problems.

**AMS Subject Classification 1991.** 65F15, 90C25; Secondary 52A41, 90C06.

## 1 Introduction

The development of interior point methods for semidefinite programming [19, 29, 1, 41] has increased interest into semidefinite modeling techniques in several fields such as control theory, eigenvalue optimization, and combinatorial optimization. In fact, interior point methods proved to be very useful and reliable solution methods for semidefinite programs of moderate size. However, if the problem is defined over large matrix variables or a huge number of constraints interior point methods grow terribly slow and consume huge amounts of memory. The most efficient methods of today [15, 22, 2, 30, 40, 27] are primal-dual methods that require, in each iteration of the interior point method, the factorization of a dense matrix of order equal to the number of constraints and one to three factorizations of the positive semidefinite matrix variables within the line search. For a typical workstation this restricts the number of constraints to 2000 and the size of the matrix variables to 500 if reasonable performance is required. For larger problems time and memory requirements are prohibitive. It is important to realize that either the primal or the dual matrix is

---

\*Konrad-Zuse-Zentrum für Informationstechnik Berlin, Takustraße 7, 14195 Berlin, Germany.  
[helmberg@zib.de](mailto:helmberg@zib.de), <http://www.zib.de/helmberg>

†Technische Universität Graz, Institut für Mathematik, Steyrergasse 30, A-8010 Graz, Austria.  
[rendl@opt.math.tu-graz.ac.at](mailto:rendl@opt.math.tu-graz.ac.at)

generically dense even if cost and coefficient matrices are very sparse. Very recently, a pure dual approach was proposed in [5] which offers some possibilities to exploit sparsity. It is too early to judge the potential of this method.

In combinatorial optimization semidefinite relaxations were introduced in [25]. At that time they were mainly considered a theoretical tool for obtaining strong bounds [12, 26, 36]. With the development of interior point methods hopes soared high that these relaxations could be of practical value. Within short time several approximation algorithms relying on semidefinite programming were published, one of the most prominent being [9]. On the implementational side [16, 14, 13] cutting plane approaches for semidefinite relaxations of constrained quadratic 0-1 programming problems proved to yield solutions of high quality. However, as mentioned above, they were very expensive to compute even for problems of small size (a few hundred 0-1 variables). Problems arising in practical applications (starting with a few thousand 0-1 variables) were out of reach. We believe that the method proposed in this paper will open the door to problems of this size.

Although combinatorial applications are our primary concern we stress that the method is not restricted to this kind of problems. In fact it will be a useful alternative to interior point methods whenever the number of constraints or the order of the matrices is quite large.

We transform a standard dual semidefinite program into an eigenvalue optimization problem by reformulating the semidefinite constraint as a non-negativity constraint on the minimal eigenvalue of the slack matrix variable and lifting this constraint into the cost function by means of a Lagrange multiplier. The correct value of the Lagrange multiplier is known in advance if the primal feasible matrices have constant trace. (This is the case for the combinatorial applications we have in mind.)

In this paper we develop a bundle method for solving the problem of minimizing the maximal eigenvalue of an affine matrix function with an additional linear objective term. These functions are well known to be convex and non-smooth. A very general method for optimizing non-smooth convex functions is the bundle method, see e.g. [20, 17, 18, 37]. In each step the function value and a subgradient of the function is computed for some specific point. By means of the collected subgradients a local cutting plane model of the function is formed. The minimizer of the cutting plane model augmented by a regularization term yields the new point. In the case of eigenvalue optimization the subgradient is formed by means of an eigenvector to the maximal eigenvalue. Extremal eigenvalues and associated eigenvectors of large symmetric matrices can be computed efficiently by Lanczos methods (see e.g. [10]). Lanczos methods need a subroutine that computes the product of the matrix with a vector. This allows to exploit any kind of structure present in the matrix.

The cutting plane model used in the bundle algorithm can be interpreted as being formed by a restricted set of  $\varepsilon$ -subgradients of the function in the current point. We will see that the  $\varepsilon$ -subdifferential of eigenvalue problems has the form of the feasible set of a semidefinite program. This suggests to use, instead of the traditional polyhedral cutting plane model, a semidefinite cutting plane model that is formed by restricting the feasible set of  $\varepsilon$ -subgradients to an appropriate face of the semidefinite cone. This specialization of the cutting plane model is the main contribution of the paper.

The semidefinite bundle approach allows for an intriguing interpretation in terms of the original semidefinite program. The cutting plane model requires that the dual slack matrix of the semidefinite program is positive semidefinite with respect to a subspace of vectors. In general the optimal solution of the relaxed semidefinite problem will produce an indefinite dual slack matrix. The negative eigenvalues and corresponding eigenvectors of the slack matrix are used to update the subspace and the process is iterated.

This process trivially provides the optimal solution if the subspace grows to the full space. However, we show that during the algorithm generically the dimension of the subspace is bounded by (roughly) the square root of the number of constraints. If this is still considered too large the introduction of an aggregate subgradient guarantees convergence even in case of one-dimensional subspaces.

In contrast to ‘classical’ algorithms [7, 32, 31] for eigenvalue optimization our method does not require the knowledge of the correct multiplicity of the maximal eigenvalue in the optimal solution. On the other hand [32, 31] show that their algorithms are quadratically convergent in the limit. A similar property cannot be expected to hold for our algorithm since it is a first order method only. In principle convergence follows from the traditional approach (see e.g. [20]) but we include a proof for completeness. We also present a primal-dual interior point code for solving the quadratic semidefinite programming problems associated with the semidefinite cutting plane models and discuss efficiency aspects. The properties of the algorithm are illustrated on several combinatorial examples.

In Section 2 some basic properties of semidefinite programs are stated. In Section 3 we transform semidefinite programs into eigenvalue optimization problems and derive their  $\varepsilon$ -subdifferential. Section 4 introduces the bundle method. The algorithm and the proof of convergence is given in Section 5. The quadratic semidefinite subproblems arising in the bundle method can be solved by interior point methods as explained in Section 6. Section 7 gives an outline of the implementation and briefly discusses the computation of the maximal eigenvalue and an associated eigenvector. Numerical examples on combinatorial problems are presented in Section 8. We conclude the paper with a summary and possible extensions and improvements in Section 9. For the convenience of the reader an appendix explaining the notation and the symmetric Kronecker product is included at the end of the paper.

## 2 Semidefinite Programs

We denote the set of symmetric matrices of order  $n$  by  $S_n$  which we regard as a space isomorphic to  $\mathbb{R}^{\binom{n+1}{2}}$ . The subset of positive semidefinite matrices  $S_n^+$  is a full-dimensional, non-polyhedral convex cone in  $S_n$  and defines a partial order on the symmetric matrices by  $A \succeq B \iff (A - B) \in S_n^+$ . Positive definite matrices are denoted by  $S_n^{++}$  or  $A > 0$ . Except for the apex  $\{0\}$ , the faces of the semidefinite cone  $S_n^+$  have dimension  $\binom{r+1}{2} \leq \binom{n+1}{2}$  for  $r \in \mathbb{N}$ , and are of the form

$$F = \{PVP^T : V \in S_r^+\},$$

where  $P \in M_{n,r}$  is a fixed matrix for each  $F$  and can be assumed to have orthonormal columns. All possible eigenvectors to non-zero eigenvalues of matrices in  $F$  are contained in the subspace spanned by the columns of  $P$ .

Consider the standard primal-dual pair of semidefinite programs,

$$(P) \quad \begin{array}{ll} \max & \langle C, X \rangle \\ \text{s.t.} & \mathcal{A}(X) = b \\ & X \succeq 0. \end{array} \quad (D) \quad \begin{array}{ll} \min & \langle b, y \rangle \\ \text{s.t.} & Z = \mathcal{A}^T(y) - C \\ & Z \succeq 0. \end{array}$$

$\mathcal{A} : S_n \rightarrow \mathbb{R}^m$  is a linear operator and  $\mathcal{A}^T : \mathbb{R}^m \rightarrow S_n$  is its adjoint operator, i.e. it satisfies  $\langle \mathcal{A}(X), y \rangle = \langle X, \mathcal{A}^T(y) \rangle$  for all  $X \in S_n$  and  $y \in \mathbb{R}^m$ . They are of the form

$$\mathcal{A}(X) = \begin{bmatrix} \langle A_1, X \rangle \\ \vdots \\ \langle A_m, X \rangle \end{bmatrix} \quad \text{and} \quad \mathcal{A}^T(y) = \sum_{i=1}^m y_i A_i$$

with  $A_i \in S_n$ ,  $i = 1, \dots, m$ .  $C \in S_n$  is the cost matrix,  $b \in \mathbb{R}^m$  the right-hand-side vector. (We assume tacitly some constraint qualification to hold, so that these problems satisfy strong duality.)

For any optimal solution  $X^*$  of (P) and any optimal solution  $(y^*, Z^*)$  of (D) we have

$$(1) \quad X^* Z^* = 0.$$

Thus  $X^*$  and  $Z^*$  are simultaneously diagonalizable by some orthonormal matrix  $P$ ,  $X^* = P\Lambda_{X^*}P^T$  and  $Z^* = P\Lambda_{Z^*}P^T$  with  $\Lambda_{X^*}\Lambda_{Z^*} = 0$ .

Pataki proved in [33] that for matrices  $X$  contained in a  $k$ -dimensional face of the primal feasible set the rank  $r$  is bounded by

$$\binom{r+1}{2} \leq m+k.$$

There is always an optimal solution in a 0-dimensional face of the primal feasible set, thus there is always an optimal solution  $X^*$  having rank  $r$  bounded by  $\binom{r+1}{2} \leq m$ . Usually the rank is even smaller, see [3]. It is also proved in [3] that in general the optimal solution is unique and therefore, in general, the rank of  $X^*$  satisfies the bound.

### 3 Eigenvalue Optimization

We reformulate the dual as an eigenvalue optimization problem by modeling the constraint  $Z \succeq 0$  as  $\lambda_{\max}(-Z) \leq 0$ . By (1) any optimal  $Z^*$  is singular unless  $X^* = 0$  is primal optimal. To exclude this case we assume that  $X = 0$  is not an optimal solution of (P). This holds for instance, if  $b \neq 0$ . In this case the dual becomes

$$\min_{\lambda_{\max}(C - \mathcal{A}^T(y))=0} \langle b, y \rangle.$$

We lift the equality constraint into the objective by means of a Lagrange multiplier  $a$ ,

$$(E) \quad \min_{y \in \mathbb{R}^m} a \lambda_{\max}(C - \mathcal{A}^T(y)) + \langle b, y \rangle.$$

In general the correct value for  $a$  is not known. Let us assume that

$$\mathcal{A}(X) = b \text{ implies } \text{tr}(X) = a_0,$$

for some constant  $a_0$ . It can easily be verified that this condition holds if and only if there exists a vector  $\alpha$  such that  $\mathcal{A}^T(\alpha) = I$ , i.e. the identity is in the range of  $\mathcal{A}^T$ . In this case  $a_0 = \text{tr}(X) = \langle \mathcal{A}^T(\alpha), X \rangle = \alpha^T b$ . We note that many semidefinite programs arising in practical applications possess this property.

Under this assumption, the correct multiplier  $a$  is  $a_0$ , the trace of any feasible  $X$ . To see this, observe that  $I = \mathcal{A}^T(\alpha)$  with  $0 \neq \alpha \in \mathbb{R}^m$  and let  $a = \alpha^T b$ . Therefore adding  $\lambda\alpha$  to the dual variables  $y$  shifts the eigenvalues of  $Z$  by  $\lambda$  and changes the cost function by  $\lambda\alpha$ . This proves the claim.

The eigenvalue problem (E) is a convex, non-smooth optimization problem. It is well studied in the literature but we will repeat the most important concepts to keep the paper self-contained. To this end we first study the function

$$g(X) = \lambda_{\max}(X).$$

$g(X)$  is differentiable if and only if the maximal eigenvalue has multiplicity one. When optimizing eigenvalue functions, the optimum is generically attained at matrices whose maximal eigenvalue has multiplicity larger than one. In this case one has to consider the subdifferential of  $g$  at  $X$ ,

$$\partial g(X) = \{W \succeq 0 : \langle W, X \rangle = \lambda_{\max}(X), \text{tr}(W) = 1\}$$

doesn't change  
the result of (E)

(see e.g. [32]). In particular, for any  $v \in \mathbb{R}^n$  belonging to the eigenspace of the maximal eigenvalue of  $X$ ,  $W = vv^T$  is contained in the subdifferential of  $g$  at  $X$ .

If, in an optimization problem, zero is contained in the subdifferential then the current point is optimal. However, the subdifferential provides only local information. A more stable concept is the  $\varepsilon$ -subdifferential; For a convex function  $f$  and  $X \in \text{dom}(f)$  a matrix  $W$  is called an  $\varepsilon$ -subgradient of  $f$  at  $X$  if

$$f(Y) \geq f(X) + \langle W, Y - X \rangle - \varepsilon \quad \forall Y \in S_n.$$

The set of all  $\varepsilon$ -subgradients at  $X$  is the  $\varepsilon$ -subdifferential at  $X$  and is denoted by  $\partial_\varepsilon f(X)$ . The  $\varepsilon$ -subdifferential of  $g(X)$  has been characterized in [42], see also [17].

**Lemma 3.1** Let  $X$  be a symmetric matrix of order  $n$  and  $\varepsilon > 0$ , then

$$\partial_\varepsilon \lambda_{\max}(X) = \{W : \langle W, X \rangle \geq \lambda_{\max}(X) - \varepsilon, \text{tr}(W) = 1, W \succeq 0\}.$$

Observe, that for  $W \succeq 0$  with  $\text{tr}(W) = 1$  the smallest  $\varepsilon$  such that  $W \in \partial_\varepsilon \lambda_{\max}(X)$  is determined by  $\varepsilon = \lambda_{\max}(X) - \langle W, X \rangle$ .

Standard rules can be used to extend this result to the eigenvalue problem in (E). For  $a\lambda_{\max}(C - A^T(y))$  we obtain (see [18], p. 117 and p.98)

$$\begin{aligned} \partial_\varepsilon a\lambda_{\max}(C - A^T(y)) &= \\ &= -a\mathcal{A}(\partial_{\frac{\varepsilon}{a}} \lambda_{\max}(C - A^T(y))) \\ &= \left\{ -a\mathcal{A}(W) : \langle W, C - A^T(y) \rangle \geq \lambda_{\max}(C - A^T(y)) - \frac{\varepsilon}{a}, \text{tr}(W) = 1, W \succeq 0 \right\}. \end{aligned}$$

After adding a linear term  $b^T y$  we end up with the  $\varepsilon$ -subdifferential of our function

$$f(y) := a\lambda_{\max}(C - A^T(y)) + b^T y.$$

$$\begin{aligned} \partial_\varepsilon f(y) &= \\ &= \partial_\varepsilon a\lambda_{\max}(C - A^T(y)) + \{b\} \\ &= \left\{ b - a\mathcal{A}(W) : \langle W, C - A^T(y) \rangle \geq \lambda_{\max}(C - A^T(y)) - \frac{\varepsilon}{a}, \text{tr}(W) = 1, W \succeq 0 \right\}. \end{aligned}$$

**Corollary 3.2** Let  $W \succeq 0$  with  $\text{tr}(W) = 1$ . Then

$$f(z) \geq a \langle W, C - A^T(y) \rangle + b^T y + \langle b - a\mathcal{A}(W), z - y \rangle$$

for all  $y, z \in \mathbb{R}^n$ .

**Proof.** The smallest  $\varepsilon$  such that  $W$  is contained in the  $\varepsilon$ -subdifferential  $\partial_\varepsilon f$  of  $f$  at  $y$  is

$$\varepsilon = a\lambda_{\max}(C - A^T(y)) - a \langle W, C - A^T(y) \rangle.$$

The statement now follows from the definition of the  $\varepsilon$ -subdifferential. ■

## 4 The Bundle Method

In this section we develop a new method for minimizing  $f(y) = a\lambda_{\max}(C - A^T(y)) + b^T y$ . We use two classical ingredients, the *proximal point idea*, and the *bundle concept*. The new contribution lies in the way that we derive the new iterate from the ‘bundle’ of subgradient information collected from previous iterates.

Let

$$F(W, y) := a \langle C - \mathcal{A}^T(y), W \rangle + \langle b, y \rangle.$$

The following formulation of  $f(y)$  as a semidefinite program is well known.

$$f(y) = \max\{F(W, y) : W \succeq 0, \text{tr}(W) = 1\}. \quad \leftarrow \text{f}(y) = a\lambda_{\max}(C - \mathcal{A}^T(y)) + b^T y$$

By restricting  $W$  to some subcone of  $S_n^+$ , we obtain a lower approximation of  $f$ . We start out with a local model of  $f$  at the current iterate  $\hat{y}$ . Let  $P^T(C - \mathcal{A}^T(\hat{y}))P = \text{Diag}(\lambda_1, \dots, \lambda_k)$  with  $P^T P = I_k$  represent the eigenvalue decomposition with respect to the  $k$  largest eigenvalues  $\lambda_1 \geq \dots \geq \lambda_k$  of  $C - \mathcal{A}^T(\hat{y})$ . Define

$$\hat{f}(y) := \max\{F(W, y) : W = PVP^T, V \succeq 0, \text{tr}V = 1\}.$$

Clearly  $f(\hat{y}) = \hat{f}(\hat{y})$  and  $f(y) \geq \hat{f}(y) \forall y$ . In view of Lemma 3.1, the matrices  $W$ , over which we optimize here, give rise to a nonpolyhedral inner approximation of the  $\varepsilon$ -subdifferential of  $f$  at  $\hat{y}$  for  $\varepsilon = (\lambda_1 - \lambda_k)/a$ . Alternatively, the function  $f$  may be interpreted as the maximal eigenvalue function with eigenvectors restricted to the span of the columns of  $P$ ,

$$\hat{f}(y) = a\lambda_{\max}(P^T(C - \mathcal{A}^T(y))P) + \langle b, y \rangle.$$

The proximal point idea consists in determining a new iterate  $\hat{y} + d^*$  from the current iterate  $\hat{y}$  by solving the following convex problem. (Here  $u > 0$  is some fixed real.)

$$\min_d \hat{f}(\hat{y} + d) + \frac{u}{2} \langle d, d \rangle$$

We note that this minimization corresponds to the Lagrangian relaxation of  $\min\{\hat{f}(\hat{y} + d) : \langle d, d \rangle \leq r^2\}$ . Thus we replace the original function  $f$  by its local minorant  $\hat{f}$  and minimize locally around  $\hat{y}$ . The parameter  $u$  controls (indirectly) the radius of the sphere around  $\hat{y}$ , over which we minimize.

Substituting the definition of  $\hat{f}$ , this problem is the same as

$$(2) \quad \min_d \max_{W=PVP^T, \text{tr}V=1, V \succeq 0} F(W, \hat{y} + d) + \frac{u}{2} \langle d, d \rangle.$$

Up to now we have used only the local information from the current point  $\hat{y}$ . In the general situation we have available a ‘bundle’ of vectors also from previous iterates. We collect all ‘interesting’ vectors in a matrix and denote by  $P$  an orthonormal basis of its span. We discuss below in detail, how we identify interesting vectors. The main point at the moment is the fact that instead of optimizing over all positive semidefinite matrices  $W$ , we constrain ourselves to a small subface.

To keep the dimension of the subface indeed small, we aggregate some information from previous iterates into a single matrix, once the number of columns of  $P$  exceeds some limit. This aggregation is the final ingredient of our approach. The aggregate subgradient can be thought of as arising from a general positive semidefinite matrix

$$\overline{W} \in S_n^+ \quad \text{with} \quad \text{tr}(\overline{W}) = 1.$$

Let us denote by  $\widehat{\mathcal{W}}$  the set of matrices  $W$  determining  $\hat{f}$ , the minorant to  $f$ . With the aggregate  $\overline{W}$  included, this set now reads

$$(3) \quad \widehat{\mathcal{W}} = \{\alpha \overline{W} + PVP^T : \alpha + \text{tr}V = 1, \alpha \geq 0, V \succeq 0\}.$$

Problem (2) can be simplified, because  $d$  is unconstrained. Note that

$$F(W, \hat{y} + d) = a \langle C - \mathcal{A}^T(\hat{y}), W \rangle + \langle b, \hat{y} \rangle + \langle b - a\mathcal{A}(W), d \rangle,$$

so  $\nabla_d F = b - a\mathcal{A}(W)$ . Therefore we get

$$\begin{aligned} & \min_d \max_{W \in \widehat{\mathcal{W}}} F(W, \hat{y} + d) + \frac{u}{2} \langle d, d \rangle \\ &= \max_{W \in \widehat{\mathcal{W}}, b - a\mathcal{A}(W) + ud = 0} F(W, \hat{y} + d) + \frac{u}{2} \langle d, d \rangle \quad d = \frac{1}{u} (a\mathcal{A}(W) - b) \\ &= \max_{W \in \widehat{\mathcal{W}}} a \langle C - \mathcal{A}^T(\hat{y}), W \rangle + \langle b, \hat{y} \rangle - \frac{1}{2u} \langle a\mathcal{A}(W) - b, a\mathcal{A}(W) - b \rangle. \end{aligned}$$

The first equality follows from interchanging min and max, using strong duality for this problem, and using first order optimality for the inner minimization with respect to  $d$ ,

$$(4) \quad d = \frac{1}{u} (a\mathcal{A}(W) - b).$$

The final problem is a semidefinite program with (concave) quadratic cost function. We will discuss in Section 6 how problems of this kind can be solved efficiently.

All the ingredients for the update have now been presented, so we summarize the main steps of one iteration. To be consistent with the notation of the algorithm given in Section 5, let us denote by  $x^k$  what was called  $\hat{y}$  so far. The algorithm may have to compute several trial points  $y^{k+1}, y^{k+2}, \dots$  while keeping the same  $x^k = x^{k+1} = \dots$  if progress is not considered satisfactory. For each  $y^{k+1}$  the function is evaluated and a subgradient (eigenvector) is computed. This information is added to  $\widehat{\mathcal{W}}^k$  to form an improved model  $\widehat{\mathcal{W}}^{k+1}$ . Therefore, we assume that the current ‘bundle’  $P^k = P$  contains an eigenvector of  $\lambda_1(C - \mathcal{A}^T(y^k))$  in its span ( $y^k$  may or may not be equal to  $x^k$ ). Other than that,  $P$  need only satisfy  $P^T P = I_r$ . The local minorant of  $f$  at  $x^k$  is denoted by  $\hat{f}^k$ :

$$\hat{f}^k(y) := \max_{W \in \widehat{\mathcal{W}}^k} F(W, y).$$

Here  $\widehat{\mathcal{W}}^k$  represents the current approximation to the  $\varepsilon$ - subdifferential, see (3). It will be convenient to introduce also the regularized version of  $\hat{f}^k$ . Expressing  $y$  through the current iterate,  $y = x^k + d$ , we define

$$f^k(y) := \hat{f}^k(y) + \frac{u}{2} \|y - x^k\|^2 = \hat{f}^k(x^k + d) + \frac{u}{2} \|d\|^2.$$

Now the new trial point  $y^{k+1}$  is obtained by minimizing  $f^k(x^k + d)$  with respect to  $d$ . As described above, this can be done in two steps. First, solve

$$(5) \quad \max_{W \in \widehat{\mathcal{W}}^k} a \langle C - \mathcal{A}^T(x^k), W \rangle + \langle b, x^k \rangle - \frac{1}{2u} \langle a\mathcal{A}(W) - b, a\mathcal{A}(W) - b \rangle,$$

yielding the (unique) maximizer  $W^* = \alpha^* \overline{W}^k + P^k V^*(P^k)^T$ . Next, use (4) to compute

$$(6) \quad y^{k+1} = x^k + \frac{1}{u} (a\mathcal{A}(W^*) - b).$$

We summarize some easy facts, which will be used in the convergence analysis of the following section.

$$\hat{f}^k(y^{k+1}) \leq f^k(y) \quad \forall y,$$

since  $y^{k+1}$  is minimizer of  $f^k$ . Because  $f^k(x^k) = \hat{f}^k(x^k)$  and  $\hat{f}^k$  minorizes  $f$  we obtain

$$(7) \quad \underline{f^k(y^{k+1})} \leq f^k(x^k) \leq f(x^k).$$

Next let

$$f^*_*(y) := F(W^*, y) + \frac{u}{2} \|y - x^k\|^2.$$

Using the Taylor series expansion of  $f_*^k$  around  $y^{k+1}$  and the optimality condition (4) it follows easily that

$$(8) \quad f_*^k(y) = f_*^k(y^{k+1}) + \frac{u}{2} \|y - y^{k+1}\|^2. \quad \text{Second-order term}$$

Finally, we explain in detail, how  $\alpha^*$  and  $V^*$  are used to update  $P^k$  and  $\bar{W}^k$  to start a new iteration.

Let  $Q\Lambda Q^T$  be an eigenvalue decomposition of  $V^*$ . Then the ‘important’ part of the spectrum of  $V^*$  is spanned by the eigenvectors associated with the ‘large’ eigenvalues of  $V^*$ . Thus the eigenvectors of  $Q$  are split into two parts  $Q = [Q_1 Q_2]$  (with corresponding spectra  $\Lambda_1$  and  $\Lambda_2$ ),  $Q_1$  containing as columns the eigenvectors associated to ‘large’ eigenvalues of  $V^*$  and  $Q_2$  containing the remaining columns. Now the next  $P^{k+1}$  is computed to contain  $P^k Q_1$  and the next aggregate matrix is defined to be

$$(9) \quad \bar{W}^{k+1} = \alpha^* \bar{W}^k + P^k Q_2 \Lambda_2 Q_2^T (P^k)^T / (\alpha^* + \text{tr}(\Lambda_2)).$$

By this strategy it is easy to provide enough space in  $P^{k+1}$  for new vectors to improve the current model. In particular at least one eigenvector to the maximal eigenvalue of  $C - \mathcal{A}^T(y^{k+1})$  has to be added to  $P^{k+1}$ .

By construction  $W^*$  is contained in  $\widehat{W}^{k+1}$ . Therefore the local model for the next iteration will satisfy

$$(10) \quad f^{k+1}(y) \geq f_*^k(y) \quad \forall y. \quad \widehat{W}^{k+1} = \{ \alpha^* \bar{W}^{k+1} + P^{k+1} V (P^{k+1})^T : \dots \}$$

## 5 Algorithm and Convergence Analysis

In the previous section we focused on the question of doing one iteration of the bundle method. Now we provide a formal description of the method and point out that except for the choice of the bundle, the nature of the subproblem, and some minor changes in parameters the algorithm and its proof are identical to the algorithm of Kiwiel as presented in [20]. To keep the paper self-contained we present and analyze a simplified algorithm for fixed  $u$ . We refer the reader to [20] for an algorithm with variable choice of  $u$ .

### Algorithm 5.1

**Input:** An initial point  $y^0 \in \mathbb{R}^m$ , a (normalized) eigenvector  $v^0$  to the maximal eigenvalue of  $C - \mathcal{A}^T(y^0)$ , an  $\varepsilon > 0$  for termination, an improvement parameter  $m_L \in (0, \frac{1}{2})$ , a weight  $u > 0$ , an upper bound  $R \geq 1$  on the number of columns of  $P$ .

1. (Initialization)  $k = 0$ ,  $x^0 = y^0$ ,  $P^0 = v^0$ ,  $\bar{W}^0 = v^0(v^0)^T$ .
2. (Direction finding) Solve (5) to get  $y^{k+1}$  from (6). Decompose  $V^*$  into  $V^* = Q_1 \Lambda_1 Q_1^T + Q_2 \Lambda_2 Q_2^T$  with  $\text{rank}(Q_1) \leq R - 1$ . Compute  $\bar{W}^{k+1}$  using (9).
3. (Evaluation) Compute  $\lambda_{\max}(C - \mathcal{A}^T(y^{k+1}))$  and an eigenvector  $v^{k+1}$ . Let the columns of  $P^{k+1}$  contain an orthonormal basis of the space spanned by the columns of  $P^k Q_1$  and  $v^{k+1}$ .
4. (Termination) If  $f(x^k) - \hat{f}^k(y^{k+1}) \leq \varepsilon$  then **stop**.
5. (Serious step) If

$$(11) \quad f(y^{k+1}) \leq f(x^k) - m_L(f(x^k) - \hat{f}^k(y^{k+1}))$$

then set  $x^{k+1} = y^{k+1}$ , continue with Step 7. Otherwise continue with Step 6.

6. (Null step) Set  $x^{k+1} = x^k$ .

7. Increase  $k$  by 1 and go to Step 2.

We prove convergence of the algorithm for  $\varepsilon = 0$ . First consider the case that the algorithm does not stop and only null steps occur after some iteration  $K$ .

**Lemma 5.2** If there is a  $K \geq 0$  such that (11) is violated for all  $k \geq K$ , then  $\lim_{k \rightarrow \infty} \hat{f}^k(y^{k+1}) = f(x^K)$  and  $0 \in \partial f(x^K)$ .

**Proof.** For convenience we set  $x = x^K = x^{K+1} = \dots$ . Using the relations (8), (10) and (7), we obtain for all  $k \geq K$

$$f_*^k(y^{k+1}) + \frac{u}{2} \|y^{k+2} - y^{k+1}\|^2 = f_*^k(y^{k+2}) \leq f^{k+1}(y^{k+2}) \leq f^{k+1}(x^{k+1}) \leq f(x).$$

Therefore the  $f^k(y^{k+1})$  converge to some  $f^* \leq f(x)$  and the  $y^k$  converge to some  $y^*$ . Let  $g^{k+1} = b - a\mathcal{A}(v^{k+1}(v^{k+1})^T)$  denote the computed gradient of  $f$  in  $y^{k+1}$  and observe that the linearization  $\bar{f}$  of  $f$  in  $y^{k+1}$

$$\begin{aligned} \bar{f}(y; y_{k+1}) &= f(y^{k+1}) + \langle g^{k+1}, y - y^{k+1} \rangle = \hat{f}^k(y^{k+1}) + \langle g^{k+1}, y - y^{k+1} \rangle \\ &\leq \hat{f}^{k+1}(y) \end{aligned}$$

*Supporting hyperplane*

minorizes  $\hat{f}^{k+1}$ , because  $(v^{k+1}(v^{k+1})^T) \in \widehat{\mathcal{W}}^{k+1}$ . Thus

$$\begin{aligned} 0 &\leq f(y^{k+1}) - \hat{f}^k(y^{k+1}) = \\ &= \bar{f}(y^{k+1}; y^{k+1}) - \hat{f}^k(y^{k+1}) \\ &= \bar{f}(y^{k+2}; y^{k+1}) - \hat{f}^k(y^{k+1}) - \langle g^{k+1}, y^{k+2} - y^{k+1} \rangle \\ &\leq \hat{f}^{k+1}(y^{k+2}) - \hat{f}^k(y^{k+1}) + \|g^{k+1}\| \cdot \|y^{k+2} - y^{k+1}\| \\ &= f^{k+1}(y^{k+2}) - f^k(y^{k+1}) - u \|y^{k+2} - x\|^2 + u \|y^{k+1} - x\|^2 + \\ &\quad \|g^{k+1}\| \cdot \|y^{k+2} - y^{k+1}\| = \delta \end{aligned}$$

The convergence of the  $f^k(y^{k+1})$  and the  $y^k$  and the fact that the gradient is locally bounded imply that the last term goes to zero for  $k \rightarrow \infty$ . Since (11) is violated for all  $k > K$ , for all  $\delta > 0$  there is an  $M \in \mathbb{N}$  such that for all  $k > M$

$$\begin{aligned} f(y^{k+1}) - \delta &\leq \hat{f}^k(y^{k+1}) \leq f(x). \quad \hat{f}^k(y^{k+1}) &> f(x^k) - m_L (\hat{f}^k(x^k) - \hat{f}^k(y^{k+1})) \\ \text{and, since (11) is violated for all } k > K & \quad \quad \quad = (1 - m_L) f(x^k) + m_L \hat{f}^k(y^{k+1}) \\ (1 - m_L) f(y^{k+1}) + m_L \delta &\geq f(y^{k+1}) - m_L \hat{f}^k(y^{k+1}) \geq (1 - m_L) f(x). \end{aligned}$$

Thus the sequences  $f(y^{k+1})$  and  $\hat{f}^k(y^{k+1})$  both converge to  $f(x)$ .  $y^{k+1}$  is the minimizer of the regularized function  $f^k$ . On the one hand this implies that  $y^{k+1} \rightarrow x$ . On the other hand 0 must be contained in the subgradient  $\partial f^k(y^{k+1}) = \partial \hat{f}^k(y^{k+1}) + u(y^{k+1} - x)$ , see (4). Therefore there is a sequence  $h^k \in \partial \hat{f}^k(y^{k+1})$  of subgradients converging to zero. The  $f^k$  minorize  $f$ , the  $\hat{f}^k(y^{k+1})$  converge to  $f(x)$  and the  $y^{k+1}$  converge to  $x$ , hence zero must be contained in  $\partial f(x)$ .  $\blacksquare$

We may concentrate on serious steps in the following. In order to simplify notation we will speak of  $x^k$  as the sequence generated by serious steps with all duplicates eliminated. By  $f^k$  (and the corresponding  $\hat{f}^k$ ) we will refer to the function whose minimization gives rise to  $x^{k+1}$ .

If the algorithm stops after a finite number of iterations it is not difficult to verify that the last  $x^k$  is indeed optimal. In case there is an infinite number of serious steps and  $f(x^k)$  diverges to minus infinity, then by the local boundedness of the gradients the  $x^k$  will diverge, as well, and there is nothing to show. Consider therefore the case of an infinite number of serious steps satisfying

$$(12) \quad f(x^k) > f(\tilde{x}) \quad \text{for some fixed } \tilde{x} \in \mathbb{R}^m \text{ and all } k.$$

**Lemma 5.3** *If (12) holds for an infinite number of serious steps then the  $x^k$  converge to a minimizer of  $f$ .*

**Proof.** First we prove the boundedness of the  $x^k$ . To this end denote by  $p^{k+1} \in \partial\hat{f}^k(x^{k+1})$  the subgradient arising from the optimal solution of the minimization problem for  $f^k$ ,

$$p^{k+1} = b - a\mathcal{A}(\bar{W}^{k+1})$$

and observe that by (4)

$$x^{k+1} - x^k = -\frac{p^{k+1}}{u}.$$

Since  $\hat{f}^k$  minorizes  $f$  we obtain

$$f(x^k) \geq f(\tilde{x}) \geq \hat{f}^k(x^{k+1}) + \langle p^{k+1}, \tilde{x} - x^{k+1} \rangle.$$

Therefore the distance of  $x^{k+1}$  to  $\tilde{x}$  can be bounded by

$$\begin{aligned} \|\tilde{x} - x^{k+1}\|^2 &= \|\tilde{x} - x^k + x^k - x^{k+1}\|^2 = \\ &\leq \|\tilde{x} - x^k\|^2 + 2\langle \tilde{x} - x^k, x^k - x^{k+1} \rangle + 2\langle x^k - x^{k+1}, x^k - x^{k+1} \rangle \\ &= \|\tilde{x} - x^k\|^2 + 2\langle \tilde{x} - x^{k+1}, p^{k+1}/u \rangle \\ &\leq \|\tilde{x} - x^k\|^2 + \frac{2}{u}(f(x^k) - \hat{f}^k(x^{k+1})). \end{aligned}$$

For any  $k > K$ , a recursive application of the bound above yields

$$(13) \quad \|\tilde{x} - x^k\|^2 \leq \|\tilde{x} - x^K\|^2 + \frac{2}{u} \sum_{i=K}^{\infty} (f(x^i) - \hat{f}^i(x^{i+1})).$$

By (11) the progress of the algorithm in each serious step is at least  $m_L(f(x^k) - \hat{f}^k(x^{k+1}))$  and together with (12) we obtain

$$\sum_{i=0}^{\infty} (f(x^i) - \hat{f}^i(x^{i+1})) \leq \frac{1}{m_L} (f(x^0) - f(\tilde{x})).$$

Therefore the sequence of the  $x^k$  remains bounded and has an accumulation point  $\bar{x}$ . By replacing  $\tilde{x}$  by  $\bar{x}$  in (13) and choosing  $K$  sufficiently large, the remaining sum can be made smaller than an arbitrary small  $\delta > 0$ , thus proving the convergence of the  $x^k$  to  $\bar{x}$ . As the  $x^{k+1}$  converge to  $\bar{x}$  the  $p^{k+1}$  converge to zero by (4), and since the sequence  $(f(x^k) - \hat{f}^k(x^{k+1}))$  has to converge to zero as well, we conclude that  $0 \in \partial f(\bar{x})$ , i.e.,  $\bar{x}$  is a minimizer of  $f$ .  $\blacksquare$

**Remark 5.4** *We have just seen that the bundle algorithm converges, even if  $P$  contains only one column. In this case the use of the aggregate subgradient is crucial.*

*To achieve convergence of the bundle algorithm without aggregate subgradients, it suffices to store in  $P$  only the subspace spanning the eigenvectors corresponding to non-zero eigenvalues of an optimal solution  $W^*$  of (2). Is there a reasonable upper bound on the number of columns we will need?*

*By the uniqueness of  $d$  it follows from (4) that for all optimal solutions  $W^*$  of (2)  $\bar{b} = \mathcal{A}(W^*)$  must be constant. All optimal solutions of (2) are also optimal solutions of the following semidefinite program (here we do not require the identity  $I$  to be in the span of  $A_1$  to  $A_m$ ),*

$$\begin{aligned} \max \quad & a \langle PVPT, C - \mathcal{A}^T(x^k) \rangle \\ \text{s.t.} \quad & \mathcal{A}(PVPT) = \bar{b} \\ & \text{tr}(V) = 1 \\ & V \succeq 0. \end{aligned}$$

Thus, independent of the rank of  $P$ , we can always find an optimal solution  $W^*$  to (2) whose rank  $r$  is bounded by  $\binom{r+1}{2} \leq m+1$  (and, of course, by  $n$ ). Therefore the maximal number of columns one has to provide to achieve convergence of the bundle method without aggregation is  $\bar{r}$  with  $\binom{\bar{r}+1}{2} \leq m+1$  plus the number of eigenvectors to be added in each iteration (this is at least one).

## 6 Solving the Subproblem

In this section we concentrate on how the minimizer of  $f^k$  can be computed efficiently. We have already seen in Section 4 that this task is equivalent to solving the quadratic semidefinite program (2) or (5). Problems of this kind can be solved by interior point methods, see e.g. [8, 22]. In practice it is useful to implement separate solvers for (2) and (5). Here, we will explain the algorithm for (5) only. Dropping the constants in (5) we obtain for  $y = x^k$

$$\begin{aligned} \min \quad & \frac{a^2}{2u} \langle \mathcal{A}(W), \mathcal{A}(W) \rangle - \frac{2a}{2u} \langle b, \mathcal{A}(W) \rangle - a \langle C - \mathcal{A}^T(y), W \rangle \\ \text{s.t.} \quad & W = \alpha \bar{W} + P V P^T \\ & \alpha + \text{tr}(V) = 1 \\ & \alpha \geq 0, V \succeq 0. \end{aligned}$$

Expanding  $W = \alpha \bar{W} + P V P^T$  into the cost function yields

$$\begin{aligned} & \langle \mathcal{A}(P V P^T), \mathcal{A}(P V P^T) \rangle \\ &= \sum_{i=1}^k \langle A_i, P V P^T \rangle^2 \quad \text{min} \quad \frac{a^2}{2u} [\langle \mathcal{A}(P V P^T), \mathcal{A}(P V P^T) \rangle + 2\alpha \langle \mathcal{A}(P V P^T), \mathcal{A}(\bar{W}) \rangle + \alpha^2 \langle \mathcal{A}(\bar{W}), \mathcal{A}(\bar{W}) \rangle] \\ &= \sum_{i=1}^k \langle P^T A_i P, V \rangle^2 \quad \text{s.t.} \quad -a \left\langle \frac{1}{u} \mathcal{A}^T(b) + C - \mathcal{A}^T(y), P^T V P \right\rangle - \alpha a \left\langle \frac{1}{u} \mathcal{A}^T(b) + C - \mathcal{A}^T(y), \bar{W} \right\rangle \\ &= \sum_{i=1}^k \text{svec}(V)^T \text{svec}(P^T A_i P) \text{svec}(P^T A_i P) \text{svec}(V) \quad \alpha \geq 0, V \succeq 0. \end{aligned}$$

similar to rec, but it is more compact, only for symmetric matrices.

Using the svec-operator (see the appendix for a definition and important properties of svec and the symmetric Kronecker product  $\otimes_s$ ) to expand symmetric matrices from  $S_r$  into column vectors of length  $\binom{r+1}{2}$  we obtain the quadratic program (recall that, for  $A, B \in S_r$ ,  $\langle A, B \rangle = \text{svec}(A)^T \text{svec}(B)$  and that  $\text{tr}(V) = \langle I, V \rangle$ )

$$(14) \quad \begin{aligned} \min \quad & \frac{1}{2} \text{svec}(V)^T Q_{11} \text{svec}(V) + \alpha q_{12}^T \text{svec}(V) + \frac{1}{2} q_{22} \alpha^2 + c_1^T \text{svec}(V) + c_2 \alpha \\ \alpha + s_I^T \text{svec}(V) = 1 \quad & \\ \alpha \geq 0, V \succeq 0, \quad & \end{aligned}$$

QP w.r.t. svec(V)

where (after some technical linear algebra)

$$\begin{aligned} & \sum_{i=1}^k \langle A_i, P V P^T \rangle \langle A_i, \bar{W} \rangle \quad (15) \quad Q_{11} = \frac{a^2}{u} \sum_{i=1}^m \text{svec}(P^T A_i P) \text{svec}(P^T A_i P)^T \\ &= \sum_{i=1}^k \langle A_i, \bar{W} \rangle \langle P^T A_i P, V \rangle \quad (16) \quad q_{12} = \frac{a^2}{u} \text{svec}(P^T \mathcal{A}^T(\mathcal{A}(\bar{W})) P) \\ &= \sum_{i=1}^k \langle P^T A_i \langle A_i, \bar{W} \rangle P, V \rangle \quad (17) \quad q_{22} = \frac{a^2}{u} \langle \mathcal{A}(\bar{W}), \mathcal{A}(\bar{W}) \rangle \\ & \quad (18) \quad c_1 = -a \text{svec}(P^T (\frac{1}{u} \mathcal{A}^T(b) + C - \mathcal{A}^T(y)) P) \\ & \quad (19) \quad c_2 = -a (\langle \frac{1}{u} b - y, \mathcal{A}(\bar{W}) \rangle + \langle C, \bar{W} \rangle) \\ &= \langle P^T \mathcal{A}^T(\mathcal{A}(\bar{W})) P, V \rangle \quad s_I = \text{svec}(I) \end{aligned}$$

for  $\bar{W}$ , only need  $\mathcal{A}(\bar{W})$  and  $\langle C, \bar{W} \rangle$

No necessity to really store  $\bar{W}$ .

At this point it is advisable to spend some thought on  $\bar{W}$ . The algorithm is designed for very large and sparse cost matrices  $C$ .  $\bar{W}$  is of the same size as  $C$ . Initially it might be possible to exploit the low rank structure of  $\bar{W}$  for efficient representations, but as the algorithm proceeds, the rank of  $\bar{W}$  grows inevitably. Thus it is impossible to store all the information of  $\bar{W}$ . However, as we can see in (15) to (19), it suffices

to have available the vector  $\mathcal{A}(\bar{W}) \in \mathbb{R}^m$  and the scalar  $\langle C, \bar{W} \rangle$  to construct the quadratic program. Furthermore, by the linearity of  $\mathcal{A}(\cdot)$  and  $\langle C, \cdot \rangle$ , these values are easily updated whenever  $\bar{W}$  is changed.

To solve (14) we employ a primal-dual interior point strategy. To formulate the defining equations for the central path we introduce a Lagrange multiplier  $t$  for the equality constraint, a dual slack matrix  $U \succeq 0$  as complementary variable to  $V$ , a dual slack scalar  $\beta \geq 0$  as complementary variable to  $\alpha$  and a barrier parameter  $\mu > 0$ . The system reads

$$\begin{aligned}\frac{\partial L}{\partial s \text{vec}(V)} = F_U &= Q_{11} \text{svec}(V) + \alpha q_{12} + c_1 - ts_I - \text{svec}(U) = 0 \\ \frac{\partial L}{\partial \alpha} = F_\beta &= \alpha q_{22} + q_{12}^T \text{svec}(V) + c_2 - t - \beta = 0 \\ \frac{\partial L}{\partial t} = F_1 &= 1 - \alpha - s_I^T \text{svec}(V) = 0 \\ UV &= \mu I \\ \alpha \beta &= \mu\end{aligned}$$

$$G \otimes_s H = \frac{1}{2} Q^T (G \otimes H + H \otimes G) Q$$

The step direction  $(\Delta\alpha, \Delta\beta, \Delta U, \Delta V, \Delta t)$  is determined via the linearized system

$$\begin{aligned}Q \in \mathbb{R}^{\frac{1}{2}N(n+1) \times r^2} \\ Q \text{ vec}(s) = \text{vec}(S) \\ Q \otimes_s H = \frac{1}{2} Q^T (G \otimes H + H \otimes G) Q \\ Q_{11} \text{svec}(\Delta V) + \Delta \alpha q_{12} - \Delta t s_I - \text{svec}(\Delta U) &= -F_U \quad Q_{11} \text{svec}(V) + \alpha q_{12} + c_1 - t s_I - \text{svec}(U) = F_u \\ q_{22} \Delta \alpha + q_{12}^T \text{svec}(\Delta V) - \Delta t - \Delta \beta &= -F_\beta \quad \Delta V = V^* - V \\ -\Delta \alpha - s_I^T \text{svec}(\Delta V) &= -F_1 \quad \Delta U = U^* - U\end{aligned}$$

$$\boxed{\Delta V + \Delta U V = 0}$$

$$(U \otimes_s V^{-1}) \text{svec}(\Delta V) + \text{svec}(\Delta U) = \mu \text{svec}(V^{-1}) - \text{svec}(U)$$

$$U \otimes_s V^{-1} + \Delta U = 0$$

In the current context we prefer the linearization  $(U \otimes_s V^{-1}) \text{svec}(\Delta V) + \text{svec}(\Delta U)$  because it makes the system easy to solve for  $\Delta V$  with relatively little computational work per iteration. The final system for  $\Delta V$  reads

$$\begin{aligned}(20) (Q_{11} + U \otimes_s V^{-1} + (\frac{\beta}{\alpha} + q_{22}) s_I s_I^T - q_{12} s_I^T - s_I q_{12}^T) \text{svec}(\Delta V) &= -F_\beta \\ &= \mu \text{svec}(V^{-1}) - \text{svec}(U) - F_U - F_1 q_{12} - (\mu \alpha^{-1} - \beta - \frac{\beta}{\alpha} F_1 - q_{22} F_1) s_I \\ \textcircled{1} + \textcircled{2} s_I + \textcircled{3} q_{12} + \textcircled{4} - \textcircled{5} s_I - \textcircled{6} \frac{\beta}{\alpha} s_I - \textcircled{7} q_{22} s_I &\end{aligned}$$

It is not too difficult to see that the system matrix is positive definite (it suffices to show that  $Q_{11} + q_{22} s_I s_I^T - q_{12} s_I^T - s_I q_{12}^T \succ 0$  using  $\begin{bmatrix} Q_{11} & q_{12} \\ q_{12} & q_{22} \end{bmatrix} \succ 0$ ). The main

work per iteration is the factorization of this matrix (with  $V \in S_r$  this is  $\binom{r+1}{2}^3 / 3$ ) and it is not possible to do much better since  $Q_{11}$  has to be inverted at some point. Because of the strong dominance of the factorization it pays to employ a predictor corrector approach, but we will not delve into this here.

For  $V \in S_r$ , a strictly feasible primal starting point is

$$\begin{aligned}V^0 &= I/(r+1) \\ \alpha^0 &= 1/(r+1)\end{aligned}$$

and a strictly feasible dual starting point can be constructed by choosing  $t^0$  sufficiently negative such that

$$\begin{aligned}U^0 &= \text{svec}^{-1}(Q_{11} \text{svec}(V^0) + \alpha^0 q_{12} + c_1) - t^0 I \succ 0 \\ \beta^0 &= \alpha q_{22} + q_{12}^T \text{svec}(V) + c_2 - t > 0.\end{aligned}$$

Starting from this strictly feasible primal-dual pair we compute a first  $\mu = (\langle U, V \rangle + \alpha \beta) / (r+1)$ , compute the step direction  $(\Delta\alpha, \Delta\beta, \Delta U, \Delta V, \Delta t)$  as indicated above, perform a line search with line search parameter  $0 < \delta \leq 1$  such that

$(\alpha + \delta\Delta\alpha, \beta + \delta\Delta\beta, U + \delta\Delta U, V + \delta\Delta V, \delta\Delta t)$  is again strictly feasible, move to this new point, compute a new  $\mu$  by

$$\mu = \min \left\{ \mu_{\text{old}}, \gamma \frac{\langle U, V \rangle + \alpha\beta}{r+1} \right\} \quad \text{with } \gamma = \begin{cases} 1 & \text{if } \delta \leq \frac{1}{5} \\ \frac{5}{10} - \frac{4}{10}\delta^2 & \text{if } \delta > \frac{1}{5}, \end{cases}$$

and iterate. We stop if  $(\langle U, V \rangle + \alpha\beta)/(r+1) < 10^{-7}$ .

## 7 Implementation

In our implementation of the algorithm we largely follow the rules outlined in [20]. In particular  $u$  is adapted during the algorithm. The first guess for  $u$  is equal to the norm of the first subgradient determined by  $v^0$ . The scheme for adapting  $u$  is the same as in [20] except for a few changes in parameters. For example the parameter  $m_L$  for accepting a step as serious is set to  $m_L = 0.2$  and the parameter  $m_R$  indicating that the model is so good (progress by the serious step is larger than  $m_R[f(x^k) - \hat{f}^k(y^{k+1})]$ ) that  $u$  can be decreased is set to  $m_R = 0.7$ .

The stopping criterion is formulated in relative precision,

$$f(x^k) - \hat{f}^k(y^{k+1}) \leq \varepsilon \cdot (|f(x^k)| + 1)$$

with  $\varepsilon = 10^{-5}$  in the implementation.

The choice of the upper bound  $R$  on the number of columns  $r$  of  $P$  and the selection of the subspace merits some additional remarks. Observe that by the considerations of Section 4 it is highly unlikely the  $r$  violates the bound  $\binom{r+1}{2} \leq m$  even if the number of columns of  $P$  is not restricted.  $\binom{r+1}{2}$  is also the order of the system matrix in (20) and is usually considerably smaller than the size of the system matrix in traditional interior point codes for semidefinite programming which is of order  $m$ . Furthermore the order of the matrix variables is  $r$  as compared to  $n$  for traditional interior point codes. Thus if the number of constraints  $m$  is roughly of the same size as  $n$  and a matrix of order  $m$  is still considered factorizable then running the algorithm without bounding the number of columns of  $P$  may turn out to be considerably faster than running an interior point method. This can be observed in practice, see Section 8.

For huge  $n$  and  $m$  primal-dual interior point methods are not applicable any more, because  $X$ ,  $Z^{-1}$ , and the system matrix are dense. In this case the proposed bundle approach allows to apply the powerful interior point approach at least on an important subspace of the problem. The correct identification of the relevant subspace in  $V^*$  is facilitated by the availability of the complementary variable  $U^*$ .  $U^*$  helps to discern between the small eigenvalues of  $V^*$  (because of the interior point approach we have  $V^* \succ 0!$ ). Eigenvectors  $v$  of  $V^*$  that are of no importance for the optimal solution of the subproblem will have a large value  $v^T U^* v$ , whereas eigenvectors, that are ambiguous, will have both, a small eigenvalue  $v^T V^* v$  and a small value  $v^T U^* v$ .

In practice we restrict the number of columns of  $P$  to 25 and provide room for at least five new vectors in each iteration (see below). Eigenvectors  $v$  that correspond to small but important eigenvalues  $\lambda$  of  $V^*$  ( $\lambda < 10^{-3} \lambda_{\max}(V^*)$  and  $\lambda > 10^{-2} v^T U^* v$ ) are added to  $\overline{W}$ , the smallest important eigenvectors ( $\lambda > 10^{-3} \lambda_{\max}(V^*)$ ) are added to  $\overline{W}$  only if more room is needed for new vectors.

For large  $m$  the computation of (15) to (19) is quite involved. A central object appearing in all constants is the projection of the constraint  $A_i$  on the space spanned by  $P$ ,  $P^T A_i P$ . Since the  $A_i$  are of the same size as  $X$  which we assume to be huge, it is important to exploit whatever structure is present in  $A_i$  to compute this projection efficiently. In combinatorial applications the  $A_i$  are of the form  $vv^T$  with

$v$  sparse and the projection can be computed efficiently. In the projection step and in particular in forming  $Q_{11}$  the size of  $r$  is again of strong influence. If we neglect the computation of  $\text{svec}(P^T A_i P)$ , the computation of  $Q_{11}$  still requires  $2m\binom{r+1}{2}+1$  flops. Indeed, if  $m$  is large then for small  $r$  the construction of  $Q_{11}$  takes longer than solving the associated quadratic semidefinite program.

The large computational costs involved in the construction and solution of the semidefinite subproblems may lead to the conviction that this model may not be worth the trouble. However, the evaluation of the eigenvalue-function is in fact much more expensive. There has been considerable work on computing eigenvalues of huge, sparse matrices, see e.g. [10] and the references therein. For extremal eigenvalues there seems to be a general consensus, that Lanczos type methods work best. Iterative methods run into difficulties if the eigenvalues are not well separated. In our context it is to be expected that in the course of the algorithm the  $\bar{r}$  largest eigenvalues will get closer and closer till all of them are identical in the optimum. For reasonable convergence block Lanczos algorithms with blocksize corresponding to the largest multiplicity of the eigenvalues have to be employed. During the first ten iterations the largest eigenvalue is usually well separated and the algorithm is fast. But soon the eigenvalues start to cluster, larger and larger blocksizes have to be used, and the eigenvalue problem gets more and more difficult to solve. In order to reduce the number of evaluations it seems worth to employ powerful methods in the cutting plane model. The increase in computation time required to solve the subproblem goes hand in hand with the difficulty of the eigenvalue problem because of the correspondence of the rank of  $P$  and the number of clustered eigenvalues.

Iterative methods for computing maximal eigenvectors generically offer approximate eigenvectors to several other large eigenvalues, as well. The space spanned by these approximate eigenvectors is likely to be a good approximation of the final optimal eigenspace. If the maximal number of columns for  $P$  is not yet attained it may be worth to include several of these approximate eigenvectors as well.

In our algorithm we use a block Lanczos code of our own that is based on a Fortran code of Hua (we guess that this is Hua Dai of [43]). It works with complete orthogonalization and employs Chebyshev iterations for acceleration. The choice of the blocksize is based on the approximate eigenvalues produced by previous evaluations but is at most 30. Four block Lanczos steps are followed by twenty Chebyshev iterations. This scheme is repeated till the maximal eigenvalue is found to the required relative precision. The relative precision depends on the distance of the maximal to the second largest eigenvalue but is bounded by  $10^{-6}$ . As starting vectors we use the complete block of eigenvectors and Lanczos-vectors from the previous evaluation.

## 8 Combinatorial Applications

The combinatorial problem we investigate is quadratic programming in  $\{-1, 1\}$  variables,

$$(\text{MC}) \quad \max x^T C x \quad \text{s.t. } x \in \{-1, 1\}^n.$$

In the case that  $C$  is the Laplace matrix of a (possible weighted) graph the problem is known to be equivalent to the max-cut problem.

The standard semidefinite relaxation is derived by observing that  $x^T C x = \langle C, xx^T \rangle$ . For all  $\{-1, 1\}^n$  vectors,  $xx^T$  is a positive semidefinite matrix with all diagonal elements equal to one. We relax  $xx^T$  to  $X \succeq 0$  and  $\text{diag}(X) = e$  and

obtain the following primal-dual pair of semidefinite programs,

$$\begin{array}{ll} \max & \langle C, X \rangle \\ (\text{PMC}) & \text{s.t. } \text{diag}(X) = e \\ & X \succeq 0 \end{array} \quad \begin{array}{ll} \min & e^T u \\ (\text{DMC}) & \text{s.t. } C + Z - \text{Diag}(u) = 0 \\ & Z \succeq 0. \end{array}$$

For non-negatively weighted graphs a celebrated result of Goemans and Williamson [9] says, that there is always a cut within .878 of the optimal value of the relaxation.

One of the first attempts to approximate (DMC) using eigenvalue optimization is contained in [35]. The authors use the Bundle code from Schramm and Zowe [37] with a limited number of bundle iterations, and so do not solve (DMC) exactly. So far the only practical algorithms for computing the optimal value were primal-dual interior point algorithms. However these are not able to exploit the sparsity of the cost function and have to cope with dense matrices  $X$  and  $Z^{-1}$ . An alternative approach by a combination of the power method with a generic optimization scheme of Plotkin, Shmoys, and Tardos [34] was proposed by [21] but seems to be purely theoretical.

In Table 1 we compare the proposed bundle method to our semidefinite code from [14] for graphs on  $n = m = 800$  nodes that were generated by **rudy**, a machine independent graph generator written by G. Rinaldi. Table 7 contains the command lines specifying the graphs. Graphs  $G_1$  to  $G_5$  are unweighted random graphs with a density of 6% (approx. 19000 edges).  $G_6$  to  $G_{10}$  are the same graphs with random edge weights from  $\{-1, 1\}$ .  $G_{11}$  to  $G_{13}$  are toroidal grids with random edge weights from  $\{-1, 1\}$  (approx. 1200 edges).  $G_{14}$  to  $G_{17}$  are unweighted ‘almost’ planar graphs having as edge set the union of two (almost maximal) planar graphs (approx. 4500 edges).  $G_{18}$  to  $G_{21}$  are the same almost planar graphs with random edge weights from  $\{-1, 1\}$ . In all cases the cost matrix  $C$  is the Laplace matrix of the graph, i.e., let  $A$  denote the (weighted) adjacency matrix of  $G$ , then

$$C = \text{Diag}(Ae) - A.$$

For a description of the semidefinite interior point code see [14], the termination criterion requires the gap between primal and dual optimal solution to be closed to a relative accuracy of  $5 \cdot 10^{-6}$ .

For the bundle algorithm the diagonal of  $C$  is removed so that in fact the algorithm works on the matrix  $C = \text{Diag}(\text{diag}(A)) - A$ . This does not change the problem because the diagonal elements of  $X$  are fixed to one. The offset  $e^T(Ae - \text{diag}(A))$  is added to the output only and has no influence on the algorithm whatsoever, in particular it has no influence on the stopping criterion. As starting vector  $y^0$  we choose the zero vector. All other parameters are as described in Section 7.

All computation times, for the interior point code as well as for the bundle code, refer to the same machine, a Sun sparc Ultra 1 with a Model 140 UltraSPARC CPU and 64 MB RAM. The time measured is the user time and it is given in the format  $hh : mm : ss$ , hours:minutes:seconds, leading zeros are dropped.

The first column of Table 1 identifies the graphs. The second and third refer to the interior point code and contain the optimal objective value produced (these can be regarded as highly accurate solutions) and the computation time. The fourth and fifth column give the same numbers for the bundle code.

On these examples the bundle code is superior to the interior point code. Although the examples do belong to the favorable class of instances having small  $m$  and relatively large  $n$ , the difference in computation time is astonishing. Note that the termination criterion used in the bundle code is quite accurate, except for  $G_{11}$  which seems to be a difficult problem for the bundle method. This deviation in accuracy is *not* caused by cancellations in connection with the offset. The difficulty of an example does not seem to depend on the number of nonzeros but rather on

the shape of the objective function. For toroidal grid graphs the maximum cut is likely to be not unique, thus the objective function will be rather flat. This flatness has its effect on the distribution of the eigenvalues in the optimal solution. Indeed, for  $G_{11}$  more eigenvalues cluster around the maximal eigenvalue than for the other problems. We illustrate this in Table 2, which gives the 30 largest eigenvalues of the solution at termination for problems  $G_1$ ,  $G_6$ ,  $G_{11}$ ,  $G_{14}$ , and  $G_{15}$ .

	IP-sol	IP-time	B-sol	B-time
$G_1$	12083.20	1:18:42	12083.22	4:11
$G_2$	12089.43	1:19:14	12089.45	5:19
$G_3$	12084.33	1:25:30	12084.34	4:38
$G_4$	12111.45	1:23:16	12111.46	3:37
$G_5$	12099.89	1:27:09	12099.91	4:38
$G_6$	2656.16	1:24:53	2656.18	3:57
$G_7$	2489.26	1:32:34	2489.29	7:54
$G_8$	2506.93	1:21:47	2506.95	3:38
$G_9$	2528.73	1:30:36	2528.75	3:51
$G_{10}$	2485.06	1:24:30	2485.08	3:56
$G_{11}$	629.16	1:28:41	629.21	45:26
$G_{12}$	623.87	1:34:55	623.89	31:14
$G_{13}$	647.13	1:37:52	647.14	18:44
$G_{14}$	3191.57	2:05:24	3191.58	14:30
$G_{15}$	3171.56	2:25:53	3171.56	23:20
$G_{16}$	3175.02	2:18:21	3175.04	16:59
$G_{17}$	3171.32	2:13:10	3171.35	16:32
$G_{18}$	1166.01	2:58:22	1166.02	18:32
$G_{19}$	1082.01	3:07:58	1082.04	12:27
$G_{20}$	1111.39	3:12:41	1111.40	11:59
$G_{21}$	1104.28	3:13:53	1104.29	13:35

Table 1: Comparison of the interior point (IP) and the bundle (B) approach. *sol* gives the computed solution value and *time* gives the computation time.

Table 3 provides additional information on the performance of the bundle algorithm on the examples of Table 1. The second column gives the accumulated time spent in the eigenvalue computation, it accounts for roughly 90% of the computation time. *serious* displays the number of serious steps, *iter* gives the total number of iterations including both, serious and null steps.  $\|g\|$  is the norm of the subgradient arising from the last optimal  $W^*$  before termination. For  $G_{11}$  the norm is considerably higher than for all other examples. Since the desired accuracy was not achieved for  $G_{11}$  by the standard stopping criterion it may be worth to consider an alternative stopping criterion taking into account the norm of the subgradient as well. Column *max-r* gives the maximal rank of  $P$  attained over all iterations. The rank of  $P$  would have been bounded by 25, but this bound never came into effect for any of these examples. Aggregation was not necessary. Observe that the theoretic bound allows for  $r$  up to 39, yet the maximal rank is only half this number. The last column gives the time when the objective value was first within  $10^{-3}$  of the optimum.

For combinatorial applications high accuracy of the optimal solution is of minor importance. An algorithm should deliver a reasonable bound fast and its solution should provide some hint on how a good feasible solution can be constructed. The bundle algorithm offers both. With respect to computation time the bundle algorithm displays the usual behavior of subgradient algorithms. Initially progress

	$G_1$	$G_6$	$G_{11}$	$G_{14}$	$G_{18}$
1	3.1190	3.2240	0.7653	1.0557	1.4175
2	.	3.2239	.	.	.
3	.	.	.	.	.
4	.	.	.	.	.
5	.	.	.	.	.
6	.	.	.	.	.
7	.	.	0.7653	.	.
8	.	.	0.7652	.	.
9	.	.	0.7652	.	.
10	.	.	0.7652	.	1.4175
11	.	.	0.7652	.	1.4155
12	.	.	0.7651	.	1.4090
13	3.1190	3.2239	0.7651	1.0557	1.4047
14	3.1135	3.2181	0.7650	1.0515	1.4007
15	3.1020	3.1968	0.7650	1.0500	1.3905
16	3.0928	3.1774	0.7649	1.0490	1.3834
17	3.0772	3.1556	0.7649	1.0450	1.3814
18	3.0594	2.7886	0.7648	1.0432	1.3798
19	2.7214	2.7716	0.7647	1.0398	1.3725
20	2.6964	2.7681	0.7646	1.0379	1.3709
21	2.6858	2.7269	0.7645	1.0358	1.3652
22	2.6834	2.6756	0.7644	1.0341	1.3583
23	2.6682	2.5851	0.7641	1.0331	1.3555
24	2.6649	2.5239	0.7639	1.0284	1.3510
25	2.6468	2.2357	0.7638	1.0266	1.3495
26	2.6274	1.8722	0.7636	1.0239	1.3480
27	2.6035	1.8473	0.7634	1.0231	1.3417
28	2.5137	1.7974	0.7633	1.0211	1.3397
29	2.4840	1.4859	0.7630	1.0180	1.3345
30	2.3281	1.4411	0.7629	1.0175	1.3291

Table 2: The 30 maximal eigenvalues after termination of examples  $G_1$ ,  $G_6$ ,  $G_{11}$ ,  $G_{14}$ , and  $G_{18}$ .

is very fast, but as the bound approaches the optimum there is a strong tailing off effect. We illustrate this by giving the objective values and computation times for the serious steps of example  $G_6$  (the diagonal offset is +77 in this example) in Table 4. After one minute the bound is within 0.1% of the optimum. For the other examples see the last column of Table 3.

With respect to a primal feasible solution observe that  $P^k V_*^k (P^k)^T$  is an successively better and better approximation to the primal optimal solution  $X^*$ . In case too much information is stored in the aggregate vector  $\mathcal{A}(\bar{W}^k)$  (remember that it is not advisable to store  $\bar{W}^k$  itself),  $P^k$  may be enriched with additional Lanczos-vectors from the eigenvalue computation. The solution of this enlarged quadratic semidefinite subproblem will be an acceptable approximation of  $X^*$ . It is not necessary to construct the whole matrix  $X^*$ . In fact, the factorized form  $(P^k \sqrt{V_*^k})(P^k \sqrt{V_*^k})^T$  is much more convenient to work with. For example the approximation algorithm of Goemans and Williamson [9] requires precisely this factorization. A particular  $x_{ij}$  element of  $X^*$  is easily computed by the inner product of row  $i$  and  $j$  of the  $n \times r$  matrix  $P^k \sqrt{V_*^k}$ . In principle this opens the door for branch and cut approaches to improve the initial relaxation. This will be the subject of

	$\lambda$ -time	serious	iter	$\ g\ $	max- $r$	0.1%-time
$G_1$	3:12	22	33	0.1639	18	48
$G_2$	4:04	21	36	0.05035	18	1:02
$G_3$	3:20	21	33	0.04107	19	57
$G_4$	2:38	19	27	0.08235	19	54
$G_5$	3:41	23	39	0.04425	17	46
$G_6$	2:42	21	35	0.0646	18	1:09
$G_7$	5:46	24	65	0.0854	17	57
$G_8$	2:39	23	38	0.1549	17	59
$G_9$	2:59	24	34	0.0711	17	1:04
$G_{10}$	2:56	23	37	0.02997	17	1:15
$G_{11}$	42:10	97	172	0.4696	15	17:04
$G_{12}$	28:47	50	130	0.2579	15	8:19
$G_{13}$	17:24	43	78	0.218	15	6:17
$G_{14}$	12:43	41	59	0.1682	18	3:09
$G_{15}$	20:16	44	89	0.1059	18	3:08
$G_{16}$	14:35	31	69	0.2246	19	3:11
$G_{17}$	14:38	41	65	0.2079	18	3:22
$G_{18}$	16:46	38	98	0.08161	15	4:15
$G_{19}$	11:24	41	71	0.1571	15	3:34
$G_{20}$	11:05	42	71	0.08226	15	3:50
$G_{21}$	12:23	44	80	0.09432	15	3:32

Table 3: Additional information about the bundle algorithm for the examples of Table 1.  $\lambda$ -time gives the total amount of time spent for computing the eigenvalues and eigenvectors, *serious* gives the number of serious steps, *iter* the total number of iterations including null steps.  $\|g\|$  refers to the norm of the gradient resulting from the optimal solution of the last semidefinite subproblem. *max- $r$*  is the maximum number of columns used in  $P$  (the limit would have been 25). 0.1%-time gives the time when the bound is within  $10^{-3}$  of the optimum in relative precision.

further work.

Table 5 gives a similar set of examples for  $n = m = 2000$ .

A last set of examples is devoted to the Lovász  $\vartheta$ -function [25] which yields an upper bound on the cardinality of a maximal independent (or stable) set of a graph. For implementational convenience we use its formulation within the quadratic  $\{-1, 1\}$  programming setting, see [23]. For a graph with  $k$  nodes and  $h$  edges we obtain a semidefinite program with matrix variables of order  $n = k + 1$  and  $m = k + 1 + h$  constraints. The examples we are going to consider have more than one thousand nodes and more than six thousand edges. For these examples interior point methods are not applicable any more because of memory requirements. It should be clear from the examples of Table 5 that there is also little hope for the bundle method to terminate within reasonable time. However, the most significant progress is achieved in the beginning and for the bundle method memory consumption is not a problem. We run these examples with a time limit of five hours. More precisely, the algorithm is terminated after the first serious step that occurs after five hours of computation time.

The graph instances are of the same type as above. The computational results are displayed in Table 6. The new columns  $n$  and  $m$  give the order of the matrix variable and the number of constraints, respectively. Observe that the toroidal grid graphs  $G_{48}$  and  $G_{49}$  are perfect with independence number 1500; the independence number of  $G_{50}$  is 1440 but  $G_{50}$  is not perfect. We do not know the independence

iter	value	time	$\ g\ $	max- $r$
0	2861.20			
1	2821.75	3	42.02	3
2	2798.60	4	28.83	7
3	2782.99	6	25.94	8
4	2736.11	9	18.58	8
5	2704.98	13	12.69	10
6	2685.84	17	10.05	12
7	2679.95	23	7.46	14
8	2670.10	28	5.992	15
9	2666.07	34	4.173	16
10	2660.31	48	3.268	16
11	2658.65	1:09	0.6295	17
12	2656.66	1:32	0.7753	17
13	2656.42	1:41	0.6636	18
14	2656.31	1:57	0.431	18
15	2656.27	2:14	0.2997	18
16	2656.22	2:45	0.1609	18
17	2656.20	3:01	0.13	18
18	2656.19	3:17	0.09695	18
19	2656.19	3:33	0.07243	18
20	2656.17	3:49	0.05042	18
21	2656.17	3:57	0.0646	18

Table 4: Detailed account of the serious steps of example  $G_6$ .

number of the other graphs. Except for  $G_{48}$  and  $G_{49}$ , which have  $\theta(G_{48}) = \theta(G_{49}) = 1500$  by perfectness, it is hard to judge the quality of the solutions. Tracing the development of the bounds the last serious steps of examples  $G_{43}$  to  $G_{47}$  and  $G_{51}$  to  $G_{54}$  still produced improvements of 0.5% to 1%. This and the rather large norm of the subgradient of  $G_{51}$  and  $G_{54}$  indicate that the values cannot be expected to be ‘good’ approximations of the  $\vartheta$ -function. Also note, that the size of the subspace required for  $G_{48}$  to  $G_{50}$  is still well below 25. In examples  $G_{51}$  to  $G_{54}$  the value of  $\alpha$  is almost negligible, but for  $G_{43}$  to  $G_{47}$  the value of  $\alpha$  is roughly 1/3 at termination. Thus for these examples the restriction to 25 columns became relevant.

The computational results of Table 6 demonstrate that the algorithm has its limits. Nonetheless the bounds obtained are still useful and the primal approximation corresponding to the subgradient is a reasonable starting point for primal heuristics.

## 9 Conclusions and Extensions

We have proposed a proximal bundle method for solving semidefinite programs with large sparse or strongly structured coefficient matrices. The semidefinite constraint is lifted into the objective function by means of a Lagrange multiplier  $a$  whose correct value is not known in general, except for problems with fixed primal trace. In the latter case  $a$  is precisely the value of the trace. The approach differs from previous bundle methods in that the subproblem is tailored for semidefinite programming. In fact the whole approach can be interpreted as semidefinite programming over subspaces where the subspace is successively corrected and improved till the optimal subspace is identified. The set of subgradients modeled by the semidefinite subproblem is a superset of the subgradients used in the traditional polyhedral cut-

	B-sol	B-time	$\lambda$ -time	serious	iter	$\ g\ $	max- $r$
$G_{22}$	14135.98	38:11	28:00	26	52	0.0781	23
$G_{23}$	14145.58	1:19:29	1:01:06	32	107	0.3707	23
$G_{24}$	14140.88	28:04	21:11	25	40	0.1565	23
$G_{25}$	14144.30	43:44	32:26	27	59	0.3452	24
$G_{26}$	14132.93	34:45	26:37	31	48	0.3066	23
$G_{27}$	4141.68	24:56	18:54	23	37	0.1423	23
$G_{28}$	4100.81	29:41	21:08	23	39	0.1954	25
$G_{29}$	4208.94	48:16	37:53	27	65	0.1725	22
$G_{30}$	4215.42	1:02:39	48:16	27	83	0.1282	22
$G_{31}$	4116.70	26:11	19:02	24	38	0.1889	24
$G_{32}$	1567.80	6:20:54	6:09:52	144	312	0.892	15
$G_{33}$	1544.42	6:04:22	5:53:37	112	305	0.6887	15
$G_{34}$	1546.82	3:46:31	3:39:26	105	198	0.8369	15
$G_{35}$	8014.81	4:31:17	3:54:01	61	209	0.2397	22
$G_{36}$	8006.04	2:56:10	2:31:09	62	115	0.2634	24
$G_{37}$	8018.68	3:10:01	2:46:35	58	130	0.254	23
$G_{38}$	8015.01	4:03:53	3:39:24	58	155	0.1937	22
$G_{39}$	2877.71	1:20:24	1:12:23	50	85	0.2194	20
$G_{40}$	2864.96	2:42:02	2:30:21	59	158	0.2737	19
$G_{41}$	2865.29	1:41:33	1:32:50	51	108	0.1954	19
$G_{42}$	2946.29	1:32:45	1:24:12	59	93	0.1535	20

Table 5: Examples for  $n = m = 2000$ .

ting plane model. Therefore convergence of the new method is a direct consequence of previous proofs for traditional bundle methods. It is not yet clear whether the specialized model admits stronger convergence results. The choice of  $u$  is still very much an open problem of high practical importance.

For (constrained) quadratic  $\{-1, 1\}$ -programming the method offers a good bound within reasonable time and allows to construct an approximate primal optimal solution (of the relaxation) in compact representation. To improve the bound by a cutting plane approach the algorithm must be able to deal with sign constraints on the  $y$ -variables. In principle it is not difficult to model the sign constraints in the semidefinite subproblem. However, as a consequence the order of the system matrix grows linearly with the number of constraints on  $y$  rendering the method impractical even for a moderate number of cutting planes. Alternatively one might consider active set methods but these entail the danger of destroying convergence. Together with K.C. Kiwiel we are currently working on alternative methods for incorporating sign constraints on  $y$ .

The backbone of the method is an efficient routine for computing the maximal eigenvalue of huge structured symmetric matrices. Although our own implementation (based on the code of Hua) seems to work sufficiently stable there is certainly much room for improvement. A straight forward approach to achieve serious speed-ups is to implement the algorithm on parallel machines, see for instance [38]. Rather recently interest in the Lanczos method has risen again, see [24, 4, 6, 11, 28] and references therein. Most of these papers are based on the concept of an implicit restart proposed in [39] which is a polynomial acceleration approach that does not require additional matrix vector multiplications. It will be interesting to test these new ideas within the bundle framework.

We thank K.C. Kiwiel for fruitful discussions.

	n	m	B-sol	B-time	$\lambda$ -time	serious	iter	$\ g\ $	max-r
$G_{43}$	1001	10991	308.47	6:02:20	3:40:57	38	104	0.5098	25
$G_{44}$	1001	10991	310.13	5:06:31	3:14:25	39	88	0.5493	25
$G_{45}$	1001	10991	309.00	5:10:38	3:19:25	40	87	0.5532	25
$G_{46}$	1001	10991	309.21	5:35:56	3:31:42	42	99	0.5535	25
$G_{47}$	1001	10991	310.84	5:24:52	3:10:35	44	100	0.5558	25
$G_{48}$	3001	9001	1526.53	5:11:31	4:59:57	54	94	0.4062	15
$G_{49}$	3001	9001	1521.24	5:06:21	4:53:45	53	102	0.3751	15
$G_{50}$	3001	9001	1536.12	5:17:51	5:01:25	50	124	0.4728	15
$G_{51}$	1001	6910	455.21	5:07:40	4:29:00	32	118	2.556	25
$G_{52}$	1001	6917	465.12	5:09:02	4:38:30	41	108	2.683	25
$G_{53}$	1001	6915	463.86	5:08:36	4:36:34	41	104	2.593	25
$G_{54}$	1001	6917	466.04	5:02:21	4:32:21	40	98	2.590	25

Table 6: Upper bound on the  $\vartheta$ -function after five hours of computation time.

## References

- [1] F. Alizadeh. Interior point methods in semidefinite programming with applications to combinatorial optimization. *SIAM J. Optim.*, 5(1):13–51, 1995.
- [2] F. Alizadeh, J.-P. A. Haeberly, and M. L. Overton. Primal-dual interior-point methods for semidefinite programming: Convergence rates, stability and numerical results. Technical Report 721, Computer Science Department, Courant Institute of Mathematical Sciences, New York University, May 1996.
- [3] F. Alizadeh, J.-P. A. Haeberly, and M. L. Overton. Complementarity and nondegeneracy in semidefinite programming. *Math. Programming*, 77(2):111–128, 1997.
- [4] J. Baglama, D. Calvetti, and L. Reichel. Iterative methods for the computation of a few eigenvalues of a large symmetric matrix. *BIT*, 36:400–421, 1996.
- [5] S. Benson, Y. Ye, and X. Zhang. Solving large-scale sparse semidefinite programs for combinatorial optimization. Working paper, Department of Management Science, University of Iowa, IA, 52242, USA, Sept. 1997.
- [6] D. Calvetti, L. Reichel, and D. Sorensen. An implicitly restarted Lanczos method for large symmetric eigenvalue problems. *Electronic Transactions on Numerical Analysis*, 2:1–21, Mar. 1994.
- [7] J. Cullum, W. E. Donath, and P. Wolfe. The minimization of certain non-differentiable sums of eigenvalues of symmetric matrices. *Math. Programming Study*, 3:35–55, 1975.
- [8] L. Faybusovich. Semidefinite programming: a path-following algorithm for a linear-quadratic functional. *SIAM J. Optim.*, pages 1007–1024, 1996.
- [9] M. X. Goemans and D. P. Williamson. Improved approximation algorithms for maximum cut and satisfiability problems using semidefinite programming. *J. ACM*, 42:1115–1145, 1995.
- [10] G. H. Golub and C. F. van Loan. *Matrix Computations*. The Johns Hopkins University Press, 2<sup>nd</sup> edition, 1989.

$G_1$	-rnd_graph 800 6 8001
$G_2$	-rnd_graph 800 6 8002
$G_3$	-rnd_graph 800 6 8003
$G_4$	-rnd_graph 800 6 8004
$G_5$	-rnd_graph 800 6 8005
$G_6$	-rnd_graph 800 6 8001 -random 0 1 8001 -times 2 -plus -1
$G_7$	-rnd_graph 800 6 8002 -random 0 1 8002 -times 2 -plus -1
$G_8$	-rnd_graph 800 6 8003 -random 0 1 8003 -times 2 -plus -1
$G_9$	-rnd_graph 800 6 8004 -random 0 1 8004 -times 2 -plus -1
$G_{10}$	-rnd_graph 800 6 8005 -random 0 1 8005 -times 2 -plus -1
$G_{11}$	-toroidal_grid_2D 100 8 -random 0 1 8004 -times 2 -plus -1
$G_{12}$	-toroidal_grid_2D 50 16 -random 0 1 8002 -times 2 -plus -1
$G_{13}$	-toroidal_grid_2D 25 32 -random 0 1 8003 -times 2 -plus -1
$G_{14}$	-planar 800 99 8001 -planar 800 99 8002 +
$G_{15}$	-planar 800 99 8003 -planar 800 99 8004 +
$G_{16}$	-planar 800 99 8005 -planar 800 99 8006 +
$G_{17}$	-planar 800 99 8007 -planar 800 99 8008 +
$G_{18}$	-planar 800 99 8001 -planar 800 99 8002 + -random 0 1 8001 -times 2 -plus -1
$G_{19}$	-planar 800 99 8003 -planar 800 99 8004 + -random 0 1 8002 -times 2 -plus -1
$G_{20}$	-planar 800 99 8005 -planar 800 99 8006 + -random 0 1 8003 -times 2 -plus -1
$G_{21}$	-planar 800 99 8007 -planar 800 99 8008 + -random 0 1 8004 -times 2 -plus -1
$G_{22}$	-rnd_graph 2000 1 20001
$G_{23}$	-rnd_graph 2000 1 20002
$G_{24}$	-rnd_graph 2000 1 20003
$G_{25}$	-rnd_graph 2000 1 20004
$G_{26}$	-rnd_graph 2000 1 20005
$G_{27}$	-rnd_graph 2000 1 20001 -random 0 1 20001 -times 2 -plus -1
$G_{28}$	-rnd_graph 2000 1 20002 -random 0 1 20002 -times 2 -plus -1
$G_{29}$	-rnd_graph 2000 1 20003 -random 0 1 20003 -times 2 -plus -1
$G_{30}$	-rnd_graph 2000 1 20004 -random 0 1 20004 -times 2 -plus -1
$G_{31}$	-rnd_graph 2000 1 20005 -random 0 1 20005 -times 2 -plus -1
$G_{32}$	-toroidal_grid_2D 100 20 -random 0 1 20003 -times 2 -plus -1
$G_{33}$	-toroidal_grid_2D 80 25 -random 0 1 20002 -times 2 -plus -1
$G_{34}$	-toroidal_grid_2D 50 40 -random 0 1 20001 -times 2 -plus -1
$G_{35}$	-planar 2000 99 20001 -planar 2000 99 20002 +
$G_{36}$	-planar 2000 99 20003 -planar 2000 99 20004 +
$G_{37}$	-planar 2000 99 20005 -planar 2000 99 20006 +
$G_{38}$	-planar 2000 99 20007 -planar 2000 99 20008 +
$G_{39}$	-planar 2000 99 20001 -planar 2000 99 20002 + -random 0 1 20001 -times 2 -plus -1
$G_{40}$	-planar 2000 99 20003 -planar 2000 99 20004 + -random 0 1 20002 -times 2 -plus -1
$G_{41}$	-planar 2000 99 20005 -planar 2000 99 20006 + -random 0 1 20003 -times 2 -plus -1
$G_{42}$	-planar 2000 99 20007 -planar 2000 99 20008 + -random 0 1 20004 -times 2 -plus -1
$G_{43}$	-rnd_graph 1000 2 10001
$G_{44}$	-rnd_graph 1000 2 10002
$G_{45}$	-rnd_graph 1000 2 10003
$G_{46}$	-rnd_graph 1000 2 10004
$G_{47}$	-rnd_graph 1000 2 10005
$G_{48}$	-toroidal_grid_2D 50 60
$G_{49}$	-toroidal_grid_2D 30 100
$G_{50}$	-toroidal_grid_2D 25 120
$G_{51}$	-planar 1000 100 10001 -planar 1000 100 10002 +
$G_{52}$	-planar 1000 100 10003 -planar 1000 100 10004 +
$G_{53}$	-planar 1000 100 10005 -planar 1000 100 10006 +
$G_{54}$	-planar 1000 100 10007 -planar 1000 100 10008 +

Table 7: Arguments for generating the graphs by the graph generator *rudy*.

- [11] R. G. Grimes, J. G. Lewis, and H. D. Simon. A shifted block Lanczos algorithm for solving sparse symmetric generalized eigenproblems. *SIAM J. Matrix Anal. Appl.*, 15(1):228–272, Jan. 1994.
- [12] M. Grötschel, L. Lovász, and A. Schrijver. *Geometric Algorithms and Combinatorial Optimization*, volume 2 of *Algorithms and Combinatorics*. Springer, 2<sup>nd</sup> edition, 1988.
- [13] C. Helmberg. Fixing variables in semidefinite relaxations. In R. Burkard and G. Woeginger, editors, *Algorithms – ESA’97*, volume 1284 of *Lecture Notes in Computer Science*, pages 259–270. Springer, Sept. 1997.
- [14] C. Helmberg and F. Rendl. Solving quadratic (0,1)-problems by semidefinite programs and cutting planes. ZIB Preprint SC-95-35, Konrad-Zuse-Zentrum für Informationstechnik Berlin, Nov. 1995. To appear in Math. Programming.
- [15] C. Helmberg, F. Rendl, R. J. Vanderbei, and H. Wolkowicz. An interior-point method for semidefinite programming. *SIAM J. Optim.*, 6(2):342–361, May 1996.
- [16] C. Helmberg, F. Rendl, and R. Weismantel. Quadratic knapsack relaxations using cutting planes and semidefinite programming. In W. H. Cunningham, S. T. McCormick, and M. Queyranne, editors, *Integer Programming and Combinatorial Optimization*, volume 1084 of *Lecture Notes in Computer Science*, pages 175–189. Springer, June 1996.

- [17] J.-B. Hiriart-Urruty and C. Lemaréchal. *Convex Analysis and Minimization Algorithms I*, volume 305 of *Grundlehren der mathematischen Wissenschaften*. Springer, Berlin, Heidelberg, 1993.
- [18] J.-B. Hiriart-Urruty and C. Lemaréchal. *Convex Analysis and Minimization Algorithms II*, volume 306 of *Grundlehren der mathematischen Wissenschaften*. Springer, Berlin, Heidelberg, 1993.
- [19] F. Jarre. An interior-point method for minimizing the maximum eigenvalue of a linear combination of matrices. *Siam J. Control and Optimization*, 31(5):1360–1377, Sept. 1993.
- [20] K. C. Kiwiel. Proximity control in bundle methods for convex nondifferentiable minimization. *Math. Programming*, 46:105–122, 1990.
- [21] P. Klein and H.-I. Lu. Efficient approximation algorithms for semidefinite programs arising from MAXCUT and COLORING. In *STOC’96*, pages 338–347, 1996.
- [22] M. Kojima, S. Shindoh, and S. Hara. Interior-point methods for the monotone linear complementarity problem in symmetric matrices. *SIAM J. Optim.*, 7(1):86–125, Feb. 1997.
- [23] M. Laurent, S. Poljak, and F. Rendl. Connections between semidefinite relaxations of the max-cut and stable set problems. *Math. Programming*, 77(2):225–246, 1997.
- [24] R. B. Lehoucq and K. J. Maschhoff. Implementation of an implicitly restarted block Arnoldi method. Technical Report MCS-P649-0297, Argonne National Laboratory, 1997.
- [25] L. Lovász. On the Shannon capacity of a graph. *IEEE Transactions on Information Theory*, IT-25(1):1–7, Jan. 1979.
- [26] L. Lovász and A. Schrijver. Cones of matrices and set-functions and 0-1 optimization. *SIAM J. Optim.*, 1(2):166–190, May 1991.
- [27] R. D. C. Monteiro. Polynomial convergence of primal-dual algorithms for semidefinite programming based on Monteiro and Zhang family of directions. Working paper, School of Industrial and Systems Engineering, Georgia Tech, Atlanta, GA 30332, July 1996.
- [28] R. B. Morgan and D. S. Scott. Preconditioning the Lanczos algorithm for sparse symmetric eigenvalue problems. *SIAM J. Sci. Comput.*, 14(3):585–593, May 1993.
- [29] Y. Nesterov and A. Nemirovskii. *Interior-Point Polynomial Algorithms in Convex Programming*. SIAM Studies in Applied Mathematics, Philadelphia, 1994.
- [30] Y. Nesterov and M. J. Todd. Self-scaled barriers and interior-point methods for convex programming. *Math. Oper. Res.*, 22:1–42, 1997.
- [31] F. Oustry. The U-Lagrangian of the maximum eigenvalue function. Technical report, INRIA-ENSTA, France, Nov. 1996.
- [32] M. L. Overton. Large-scale optimization of eigenvalues. *SIAM J. Optim.*, 2(1):88–120, Feb. 1992.

- [33] G. Pataki. On the rank of extreme matrices in semidefinite programming and the multiplicity of optimal eigenvalues. Technical Report #MSRR-604, Carnegie Mellon University, Apr. 1994. Revised Aug. 1995.
- [34] S. A. Plotkin, D. B. Shmoys, and E. Tardos. Fast approximation algorithms for fractional packing and covering problems. *Math. Oper. Res.*, 20(2):257–301, May 1995.
- [35] S. Poljak and F. Rendl. Node and edge relaxations of the Max-Cut problem. *Computing*, 52:123–137, 1994.
- [36] S. Poljak and F. Rendl. Nonpolyhedral relaxations of graph-bisection problems. *SIAM J. Optim.*, 5(3):467–487, 1995.
- [37] H. Schramm and J. Zowe. A version of the bundle idea for minimizing a nonsmooth function: conceptual idea, convergence analysis, numerical results. *SIAM J. Optim.*, 2:121–152, 1992.
- [38] D. S. Scott. Implementing Lanczos-like algorithms on hypercube architectures. *Computer Physics Communications*, 53:271–281, 1989.
- [39] D. C. Sorensen. Implicit application of polynomial filters in a k-step Arnoldi method. *SIAM J. Matrix Anal. Appl.*, 13(1):357–385, 1992.
- [40] M. J. Todd, K. C. Toh, and R. H. Tütüncü. On the Nesterov-Todd direction in semidefinite programming. Technical Report TR 1154, School of Operations Research and Industrial Engineering, Cornell University, Ithaca, New York 14853, Mar. 1996. Revised May 1996.
- [41] L. Vandenberghe and S. Boyd. A primal-dual potential reduction method for problems involving matrix inequalities. *Mathematical Programming, Series B*, 69(1):205–236, 1995.
- [42] D. Y. Ye. Sensitivity analysis of the greatest eigenvalue of a symmetric matrix via the  $\varepsilon$ -subdifferential of the associated convex quadratic form. *J. Optimization Theory and Applications*, 76(2):287–304, Feb. 1993.
- [43] S. Zhou and H. Dai. The block Chebyshev-Lanczos method for solving large symmetric eigenvalue problems. *J. Nanjing Aeronaut. Inst.*, 21(4):22–28, 1989.

## A Notation

$\mathbb{R}^n$	real column vector of dimension $n$
$M_{m,n}$	$m \times n$ real matrices
$S_n$	$n \times n$ symmetric real matrices
$S_n^{++}$	$n \times n$ symmetric positive definite matrices
$S_n^+$	$n \times n$ symmetric positive semidefinite matrices
$A \succ 0$	$A$ is positive definite
$A \succeq 0$	$A$ is positive semidefinite
$I, I_n$	identity of appropriate size or of size $n$
$e$	vector of all ones of appropriate dimension
$e_i$	i-th column of $I$
$\lambda_i(A)$	i-th eigenvalue of $A \in M_n$ , usually $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$
$\lambda_{\min}(A), \lambda_{\max}(A)$	minimal and maximal eigenvalue of $A$
$\Lambda_A$	diagonal matrix with $(\Lambda_A)_{ii} = \lambda_i(A)$
$\text{tr}(A)$	trace of $A \in M_n$ , $\text{tr}(A) = \sum_{i=1}^n a_{ii} = \sum_{i=1}^n \lambda_i(A)$
$\langle A, B \rangle$	inner product in $M_{m,n}$ , $\langle A, B \rangle = \text{tr}(B^T A)$
$\text{svec}(A)$	$\binom{n+1}{2}$ dimensional vector representation of $A \in S_n$
$A \otimes_s B$	symmetric Kronecker product of $A \in M_{m,n}, B \in M_{k,l}$
$\text{diag}(A)$	the diagonal of $A \in M_n$ as a column vector
$\text{Diag}(v)$	diagonal matrix with $v$ on its main diagonal

$S_n$  is isomorphic to  $\mathbb{R}^{\binom{n+1}{2}}$  via the map  $\text{svec}(A)$  defined by stacking the columns of the lower triangle of  $A$  on top of each other and multiplying the offdiagonal elements with  $\sqrt{2}$ ,

$$\text{svec}(A) := [a_{11}, \sqrt{2}a_{21}, \dots, \sqrt{2}a_{n1}, a_{22}, \sqrt{2}a_{32}, \dots, a_{nn}]^T.$$

The factor  $\sqrt{2}$  for offdiagonal elements ensures that, for  $A, B \in S_n$ ,

$$\langle A, B \rangle = \text{tr}(AB) = \text{svec}(A)^T \text{svec}(B).$$

The symmetric Kronecker product  $\otimes_s$  is defined for arbitrary square matrices  $A, B \in M_{n,n}$  by its action on a vector  $\text{svec}(C)$  for a symmetric matrix  $C \in S_n$ ,

$$(A \otimes_s B)\text{svec}(C) := \frac{1}{2}\text{svec}(BCA^T + ACB^T).$$

Both concepts were first introduced in [2]. Here we use the notation introduced in [40]. From the latter paper we also cite some properties of the symmetric Kronecker product for the convenience of the reader.

1.  $A \otimes_s B = B \otimes_s A$
2.  $(A \otimes_s B)^T = B^T \otimes_s A^T$
3.  $A \otimes_s I$  is symmetric if and only if  $A$  is.
4.  $(A \otimes_s A)^{-1} = A^{-1} \otimes_s A^{-1}$
5.  $(A \otimes_s B)(C \otimes_s D) = \frac{1}{2}(AC \otimes_s BD + AD \otimes_s BC)$
6. If  $A \succ 0$  and  $B \succ 0$  then  $(A \otimes_s B) \succ 0$
7.  $\text{svec}(A)^T \text{svec}(B) = \langle A, B \rangle = \text{tr}(AB)$