



# Convergence rate of block-coordinate maximization Burer–Monteiro method for solving large SDPs

Murat A. Erdogdu<sup>1,2</sup> · Asuman Ozdaglar<sup>3</sup> · Pablo A. Parrilo<sup>3</sup> ·  
Nuri Denizcan Vanli<sup>3</sup>

Received: 10 December 2019 / Accepted: 26 June 2021

© Springer-Verlag GmbH Germany, part of Springer Nature and Mathematical Optimization Society 2021

## Abstract

Semidefinite programming (SDP) with diagonal constraints arise in many optimization problems, such as Max-Cut, community detection and group synchronization. Although SDPs can be solved to arbitrary precision in polynomial time, generic convex solvers do not scale well with the dimension of the problem. In order to address this issue, Burer and Monteiro (Math Program 95(2):329–357, 2003) proposed to reduce the dimension of the problem by appealing to a low-rank factorization and solve the subsequent non-convex problem instead. In this paper, we present coordinate ascent based methods to solve this non-convex problem with provable convergence guarantees. More specifically, we prove that the block-coordinate maximization algorithm applied to the non-convex Burer–Monteiro method globally converges to a first-order stationary point with a sublinear rate without any assumptions on the problem. We further show that this algorithm converges linearly around a local maximum provided that the objective function exhibits quadratic decay. We establish that this condition generically holds when the rank of the factorization is sufficiently large. Furthermore,

---

Part of this work has previously appeared in ICML 2018 Workshop on Modern Trends in Nonconvex Optimization for Machine Learning.

---

✉ Nuri Denizcan Vanli  
denizcan@mit.edu

Murat A. Erdogdu  
erdogdu@cs.toronto.edu

Asuman Ozdaglar  
asuman@mit.edu

Pablo A. Parrilo  
parrilo@mit.edu

<sup>1</sup> Department of Computer Science, University of Toronto, Toronto, Canada

<sup>2</sup> Department of Statistical Sciences, University of Toronto, Toronto, Canada

<sup>3</sup> Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge, USA

incorporating Lanczos method to the block-coordinate maximization, we propose an algorithm that is guaranteed to return a solution that provides  $1 - \mathcal{O}(1/r)$  approximation to the original SDP without any assumptions, where  $r$  is the rank of the factorization. This approximation ratio is known to be optimal (up to constants) under the unique games conjecture, and we can explicitly quantify the number of iterations to obtain such a solution.

**Keywords** Semidefinite programming · Burer–Monteiro method · Coordinate descent · Non-convex optimization · Large-scale optimization

**Mathematics Subject Classification** 65K05, 90C22, 90C25, 90C26, 90C27, 90C30, 58C05, 49M37

## 1 Introduction

A variety of problems in statistical estimation and machine learning require solving a combinatorial optimization problem, which are often intractable [42]. Semidefinite programs (SDP) are commonly used as convex relaxations for these problems, providing efficient algorithms with approximate optimality [34]. A generic SDP in this framework can be written as

$$\begin{aligned} & \text{maximize } \langle \mathbf{A}, \mathbf{X} \rangle && \text{(CVX)} \\ & \text{subject to } X_{ii} = 1, \text{ for } i \in [n], \\ & \mathbf{X} \succeq 0, \end{aligned}$$

where  $\mathbf{A}, \mathbf{X} \in \text{Sym}_n$  (real symmetric matrices of size  $n \times n$ ) and  $[n] = \{1, 2, \dots, n\}$ . This problem appears as a convex relaxation to the celebrated Max-Cut problem [22], graphical model inference [19], community detection problems [6], and group synchronization [32].

Although SDPs serve as reliable relaxations to many combinatorial problems, the resulting convex problem is still computationally challenging. Interior point methods can solve SDPs to arbitrary accuracy in polynomial-time, but they do not scale well with the problem dimension  $n$ . A popular approach to remedy these limitations is to introduce a low-rank factorization  $\mathbf{X} = \boldsymbol{\sigma} \boldsymbol{\sigma}^\top$ , where  $\boldsymbol{\sigma} \in \mathbb{R}^{n \times r}$  with  $r$  denoting the rank. This reformulation removes the positive semidefinite cone constraint in (CVX) since  $\mathbf{X} = \boldsymbol{\sigma} \boldsymbol{\sigma}^\top$  is guaranteed to be a positive semidefinite matrix, and choosing  $r \ll n$  provides computational efficiency as well as storage benefits. This method is often referred to as Burer–Monteiro approach [15]. Denoting  $i$ th row of  $\boldsymbol{\sigma}$  by  $\sigma_i$ , i.e.,  $\boldsymbol{\sigma} = [\sigma_1, \sigma_2, \dots, \sigma_n]^\top$ , the resulting non-convex problem can be written as follows

$$\begin{aligned} & \text{maximize } \langle \mathbf{A}, \boldsymbol{\sigma} \boldsymbol{\sigma}^\top \rangle && \text{(Non-CVX)} \\ & \text{subject to } \|\sigma_i\| = 1, \text{ for } i \in [n]. \end{aligned}$$

In the original Burer–Monteiro approach [15], the authors proposed to use an augmented Lagrangian method for a general form SDP. However, it has been recently

observed that feasible methods (such as block-coordinate maximization [25,43], Riemannian gradient [25,32] and Riemannian trust-region methods [1,11,26]) provide empirically faster rates since feasibility can be efficiently guaranteed via projection onto the Cartesian product of spheres. Despite overwhelming empirical evidence [25,32,43], convergence properties of these feasible methods are not well-understood (except the Riemannian trust-region method, for which a sublinear convergence rate is shown in [11] and a local superlinear convergence is shown in [1] with no rate estimate). Among these methods, block-coordinate maximization and Riemannian gradient ascent are simpler to implement and have computational complexity of  $\mathcal{O}(nr)$  and  $\mathcal{O}(n^2r)$ , respectively, whereas Riemannian trust-region requires to solve the trust-region subproblem at each iteration, which is usually solved iteratively using the Lanczos method in a few iterations, whose per iteration requires  $\mathcal{O}(n^2r)$  arithmetic operations. Furthermore, block-coordinate maximization does not have any step size or tuning parameters, unlike Riemannian gradient ascent and Riemannian trust-region methods. Empirical studies further motivate the use of block-coordinate maximization by presenting superior performance compared to existing methods on large-scale problems, often with linear convergence [25]. In this paper, we provide the first local and global convergence rate guarantees for the block-coordinate maximization method (applied to Burer–Monteiro approach) in the literature, which are consistent with the empirical performance of the algorithm. Our contributions can be summarized as follows:

- We establish the global sublinear convergence of the block-coordinate maximization algorithm applied to (Non-CVX) without any assumptions on the cost matrix  $A$ .
- We show that this algorithm enjoys a linear rate around a neighborhood of any local maximum when the objective function satisfies the quadratic decay assumption.
- We establish that the quadratic decay condition that leads to local linear convergence generically holds when the rank of the factorization satisfies  $r \geq \sqrt{2n}$ .
- Incorporating Lanczos methods into the block-coordinate maximization procedure, we propose an algorithm that returns an approximate second-order stationary point of (Non-CVX). By choosing the rank of the factorization sufficiently large and selecting the parameters of the algorithm according to the cost matrix  $A$ , we show that the solution returned by this algorithm is not only an approximate local maximum to (Non-CVX), but also provides  $1 - \mathcal{O}(1/r)$  approximation to (CVX). We highlight that this approximation ratio is optimal under the unique games conjecture.
- We validate our theoretical results via numerical examples and compare the performance of the block-coordinate maximization algorithm with various manifold optimization methods to demonstrate its performance.

## 1.1 Related work

There are numerous papers that analyze the landscape of the solution space of (Non-CVX). In particular, it is known that (CVX) admits an optimal solution of rank  $r$  such that  $r(r+1)/2 \leq n$  [7,35]. Using this observation, it has been shown in

[15,16,26] that when  $r \geq \sqrt{2n}$ , if  $\sigma$  is a rank deficient second-order stationary point of (Non-CVX), then  $\sigma$  is a global maximum for (Non-CVX) and  $X = \sigma\sigma^\top$  is a global maximum for (CVX). The recent paper [12] showed that when  $r \geq \sqrt{2n}$ , for almost all  $A$ , every  $\sigma$  that is a first-order stationary point is rank deficient. For arbitrary rank  $r$ , it is shown that all local maxima are within a  $n\|A\|_2/\sqrt{r}$  gap from the optimum of (CVX) [33], and any  $\varepsilon$ -approximate concave point is within a  $\text{Rg}(\text{Non-CVX})/(r-1) + n\varepsilon/2$  gap from the optimum of (CVX) [32], where  $\text{Rg}(\text{Non-CVX})$  is the range of the problem (Non-CVX), i.e., the difference between the maximum and the minimum values of the objective in (Non-CVX).

Javanmard *et al.* [25] showed that when applied to solve (Non-CVX), Riemannian gradient ascent and block-coordinate maximization methods provide excellent numerical results, yet no convergence guarantee is provided. Similar experimental results are also observed in [43] for the block-coordinate maximization algorithm and in [32] for the Riemannian gradient ascent algorithm. Concurrent to this work, in [43], the authors analyzed the convergence of the deterministic block-coordinate maximization algorithm. In particular, they showed that the deterministic block-coordinate maximization algorithm is asymptotically convergent (see [43, Theorem 3.2]) and enjoys a local linear convergence with no explicit rate estimates (see [43, Theorem 3.5]). They also proved that the deterministic block-coordinate maximization approach converges to a local maximum generically under random initialization using the center-stable manifold theorem similar to [30]. These results hold under the assumption that the iterates generated by the algorithm satisfy a certain condition that is seemingly impossible to verify without actually running the algorithm. To alleviate this issue, the authors suggested using a coordinate ascent method with a sufficiently small step size, for which the aforementioned convergence results hold without this precarious assumption. In [9], the authors provided a global sublinear convergence rate for the Riemannian trust-region method for general non-convex problems and these results have been used in [11,32] for the non-convex Burer–Monteiro approach. Augmented Lagrangian methods have been proposed to solve (Non-CVX) as well [15,16], however these methods do not benefit from separability of the manifold constraints, and hence are usually slower [12].

There also exist methods that solve (CVX) by exploiting its special structure [5,21,27,39]. In particular, [27] reduces (CVX) to a sequence of approximate eigenpair computations that is efficiently solved using the power method. In [5,39], matrix multiplicative weights algorithm is used to approximately solve (CVX), and these ideas are extended in [21] using sketching techniques [40]. However, these methods require constructing  $X \in \text{Sym}_n$  explicitly, which is prohibitive when  $n$  goes beyond a few thousands, whereas the Burer–Monteiro approach we consider in this paper easily scales to very large instances as the low-rank factorization decreases the dimension of the problem from  $\mathcal{O}(n^2)$  to  $\mathcal{O}(nr)$  with  $r \ll n$ . For time complexity comparison between these methods that are based on Lagrangian relaxation and the Burer–Monteiro approach in this paper, we refer to Corollary 4.

Coordinate descent methods have been successfully applied to non-convex differentiable optimization problems in several papers [31,36,38,41]. In [41], the authors propose a coordinate gradient descent approach that may be viewed as a hybrid of gradient-projection and coordinate descent to minimize the sum of a smooth function

and a convex separable function. They analyze the greedy coordinate selection rule and present local linear convergence, although no rate estimates are provided. [38] considers a similar composite but convex optimization problem and provides explicit rate estimates. These results are then generalized to non-convex problems by [36] and [31]. However, these approaches heavily rely on the Euclidean geometry and cannot handle non-convex constraints, which is the main focus of our paper.

## 1.2 Notations and preliminaries

Throughout the paper, matrices are denoted with a boldface font, and all vectors are column vectors. The superscripts are used to denote iteration counters, i.e.,  $\sigma^k$  denotes the value of  $\sigma$  at iteration  $k$ . For a vector  $g$ ,  $\|g\|$  represents its Euclidean norm. For matrices  $\mathbf{A}$ ,  $\mathbf{B}$ , we write  $\langle \mathbf{A}, \mathbf{B} \rangle = \text{trace}(\mathbf{A}\mathbf{B}^\top)$  for the inner product associated to the Frobenius norm  $\|\mathbf{A}\|_F = \sqrt{\langle \mathbf{A}, \mathbf{A} \rangle}$ .  $A_{ij}$  represents the entry at the  $i$ th row and  $j$ th column of  $\mathbf{A}$ ,  $A_i$  represents its  $i$ th row as a column vector, and  $\|\mathbf{A}\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^n |A_{ij}|$  represents its 1-norm, and  $\|\mathbf{A}\|_{1,1} = \sum_{i,j=1}^n |A_{ij}|$  represents its  $L_{1,1}$ -norm. For a function  $h$ ,  $\nabla h$  and  $\text{grad}h$  represent its Euclidean and Riemannian gradients, respectively. Similarly,  $\nabla^2 h$  and  $\text{Hess}h$  represent its Euclidean and Riemannian Hessians, respectively. We let  $S^{m-1}$  denote the unit sphere in  $\mathbb{R}^m$ . For a vector  $y$ ,  $\text{Diag}(y)$  represents the diagonal matrix whose  $i$ th diagonal entry is  $y_i$ . Similarly for a matrix  $\mathbf{A}$ ,  $\text{diag}(\mathbf{A})$  represents the vector whose  $i$ th entry is  $A_{ii}$ .

For brevity, we assume without loss of generality that  $\mathbf{A}$  is a symmetric matrix and  $A_{ii} = 0$ , for all  $i \in [n]$  (the latter assumption is removed in Sect. 4 to keep our presentation consistent with the existing works in the literature). Indeed, if  $\mathbf{A}$  is not symmetric, then we can replace  $\mathbf{A}$  by  $(\mathbf{A} + \mathbf{A}^\top)/2$ , which is a symmetric matrix, and the objective value (**Non-CVX**) remains the same for all  $\sigma \in \mathbb{R}^{n \times r}$  since  $\sigma\sigma^\top$  is symmetric. Similarly, replacing the diagonal entries of  $\mathbf{A}$  by zeros decreases the objective value by the constant  $\text{Tr}(\mathbf{A})$  for all feasible  $\sigma$ , since the diagonal entries of  $\sigma\sigma^\top$  are equal to 1.

The rest of the paper is organized as follows. In Sect. 2, we present the algorithm and discuss its per iteration cost. In Sect. 3, we prove the global sublinear convergence and local linear convergence of the algorithm with explicit rate estimates. In Sect. 4, we introduce a second-order method based on block-coordinate maximization and Lanczos method that is guaranteed to return solutions with global optimality guarantees. We also provide a global sublinear convergence rate estimate for this algorithm. We perform numerical experiments to validate our theoretical results in Sect. 5 and conclude the paper in Sect. 6.

## 2 Block-coordinate maximization algorithm

In this section, we discuss the block-coordinate maximization (BCM) algorithm, its update rule and per iteration computational cost. Throughout the paper, we let  $f : \mathbb{R}^{n \times r} \rightarrow \mathbb{R}$  denote the objective function of (**Non-CVX**), i.e.,

**Algorithm 1** Block-Coordinate Maximization (BCM)

Initialize  $\sigma^0 \in \mathbb{R}^{n \times r}$  and calculate  $g_i^0 = \sum_{j \neq i} A_{ij} \sigma_j^0$ , for all  $i \in [n]$ .  
**for**  $k = 0, 1, 2, \dots$  **do**  
    Choose block  $i_k = i$  using one of the coordinate selection rules.  
     $\sigma_{i_k}^{k+1} \leftarrow g_{i_k}^k / \|g_{i_k}^k\|$ .  
     $g_i^{k+1} \leftarrow g_i^k - A_{i i_k} \sigma_{i_k}^k + A_{i i_k} \sigma_{i_k}^{k+1}$ , for all  $i \neq i_k$ .  
**end for**

$$f(\sigma) = \langle A, \sigma \sigma^\top \rangle.$$

Given the current iterate  $\sigma^k$ , the BCM algorithm chooses a row  $i_k \in [n]$  of the matrix  $\sigma^k$  and maximizes the following objective

$$f(\sigma^k) = \sum_{i=1}^n \langle \sigma_i^k, g_i^k \rangle, \quad \text{where } g_i^k := \sum_{j \neq i} A_{ij} \sigma_j^k,$$

over the block  $\sigma_{i_k}^k \in \mathbb{S}^{r-1}$ . More formally, we can write the update rule of the algorithm as follows

$$\begin{aligned} \sigma_{i_k}^{k+1} &= \arg \max_{\|\zeta\|=1} f(\sigma_1^k, \dots, \sigma_{i_k-1}^k, \zeta, \sigma_{i_k+1}^k, \dots, \sigma_n^k), \\ &= \arg \max_{\|\zeta\|=1} 2\langle \zeta, g_{i_k}^k \rangle + \sum_{i \neq i_k} \sum_{j \neq i, i_k} A_{ij} \langle \sigma_i^k, \sigma_j^k \rangle, \tag{1} \\ &= \arg \max_{\|\zeta\|=1} \langle \zeta, g_{i_k}^k \rangle = \frac{g_{i_k}^k}{\|g_{i_k}^k\|}, \tag{2} \end{aligned}$$

with the convention that  $\sigma_{i_k}^{k+1} = \sigma_{i_k}^k$  when  $\|g_{i_k}^k\| = 0$ . Blocks  $\sigma_{i_k}^k$  that are updated at each iteration can be chosen through any deterministic or randomized rule, and in this paper we focus on three coordinate selection rules:

- Uniform sampling:  $i_k = i$  with probability  $p_i = 1/n$ .
- Importance sampling:  $i_k = i$  with probability  $p_i = \|g_i^k\| / \sum_{j=1}^n \|g_j^k\|$ .
- Greedy coordinate selection:  $i_k = \arg \max_{i \in [n]} (\|g_i^k\| - \langle \sigma_i^k, g_i^k \rangle)$ .

Per iteration computational cost of the BCM algorithm with uniform sampling is  $\mathcal{O}(nr)$  as after  $i_k$  is chosen uniformly at random,  $g_{i_k}^k$  can be computed in  $2(n-1)r$  floating point operations. On the other hand, the BCM algorithm with importance sampling and greedy coordinate selection requires all  $\{\|g_i^k\|\}_{i=1}^n$ , which can be naively computed in  $\mathcal{O}(n^2r)$  floating point operations per iteration. Instead, a smarter implementation is to keep both  $\{\sigma_i^k\}_{i=1}^n$  and  $\{g_i^k\}_{i=1}^n$ 's in the memory (only the current iterates, not all the past ones) and update them as presented in Algorithm 1, which can be done in  $2(n-1)r$  floating point operations. Therefore, per iteration computational cost of the BCM method with all three coordinate selection rules is  $\mathcal{O}(nr)$  for dense  $A$  (i.e., when no structure is available on  $A$ ). However, in many SDP applications (such as

Max-Cut and graphical model inference),  $A$  is induced by a graph and letting  $d$  denote the maximum degree of the graph that induces  $A$ , the computational cost of the BCM algorithm becomes  $\mathcal{O}(dr)$ . In comparison, per iteration computational complexity of the Riemannian gradient ascent algorithm is  $\mathcal{O}(n^2r)$ , whereas the Riemannian trust-region algorithm runs a few iterations of a subroutine (e.g., power method) to solve the trust-region subproblem, whose per iteration cost is typically  $\mathcal{O}(n^2r)$ .

Using any of the coordinate selection rules, the function values of the iterates generated by the BCM algorithm is a non-decreasing sequence by the definition of the algorithm. The increase in the function value per iteration (before reaching to stationarity) can be explicitly computed as we present in the following lemma.

**Lemma 1** *Suppose at the  $k$ th iteration of the BCM algorithm,  $i_k$ th block is chosen (with some coordinate selection rule). Then, the BCM algorithm yields the following ascent on the objective value:*

$$f(\sigma^{k+1}) - f(\sigma^k) = 2 \left( \|g_{i_k}^k\| - \langle \sigma_{i_k}^k, g_{i_k}^k \rangle \right) \geq 0.$$

**Proof** According to the decomposition in (1), we can compute the objective function as follows:

$$\begin{aligned} f(\sigma^{k+1}) &= 2\langle \sigma_{i_k}^{k+1}, g_{i_k}^{k+1} \rangle + \sum_{i \neq i_k} \sum_{j \neq i, i_k} A_{ij} \langle \sigma_i^{k+1}, \sigma_j^{k+1} \rangle, \\ &= 2\langle \sigma_{i_k}^{k+1}, g_{i_k}^k \rangle + \sum_{i \neq i_k} \sum_{j \neq i, i_k} A_{ij} \langle \sigma_i^k, \sigma_j^k \rangle, \end{aligned} \tag{3}$$

where the latter equality follows since  $g_{i_k}^{k+1} = g_{i_k}^k$  and all the terms in the sum are independent of  $\sigma_{i_k}^{k+1}$ . After adding and subtracting  $2\langle \sigma_{i_k}^k, g_{i_k}^k \rangle$  to the right-hand side of (3), we obtain

$$f(\sigma^{k+1}) = f(\sigma^k) + 2 \left( \langle \sigma_{i_k}^{k+1}, g_{i_k}^k \rangle - \langle \sigma_{i_k}^k, g_{i_k}^k \rangle \right).$$

By the update rule of the algorithm, we have  $\sigma_{i_k}^{k+1} = g_{i_k}^k / \|g_{i_k}^k\|$ , and plugging this value in the above equation concludes the proof.  $\square$

### 3 Convergence rate of BCM

In this section, we analyze the convergence rate of the BCM algorithm. As the feasible set of the problem in (Non-CVX) defines a smooth manifold, we will use certain tools from manifold optimization throughout the paper, which are highlighted in the following subsection. We refer to [2, Section 5.4] for a more detailed treatment of this topic.

### 3.1 Riemannian geometry of the problem

We define the following submanifold of matrices  $\mathbb{R}^{n \times r}$  that corresponds to the Riemannian geometry induced by the constraints of the problem (Non-CVX) in the Euclidean space:

$$\mathcal{M}_r := \left\{ \boldsymbol{\sigma} = (\sigma_1, \dots, \sigma_n)^\top \in \mathbb{R}^{n \times r} : \|\sigma_i\| = 1, \forall i \in [n] \right\}.$$

This manifold represents the Cartesian product of  $n$  unit spheres in  $\mathbb{R}^r$ . For any given point  $\boldsymbol{\sigma} \in \mathcal{M}_r$ , its tangent space can be found by taking the differential of the equality constraints as follows

$$T_{\boldsymbol{\sigma}}\mathcal{M}_r := \left\{ \mathbf{u} = (u_1, \dots, u_n)^\top \in \mathbb{R}^{n \times r} : \langle u_i, \sigma_i \rangle = 0, \forall i \in [n] \right\}.$$

The Riemannian gradient of  $f$  on this manifold can be computed by the projection of its Euclidean gradient onto the tangent bundle. In particular, let  $\mathbf{P}_{\boldsymbol{\sigma}}^\perp : \mathbb{R}^{n \times r} \rightarrow T_{\boldsymbol{\sigma}}\mathcal{M}_r$  denote the projection operator from the Euclidean space to the tangent space of  $\boldsymbol{\sigma}$ . When applied to a given matrix  $\mathbf{w} = (w_1, \dots, w_n)^\top \in \mathbb{R}^{n \times r}$ , this projection operator yields

$$\begin{aligned} \mathbf{P}_{\boldsymbol{\sigma}}^\perp(\mathbf{w}) &= (w_1 - \langle \sigma_1, w_1 \rangle \sigma_1, \dots, w_n - \langle \sigma_n, w_n \rangle \sigma_n)^\top, \\ &= \mathbf{w} - \text{Diag}(\text{diag}(\mathbf{w}\boldsymbol{\sigma}^\top))\boldsymbol{\sigma}. \end{aligned}$$

Therefore, the Riemannian gradient of  $f$  at  $\boldsymbol{\sigma}$  can be computed as follows

$$\text{grad} f(\boldsymbol{\sigma}) = \mathbf{P}_{\boldsymbol{\sigma}}^\perp(\nabla f(\boldsymbol{\sigma})) = 2(\mathbf{A} - \boldsymbol{\Lambda})\boldsymbol{\sigma},$$

where  $\boldsymbol{\Lambda} = \text{Diag}(\text{diag}(\mathbf{A}\boldsymbol{\sigma}\boldsymbol{\sigma}^\top))$ . Or equivalently, the Riemannian gradient of  $f$  at  $\boldsymbol{\sigma}$  can be explicitly expressed as follows

$$\text{grad} f(\boldsymbol{\sigma}) = 2(g_1 - \langle \sigma_1, g_1 \rangle \sigma_1, \dots, g_n - \langle \sigma_n, g_n \rangle \sigma_n)^\top,$$

and its magnitude is given by

$$\|\text{grad} f(\boldsymbol{\sigma})\|_{\mathbb{F}}^2 = 2 \sum_{i=1}^n \|g_i - \langle \sigma_i, g_i \rangle \sigma_i\|^2 = 2 \sum_{i=1}^n \left( \|g_i\|^2 - \langle \sigma_i, g_i \rangle^2 \right). \quad (4)$$

Using the same approach, we can calculate the Riemannian Hessian of  $f$  at  $\boldsymbol{\sigma}$  along the direction of a vector  $\mathbf{u} \in T_{\boldsymbol{\sigma}}\mathcal{M}_r$  by projecting the directional derivative of the gradient vector field onto the tangent space of  $\boldsymbol{\sigma}$  as follows

$$\text{Hess} f(\boldsymbol{\sigma})[\mathbf{u}] = \mathbf{P}_{\boldsymbol{\sigma}}^\perp(\mathbf{D} \text{grad} f(\boldsymbol{\sigma})[\mathbf{u}]),$$

where  $D \operatorname{grad} f(\boldsymbol{\sigma})[\mathbf{u}]$  denotes the directional gradient of  $\operatorname{grad} f(\boldsymbol{\sigma})$  along the direction  $\mathbf{u}$ . This yields

$$\begin{aligned} \operatorname{Hess} f(\boldsymbol{\sigma})[\mathbf{u}] &= \mathbf{P}^\perp \left( 2(\mathbf{A} - \boldsymbol{\Lambda})\mathbf{u} - 2\operatorname{Diag}(\operatorname{diag}(\mathbf{A}\boldsymbol{\sigma}\mathbf{u}^\top + \mathbf{A}\mathbf{u}\boldsymbol{\sigma}^\top))\boldsymbol{\sigma} \right) \\ &= \mathbf{P}^\perp (2(\mathbf{A} - \boldsymbol{\Lambda})\mathbf{u}), \end{aligned} \tag{5}$$

and in particular, for any  $\mathbf{u} \in T_\sigma \mathcal{M}_r$ , we have

$$\langle \mathbf{u}, \operatorname{Hess} f(\boldsymbol{\sigma})[\mathbf{u}] \rangle = 2\langle \mathbf{u}, (\mathbf{A} - \boldsymbol{\Lambda})\mathbf{u} \rangle. \tag{6}$$

The geodesics  $t \mapsto \boldsymbol{\sigma}(t)$  (i.e., curves of shortest path with zero acceleration) can be expressed as a function of  $\boldsymbol{\sigma} = \boldsymbol{\sigma}(0) \in \mathcal{M}_r$  and  $\mathbf{u} \in T_\sigma \mathcal{M}_r$  as follows

$$\sigma_i(t) = \sigma_i \cos(\|\mathbf{u}_i\|t) + \frac{u_i}{\|\mathbf{u}_i\|} \sin(\|\mathbf{u}_i\|t). \tag{7}$$

This geodesic can be thought as the curve on the manifold that are obtained by moving from  $\boldsymbol{\sigma} \in \mathcal{M}_r$  towards the direction pointed by  $\mathbf{u} \in T_\sigma \mathcal{M}_r$  (Table 1). According to this definition, the exponential map  $\operatorname{Exp}_\sigma : T_\sigma \mathcal{M}_r \rightarrow \mathcal{M}_r$  corresponds to evaluating the point at  $t = 1$  on the geodesic function, i.e., letting  $\boldsymbol{\sigma}' = \operatorname{Exp}_\sigma(\mathbf{u})$ , where  $\mathbf{u} \in T_\sigma \mathcal{M}_r$ , we have

$$\sigma'_i = \sigma_i \cos(\|\mathbf{u}_i\|) + \frac{u_i}{\|\mathbf{u}_i\|} \sin(\|\mathbf{u}_i\|).$$

According to this geodesic map, we can also define the following geodesic distance between two points  $\boldsymbol{\sigma}$  and  $\boldsymbol{\sigma}'$  on the manifold:

$$\operatorname{dist}(\boldsymbol{\sigma}, \boldsymbol{\sigma}') = \left( \sum_{i=1}^n (\arccos \langle \sigma_i, \sigma'_i \rangle)^2 \right)^{1/2}. \tag{8}$$

More specifically, letting  $\boldsymbol{\sigma}' = \operatorname{Exp}_\sigma(\mathbf{u})$ , we obtain

$$\operatorname{dist}(\boldsymbol{\sigma}, \boldsymbol{\sigma}') = \left( \sum_{i=1}^n (\arccos \langle \sigma_i, \sigma_i \cos \|\mathbf{u}_i\| \rangle)^2 \right)^{1/2} = \|\mathbf{u}\|_F.$$

Similarly, the distance between a point  $\boldsymbol{\sigma}$  and a non-empty, closed and (geodesically) convex set  $\Omega$  can be found as

$$\operatorname{dist}(\boldsymbol{\sigma}, \Omega) = \min_{\boldsymbol{\sigma}' \in \Omega} \operatorname{dist}(\boldsymbol{\sigma}, \boldsymbol{\sigma}').$$

**Table 1** Summary of certain definitions stated in Sect. 3.1

Projection to the tangent space $T_{\sigma} \mathcal{M}_r$ at $\sigma$	$P_{\sigma}^{\perp}(\mathbf{w}) = \mathbf{w} - \text{Diag}(\text{diag}(\mathbf{w}\sigma^{\top}))\sigma$
Riemannian gradient at $\sigma$	$\text{grad} f(\sigma) = 2(\mathbf{A} - \mathbf{A})\sigma$ where $\mathbf{A} = \text{Diag}(\text{diag}(\mathbf{A}\sigma\sigma^{\top}))$
Riemannian Hessian at $\sigma$ along $\mathbf{u} \in T_{\sigma} \mathcal{M}_r$	$\text{Hess} f(\sigma)[\mathbf{u}] = P^{\perp}(2(\mathbf{A} - \mathbf{A})\mathbf{u})$
Geodesic $t \rightarrow \sigma(t)$	$\sigma_i(t) = \sigma_i \cos(\ u_i\ t) + \frac{u_i}{\ u_i\ } \sin(\ u_i\ t)$
Exponential map $\sigma' = \text{Exp}_{\sigma}(\mathbf{u})$	$\sigma'_i = \sigma_i \cos(\ u_i\ ) + \frac{u_i}{\ u_i\ } \sin(\ u_i\ )$

### 3.2 Global rate of convergence

In this section, we show that the BCM algorithm is globally convergent to a first-order stationary point of the problem with a sublinear rate. In particular, in the following theorem, we consider the BCM algorithm with greedy coordinate selection and show that its functional ascent (see Lemma 1) can be related to the norm of the Riemannian gradient of the function evaluated at the current iterate. By doing so, we prove that the BCM algorithm returns a solution with arbitrarily small Riemannian gradient.

**Theorem 1** *Let  $f^* = \max_{\|\sigma_i\|=1, \forall i \in [n]} f(\sigma)$ . Then, for any  $K \geq 1$ , BCM with greedy coordinate selection yields the following guarantee*

$$\min_{k \in [K-1]} \|\text{grad} f(\sigma^k)\|_{\mathbb{F}}^2 \leq \frac{2n\|\mathbf{A}\|_1(f^* - f(\sigma^0))}{K}. \tag{9}$$

**Proof** From Lemma 1, we have

$$\begin{aligned} f(\sigma^{k+1}) - f(\sigma^k) &= 2 \left( \|g_{i_k}^k\| - \langle \sigma_{i_k}^k, g_{i_k}^k \rangle \right) \\ &= 2 \max_{i \in [n]} \left( \|g_i^k\| - \langle \sigma_i^k, g_i^k \rangle \right), \end{aligned}$$

where the latter equality follows by the greedy coordinate selection rule. We can rewrite this equation as follows:

$$\begin{aligned} f(\sigma^{k+1}) - f(\sigma^k) &= \max_{i \in [n]} \frac{2\|g_i^k\| \left( \|g_i^k\| - \langle \sigma_i^k, g_i^k \rangle \right)}{\|g_i^k\|}, \\ &\geq \max_{i \in [n]} \frac{\|g_i^k\|^2 - \langle \sigma_i^k, g_i^k \rangle^2}{\|g_i^k\|}, \end{aligned}$$

where the inequality follows since  $\|g_i^k\| \geq \langle \sigma_i^k, g_i^k \rangle$  for all  $\sigma_i^k \in \mathbb{R}^r$ . Lower bounding the maximum with the mean of its arguments, we get

$$f(\sigma^{k+1}) - f(\sigma^k) \geq \frac{1}{n} \sum_{i=1}^n \frac{\|g_i^k\|^2 - \langle \sigma_i^k, g_i^k \rangle^2}{\|g_i^k\|}. \tag{10}$$

The  $\|g_i^k\|$  term in the denominator in (10) can be upper bounded as follows

$$\|g_{i_k}^k\| \leq \sum_{j \neq i_k} |A_{i_k j}| \|\sigma_j^k\| \leq \|\mathbf{A}\|_1. \tag{11}$$

Using this bound in (10), we get

$$\begin{aligned} f(\sigma^{k+1}) - f(\sigma^k) &\geq \frac{1}{n\|\mathbf{A}\|_1} \sum_{i=1}^n \left( \|g_i^k\|^2 - \langle \sigma_i^k, g_i^k \rangle \right) \\ &= \frac{\|\text{grad} f(\sigma^k)\|_{\mathbb{F}}^2}{2n\|\mathbf{A}\|_1}. \end{aligned} \tag{12}$$

In order to conclude (9), we assume the contrary that  $\|\text{grad} f(\sigma^k)\|_{\mathbb{F}}^2 > \epsilon$  for all  $k \in [K - 1]$ . Then, using the boundedness of  $f$ , we observe that

$$f^* - f(\sigma^0) \geq f(\sigma^K) - f(\sigma^0) = \sum_{k=0}^{K-1} \left[ f(\sigma^{k+1}) - f(\sigma^k) \right].$$

Using the functional ascent bound of BCM in (12), we get

$$f^* - f(\sigma^0) \geq \sum_{k=0}^{K-1} \frac{\|\text{grad} f(\sigma^k)\|_{\mathbb{F}}^2}{2n\|\mathbf{A}\|_1} > \frac{K\epsilon}{2n\|\mathbf{A}\|_1},$$

where the latter inequality follows by the assumption. Then, by contradiction, the algorithm returns a solution with  $\|\text{grad} f(\sigma^k)\|_{\mathbb{F}}^2 \leq \epsilon$ , for some  $k \in [K - 1]$ , provided that

$$K \geq \frac{2n\|\mathbf{A}\|_1(f^* - f(\sigma^0))}{\epsilon}.$$

□

Using a similar approach to Theorem 1, we show in the following corollary that the BCM algorithm with uniform and importance sampling attains a similar sublinear convergence rate in expectation. The proof of this corollary follows similar lines to the proof of Theorem 1, hence is deferred to ‘‘Appendix A’’.

**Corollary 1** *Let  $f^* = \max_{\|\sigma_i\|=1, \forall i \in [n]} f(\sigma)$ . Then, for any  $K \geq 1$ , randomized BCM yields the following guarantee*

$$\min_{k \in [K-1]} \mathbb{E} \|\text{grad} f(\sigma^k)\|_{\mathbb{F}}^2 \leq \frac{2L(f^* - f(\sigma^0))}{K}, \tag{13}$$

where

$$L = \begin{cases} n\|\mathbf{A}\|_1, & \text{for uniform sampling,} \\ \|\mathbf{A}\|_{1,1}, & \text{for importance sampling.} \end{cases} \quad (14)$$

We can observe from (9), (13) and (14) that BCM with uniform sampling attains the same sublinear rate as BCM with greedy coordinate selection in expectation as they both require at most  $\lceil (2n\|\mathbf{A}\|_1(f^* - f(\boldsymbol{\sigma}^0)))/\epsilon \rceil$  iterations to return a solution  $\boldsymbol{\sigma}$  satisfying  $\|\text{grad} f(\boldsymbol{\sigma})\|_{\mathbb{F}}^2 \leq \epsilon$ . On the other hand, we see that BCM with importance sampling enjoys a tighter convergence rate compared to BCM with uniform sampling, as  $\|\mathbf{A}\|_{1,1} \leq n\|\mathbf{A}\|_1$  for all  $\mathbf{A} \in \mathbb{R}^{n \times n}$ .

### 3.3 Local rate of convergence

Although the BCM algorithm enjoys the sublinear convergence rates presented in Sect. 3.2, it is numerically observed that the rate of convergence is linear when  $\boldsymbol{\sigma}^k$  is close to a local maximum [25,43]. In this section, we investigate this behavior and prove that indeed BCM attains a linear convergence rate around a local maximum under the quadratic decay condition on the objective function, which is classically defined as follows [4,8]: Consider the unconstrained maximization problem:  $\max_x \varphi(x)$ , and let  $\Omega_{\bar{x}}$  denote the set of local maximizers with objective value  $\varphi(\bar{x})$ . Then, the quadratic decay condition is said to be satisfied at  $\bar{x}$  for  $\varphi$ , if there exists constants  $\mu, \delta > 0$  such that  $\varphi(x) \leq \varphi(\bar{x}) - \mu \text{dist}^2(x, \Omega_{\bar{x}})$ , for all  $x$  such that  $\|x - \bar{x}\| \leq \delta$ , where  $\text{dist}$  measures the distance between point  $x$  and set  $\Omega_{\bar{x}}$ .

For the constrained optimization problem that we are considering in (Non-CVX), this definition needs to be slightly reworked. In particular, let  $\boldsymbol{\sigma}$  be a local maximum of (Non-CVX) and consider the Taylor expansion of  $\text{Exp}_{\boldsymbol{\sigma}}(\mathbf{u})$  around  $\boldsymbol{\sigma}$ :

$$f(\text{Exp}_{\boldsymbol{\sigma}}(\mathbf{u})) = f(\boldsymbol{\sigma}) + \frac{1}{2}\langle \mathbf{u}, \text{Hess} f(\boldsymbol{\sigma})[\mathbf{u}] \rangle + \mathcal{O}(\|\mathbf{u}\|_{\mathbb{F}}^3),$$

where the first-order term is zero as  $\boldsymbol{\sigma}$  is a local maximum. Then, for a sufficiently small neighborhood of  $\boldsymbol{\sigma}$ , the quadratic decay condition is satisfied if and only if there exists a constant  $\mu > 0$  such that  $\langle \mathbf{u}, \text{Hess} f(\boldsymbol{\sigma})[\mathbf{u}] \rangle \leq -\mu \text{dist}^2(\text{Exp}_{\boldsymbol{\sigma}}(\mathbf{u}), \Omega_{\boldsymbol{\sigma}})$ , for all  $\text{Exp}_{\boldsymbol{\sigma}}(\mathbf{u})$  sufficiently close to  $\boldsymbol{\sigma}$ , where  $\Omega_{\boldsymbol{\sigma}}$  is the set on which  $f$  has constant value  $f(\boldsymbol{\sigma})$ . Assume for the sake of simplicity that  $\boldsymbol{\sigma}$  is a strict local maximum, i.e.,  $\Omega_{\boldsymbol{\sigma}} = \{\boldsymbol{\sigma}\}$ . Then, the distance between  $\text{Exp}_{\boldsymbol{\sigma}}(\mathbf{u})$  and  $\boldsymbol{\sigma}$  can be found as the norm of the tangent vector that connects these two points via the geodesic curve, i.e.,  $\text{dist}(\text{Exp}_{\boldsymbol{\sigma}}(\mathbf{u}), \boldsymbol{\sigma}) = \|\mathbf{u}\|_{\mathbb{F}}$ . Therefore, the quadratic decay condition is satisfied if and only if there exists a constant  $\mu > 0$  such that  $\langle \mathbf{u}, \text{Hess} f(\boldsymbol{\sigma})[\mathbf{u}] \rangle \leq -\mu \|\mathbf{u}\|_{\mathbb{F}}^2$  for all  $\mathbf{u} \in T_{\boldsymbol{\sigma}}\mathcal{M}_r$ , where we note that the condition that  $\text{Exp}_{\boldsymbol{\sigma}}(\mathbf{u})$  is sufficiently close to  $\boldsymbol{\sigma}$  is dropped considering the limit as  $\mathbf{u} \rightarrow \mathbf{0}$ .

Unfortunately, no local maximum is a strict local maximum for the problem (Non-CVX). To observe this, let  $\text{O}(r) = \{\mathbf{Q} \in \mathbb{R}^{r \times r} : \mathbf{Q}^{\top} \mathbf{Q} = \mathbf{Q} \mathbf{Q}^{\top} = \mathbf{I}\}$  denote the orthogonal group in dimension  $r$ . Then, it can be observed that  $f(\boldsymbol{\sigma} \mathbf{Q}) =$

$\langle \mathbf{A}, \boldsymbol{\sigma} \mathbf{Q} \mathbf{Q}^\top \boldsymbol{\sigma}^\top \rangle = \langle \mathbf{A}, \boldsymbol{\sigma} \boldsymbol{\sigma}^\top \rangle = f(\boldsymbol{\sigma})$ , for any  $\mathbf{Q} \in O(r)$ . Therefore, in order to measure the distance between  $\text{Exp}_{\boldsymbol{\sigma}}(\mathbf{u})$  and  $\Omega_{\boldsymbol{\sigma}}$ , we define the following equivalence relation  $\sim$ :

$$\boldsymbol{\sigma} \sim \boldsymbol{\sigma}' \iff \exists \mathbf{Q} \in O(r) : \boldsymbol{\sigma} = \boldsymbol{\sigma}' \mathbf{Q}. \tag{15}$$

This equivalence relation induces a quotient space denoted by  $\mathcal{M}_r / \sim$  and we let  $[\boldsymbol{\sigma}]$  denote the equivalence class of a given matrix  $\boldsymbol{\sigma} \in \mathcal{M}_r$ . According to this definition,  $f$  has constant value of  $f(\boldsymbol{\sigma})$  on the set  $[\boldsymbol{\sigma}]$ , i.e.,  $\Omega_{\boldsymbol{\sigma}} = [\boldsymbol{\sigma}]$ . We let  $\mathcal{V}_{\boldsymbol{\sigma}} \subset T_{\boldsymbol{\sigma}} \mathcal{M}_r$  denote the tangent space to the equivalence class  $[\boldsymbol{\sigma}]$ , which can be found as  $\mathcal{V}_{\boldsymbol{\sigma}} = \{\boldsymbol{\sigma} \mathbf{B} : \mathbf{B} \in \mathbb{R}^{r \times r} \text{ and } \mathbf{B}^\top = -\mathbf{B}\}$ .<sup>1</sup> Therefore,  $\text{dist}(\text{Exp}_{\boldsymbol{\sigma}}(\mathbf{u}), [\boldsymbol{\sigma}]) = \|\mathbf{u}\|_F$  if the closest point to  $\text{Exp}_{\boldsymbol{\sigma}}(\mathbf{u})$  in  $[\boldsymbol{\sigma}]$  is  $\boldsymbol{\sigma}$ , or equivalently  $\text{dist}(\text{Exp}_{\boldsymbol{\sigma}}(\mathbf{u}), [\boldsymbol{\sigma}]) = \|\mathbf{u}\|_F$  if  $\mathbf{u} \in T_{\boldsymbol{\sigma}} \mathcal{M}_r \setminus \mathcal{V}_{\boldsymbol{\sigma}}$ . Consequently, we say that quadratic decay is satisfied at  $\boldsymbol{\sigma}$  for  $f$  if  $\text{Hess} f(\boldsymbol{\sigma})$  is negative definite on the orthogonal complement of  $\mathcal{V}_{\boldsymbol{\sigma}}$  in  $T_{\boldsymbol{\sigma}} \mathcal{M}_r$ . The formal statement of this definition is as follows.

**Definition 1** (*Quadratic decay*) Let  $\boldsymbol{\sigma}$  be a local maximum of (Non-CVX). Quadratic decay condition is said to be satisfied at  $\boldsymbol{\sigma}$  for  $f$  if there exists a constant  $\mu > 0$  such that

$$\langle \mathbf{u}, \text{Hess} f(\boldsymbol{\sigma})[\mathbf{u}] \rangle \leq -\mu \|\mathbf{u}\|_F^2, \text{ for all } \mathbf{u} \in T_{\boldsymbol{\sigma}} \mathcal{M}_r \setminus \mathcal{V}_{\boldsymbol{\sigma}}, \tag{16}$$

where  $\mathcal{V}_{\boldsymbol{\sigma}}$  is the tangent space to the equivalence class  $[\boldsymbol{\sigma}]$ .

In the following theorem, we present the linear convergence rate of the BCM algorithm under the quadratic decay condition. We defer the validity of this condition to Sect. 3.4 where we show that quadratic decay generically (over the set of matrices  $\mathbf{A}$ ) holds for  $f$  when  $r$  is sufficiently large.

**Theorem 2** Let  $\bar{\boldsymbol{\sigma}}$  be a limit point of the BCM algorithm and assume that  $\bar{\boldsymbol{\sigma}}$  is a local maximum that satisfies the quadratic decay condition. If  $\boldsymbol{\sigma}^0$  is sufficiently close to the equivalent class  $[\bar{\boldsymbol{\sigma}}]$ , then the iterates generated by the BCM algorithm with greedy coordinate selection enjoy the following linear convergence rate

$$f(\bar{\boldsymbol{\sigma}}) - f(\boldsymbol{\sigma}^{k+1}) \leq \left(1 - \frac{\mu}{4n^2 \|\mathbf{A}\|_1}\right) \left(f(\bar{\boldsymbol{\sigma}}) - f(\boldsymbol{\sigma}^k)\right). \tag{17}$$

**Proof** We first discuss the outline of the proof for clarity. By (12), we have the following functional ascent bound on the iterates of the algorithm

$$f(\boldsymbol{\sigma}^{k+1}) - f(\boldsymbol{\sigma}^k) \geq \frac{\|\text{grad} f(\boldsymbol{\sigma}^k)\|_F^2}{2n \|\mathbf{A}\|_1}. \tag{18}$$

In order to prove linear convergence, our aim is to show that  $\|\text{grad} f(\boldsymbol{\sigma}^k)\|_F^2 \geq c(f(\bar{\boldsymbol{\sigma}}) - f(\boldsymbol{\sigma}^k))$  for some positive constant  $c$  such that  $c < 2n \|\mathbf{A}\|_1$ , in a neighborhood around the limit points of the iterates generated by the algorithm. To prove this,

<sup>1</sup> Note that the dimension of  $\mathcal{V}_{\boldsymbol{\sigma}}$  depends on the rank of  $\boldsymbol{\sigma}$ , and hence the quotient space is not a manifold.

we consider the Taylor approximation of  $\|\text{grad} f(\sigma^k)\|_{\mathbb{F}}^2$  and  $f(\sigma^k)$  around  $\sigma \in [\bar{\sigma}]$ , where  $\sigma$  is the closest point to  $\sigma^k$  in the set  $\bar{\sigma}$ . In the remainder of this proof, we show that the desired inequality holds by relating the most significant terms in these Taylor expansions. We defer bounding the higher-order terms to “Appendix B” in order not to distract the reader from the content of the paper.

Let  $\bar{\sigma}$  be the limit point of a subsequence  $\{\sigma^{k_\ell}\}_{k_\ell \geq 0}$  that contains  $\sigma^k$ . Then, we consider the solution  $\sigma \in [\bar{\sigma}]$  such that  $\sigma$  is the projection of  $\sigma^k$  onto  $[\bar{\sigma}]$ , i.e.,  $\text{dist}(\sigma, \sigma^k) \leq \text{dist}(\sigma', \sigma^k)$  for all  $\sigma' \in [\bar{\sigma}]$ . Then, by construction there exists  $\bar{u} \in T_\sigma \mathcal{M}_r \setminus \mathcal{V}_\sigma$  such that  $\text{Exp}_\sigma(\bar{u}) = \sigma^k$ . For ease of presentation, we let  $u = \bar{u}/\|\bar{u}\|_{\mathbb{F}}$  denote the normalized tangent vector and consider the following geodesic to describe  $\sigma^k$ :

$$\sigma_i^k = \sigma_i \cos(\|u_i\|t) + \frac{u_i}{\|u_i\|} \sin(\|u_i\|t), \quad (19)$$

where it can be observed that  $t = \|\bar{u}\|_{\mathbb{F}}$  recovers the original exponential map  $\sigma^k = \text{Exp}_\sigma(\bar{u})$ . The second order Taylor approximation to (19) yields (note that  $t = \|\bar{u}\|_{\mathbb{F}} < 1$ , when  $\sigma$  and  $\sigma^k$  are sufficiently close):

$$\sigma_i^k = \sigma_i + t u_i - \frac{t^2}{2} \|u_i\|^2 \sigma_i + \mathcal{O}(t^3),$$

and using this approximation, we obtain

$$g_i^k = g_i + t v_i - \frac{t^2}{2} \tilde{g}_i + \mathcal{O}(t^3),$$

where

$$v_i^k = \sum_{j \neq i} A_{ij} u_j \quad \text{and} \quad \tilde{g}_i = \sum_{j \neq i} A_{ij} \|u_j\|^2 \sigma_j.$$

This yields the following Taylor approximation to  $\|\text{grad} f(\sigma^k)\|_{\mathbb{F}}^2$ :

$$\begin{aligned} \|\text{grad} f(\sigma^k)\|_{\mathbb{F}}^2 &= 2 \sum_{i=1}^n \left( \|g_i^k\|^2 - \langle \sigma_i^k, g_i^k \rangle^2 \right) \\ &= 2 \sum_{i=1}^n \left( \|g_i + t v_i - \frac{t^2}{2} \tilde{g}_i\|^2 \right. \\ &\quad \left. - \langle \sigma_i + t u_i - \frac{t^2}{2} \|u_i\|^2 \sigma_i, g_i + t v_i - \frac{t^2}{2} \tilde{g}_i \rangle^2 \right) + \mathcal{O}(t^3), \\ &= 2 \sum_{i=1}^n \left\{ \|g_i\|^2 + 2t \langle g_i, v_i \rangle - t^2 \langle g_i, \tilde{g}_i \rangle + t^2 \|v_i\|^2 \right. \end{aligned}$$

$$\begin{aligned}
 & - \left( \langle \sigma_i, g_i \rangle + t \langle \sigma_i, v_i \rangle - \frac{t^2}{2} \langle \sigma_i, \tilde{g}_i \rangle \right. \\
 & \left. + t \langle u_i, g_i \rangle + t^2 \langle u_i, v_i \rangle - \frac{t^2}{2} \|u_i\|^2 \langle \sigma_i, g_i \rangle \right)^2 \Big\} + \mathcal{O}(t^3).
 \end{aligned}$$

Observe that as  $\sigma$  is a local maximum, we have  $\sigma_i = g_i/\|g_i\|$  for all  $i \in [n]$ . This follows since the first-order stationarity condition implies  $\sigma_i = \pm g_i/\|g_i\|$  for all  $i \in [n]$ ; and having  $\sigma_i = -g_i/\|g_i\|$  for some  $i \in [n]$  conflicts with the assumption that  $\sigma$  is a local maximum as replacing  $\sigma_i$  with any other feasible point on the sphere increases the objective function. We also have that  $\langle \sigma_i, u_i \rangle = 0$  for all  $i \in [n]$ , as  $u \in T_\sigma \mathcal{M}_r$ . Using these facts in the above equality, we get

$$\begin{aligned}
 \|\text{grad} f(\sigma^k)\|_{\mathbb{F}}^2 &= 2 \sum_{i=1}^n \left[ \|g_i\|^2 + 2t \|g_i\| \langle \sigma_i, v_i \rangle - t^2 \|g_i\| \langle \sigma_i, \tilde{g}_i \rangle + t^2 \|v_i\|^2 \right. \\
 & \quad - \left( \|g_i\| + t \langle \sigma_i, v_i \rangle - \frac{t^2}{2} \langle \sigma_i, \tilde{g}_i \rangle \right. \\
 & \quad \left. \left. + t^2 \langle u_i, v_i \rangle - \frac{t^2}{2} \|u_i\|^2 \|g_i\| \right)^2 \right] + \mathcal{O}(t^3), \\
 &= 2 \sum_{i=1}^n \left[ \|g_i\|^2 + 2t \|g_i\| \langle \sigma_i, v_i \rangle - t^2 \|g_i\| \langle \sigma_i, \tilde{g}_i \rangle + t^2 \|v_i\|^2 \right. \\
 & \quad - \left( \|g_i\|^2 + 2t \|g_i\| \langle \sigma_i, v_i \rangle - t^2 \|g_i\| \langle \sigma_i, \tilde{g}_i \rangle + 2t^2 \|g_i\| \langle u_i, v_i \rangle \right. \\
 & \quad \left. \left. - t^2 \|u_i\|^2 \|g_i\|^2 + t^2 \langle \sigma_i, v_i \rangle^2 \right) \right] + \mathcal{O}(t^3), \\
 &= 2t^2 \sum_{i=1}^n \left( \|v_i\|^2 - \langle \sigma_i, v_i \rangle^2 \right. \\
 & \quad \left. - 2 \|g_i\| \langle u_i, v_i \rangle + \|u_i\|^2 \|g_i\|^2 \right) + \mathcal{O}(t^3). \tag{20}
 \end{aligned}$$

Since  $\langle \sigma_i, u_i \rangle = 0$  for all  $i \in [n]$ , we have by the Pythagorean theorem that

$$\|v_i\|^2 - \langle \sigma_i, v_i \rangle^2 - \left\langle \frac{u_i}{\|u_i\|}, v_i \right\rangle^2 \geq 0.$$

Using this inequality in (20), we get

$$\begin{aligned}
 \|\text{grad} f(\sigma^k)\|_{\mathbb{F}}^2 &\geq 2t^2 \sum_{i=1}^n \left( \left\langle \frac{u_i}{\|u_i\|}, v_i \right\rangle^2 - 2 \|g_i\| \langle u_i, v_i \rangle + \|u_i\|^2 \|g_i\|^2 \right) + \mathcal{O}(t^3), \\
 &= 2t^2 \sum_{i=1}^n \left( \|u_i\| \|g_i\| - \left\langle \frac{u_i}{\|u_i\|}, v_i \right\rangle \right)^2 + \mathcal{O}(t^3). \tag{21}
 \end{aligned}$$

In order to lower bound (21) by  $c(f(\boldsymbol{\sigma}) - f(\boldsymbol{\sigma}^k))$ , we consider the second order Taylor approximation of  $f(\boldsymbol{\sigma}^k)$ , which can be written as follows

$$\begin{aligned} f(\boldsymbol{\sigma}^k) &= \sum_{i=1}^n \langle \sigma_i^k, g_i^k \rangle, \\ &= \sum_{i=1}^n \langle \sigma_i + t u_i - \frac{t^2}{2} \|u_i\|^2 \sigma_i, g_i + t v_i - \frac{t^2}{2} \tilde{g}_i \rangle + \mathcal{O}(t^3), \\ &= \sum_{i=1}^n \left( \langle \sigma_i, g_i \rangle + t \langle \sigma_i, v_i \rangle - \frac{t^2}{2} \langle \sigma_i, \tilde{g}_i \rangle + t \langle u_i, g_i \rangle + t^2 \langle u_i, v_i \rangle \right. \\ &\quad \left. - \frac{t^2}{2} \|u_i\|^2 \langle \sigma_i, g_i \rangle \right) + \mathcal{O}(t^3). \end{aligned}$$

Similar to the previous derivations, using the fact that  $\sigma_i = g_i / \|g_i\|$  and  $\langle \sigma_i, u_i \rangle = 0$  for all  $i \in [n]$ , we obtain

$$\begin{aligned} f(\boldsymbol{\sigma}^k) &= f(\boldsymbol{\sigma}) + \sum_{i=1}^n \left( t \langle \sigma_i, v_i \rangle - \frac{t^2}{2} \langle \sigma_i, \tilde{g}_i \rangle + t^2 \langle u_i, v_i \rangle - \frac{t^2}{2} \|u_i\|^2 \langle \sigma_i, g_i \rangle \right) + \mathcal{O}(t^3), \\ &= f(\boldsymbol{\sigma}) + \sum_{i=1}^n \left( t \sum_{j \neq i} A_{ij} \langle \sigma_i, u_j \rangle - \frac{t^2}{2} \sum_{j \neq i} A_{ij} \|u_j\|^2 \langle \sigma_i, \sigma_j \rangle \right. \\ &\quad \left. + t^2 \langle u_i, v_i \rangle - \frac{t^2}{2} \|u_i\|^2 \langle \sigma_i, g_i \rangle \right) + \mathcal{O}(t^3), \\ &= f(\boldsymbol{\sigma}) + t \sum_{j=1}^n \sum_{i \neq j} A_{ji} \langle \sigma_i, u_j \rangle - \frac{t^2}{2} \sum_{j=1}^n \sum_{i \neq j} A_{ji} \|u_j\|^2 \langle \sigma_i, \sigma_j \rangle \\ &\quad + t^2 \sum_{i=1}^n \left( \langle u_i, v_i \rangle - \frac{1}{2} \|u_i\|^2 \langle \sigma_i, g_i \rangle \right) + \mathcal{O}(t^3), \end{aligned}$$

where the last line follows since  $\mathbf{A}$  is symmetric. Using the definition  $g_j = \sum_{i \neq j} A_{ji} \sigma_i$  and  $\sigma_i = g_i / \|g_i\|$  in the above inequality yields

$$\begin{aligned} f(\boldsymbol{\sigma}^k) &= f(\boldsymbol{\sigma}) + t \sum_{j=1}^n \langle g_j, u_j \rangle - \frac{t^2}{2} \sum_{j=1}^n \|u_j\|^2 \langle g_j, \sigma_j \rangle \\ &\quad + t^2 \sum_{i=1}^n \left( \langle u_i, v_i \rangle - \frac{1}{2} \|u_i\|^2 \langle \sigma_i, g_i \rangle \right) + \mathcal{O}(t^3), \\ &= f(\boldsymbol{\sigma}) + t^2 \sum_{i=1}^n \left( \langle u_i, v_i \rangle - \|u_i\|^2 \|g_i\| \right) + \mathcal{O}(t^3). \end{aligned} \quad (22)$$

Reorganizing terms, we get

$$f(\bar{\sigma}) - f(\sigma^k) = f(\sigma) - f(\sigma^k) = t^2 \sum_{i=1}^n \left( \|u_i\|^2 \|g_i\| - \langle u_i, v_i \rangle \right) + \mathcal{O}(t^3). \tag{23}$$

Turning back our attention to (21), we can lower bound the right-hand side as follows

$$\begin{aligned} \|\text{grad} f(\sigma^k)\|_{\mathbb{F}}^2 &\geq 2t^2 \sum_{i=1}^n \frac{1}{\|u_i\|^2} \left( \|u_i\|^2 \|g_i\| - \langle u_i, v_i \rangle \right)^2 + \mathcal{O}(t^3), \\ &\geq 2t^2 \sum_{i=1}^n \left( \|u_i\|^2 \|g_i\| - \langle u_i, v_i \rangle \right)^2 + \mathcal{O}(t^3), \\ &\geq \frac{2t^2}{n} \left( \sum_{i=1}^n \left( \|u_i\|^2 \|g_i\| - \langle u_i, v_i \rangle \right) \right)^2 + \mathcal{O}(t^3), \end{aligned}$$

where the second inequality follows since  $\|u_i\|^2 \leq \|u\|_{\mathbb{F}}^2 = 1$  and the last inequality follows since  $(\sum_{i=1}^n a_i)^2 \leq n \sum_{i=1}^n a_i^2$ , for all  $a_i \in \mathbb{R}, i \in [n]$ . Using the second order approximation derived in (23) in the above inequality, we obtain

$$\begin{aligned} \|\text{grad} f(\sigma^k)\|_{\mathbb{F}}^2 &\geq \frac{f(\bar{\sigma}) - f(\sigma^k)}{n} \sum_{i=1}^n 2 \left( \|u_i\|^2 \|g_i\| - \langle u_i, v_i \rangle \right) + \mathcal{O}(t^3), \\ &= \frac{2\langle u, (\mathbf{A} - \mathbf{A})u \rangle}{n} \left( f(\bar{\sigma}) - f(\sigma^k) \right) + \mathcal{O}(t^3), \end{aligned}$$

where  $\mathbf{A} = \text{Diag}(\|g_1\|, \dots, \|g_n\|)$ . Since we have  $2\langle u, (\mathbf{A} - \mathbf{A})u \rangle \leq -\mu \|u\|_{\mathbb{F}}^2$  by the quadratic decay condition, we conclude that

$$\|\text{grad} f(\sigma^k)\|_{\mathbb{F}}^2 \geq \frac{\mu}{n} \left( f(\bar{\sigma}) - f(\sigma^k) \right) + \mathcal{O}(t^3). \tag{24}$$

This implies that whenever  $\sigma^k$  is sufficiently close to  $\sigma$ , i.e., whenever  $t$  is sufficiently small (cf. (19)), the remainder in the Taylor approximation, i.e., the  $\mathcal{O}(t^3)$  terms, will be dominated by  $\frac{\mu}{n} (f(\bar{\sigma}) - f(\sigma^k))$ . In particular, if  $\sigma^0$  is sufficiently close to  $\bar{\sigma}$  to satisfy  $\mathcal{O}(t^3) \geq -\frac{\mu}{2n} (f(\bar{\sigma}) - f(\sigma^k))$  in the above inequality (see Appendix B for a proof of this), we then have

$$\|\text{grad} f(\sigma^k)\|_{\mathbb{F}}^2 \geq \frac{\mu}{2n} \left( f(\bar{\sigma}) - f(\sigma^k) \right). \tag{25}$$

Combining this inequality with (18), we get

$$f(\sigma^{k+1}) - f(\sigma^k) \geq \frac{\mu}{4n^2 \|\mathbf{A}\|_1} \left( f(\bar{\sigma}) - f(\sigma^k) \right). \tag{26}$$

Rearranging terms in the above inequality concludes the proof.  $\square$

The linear convergence rate of the BCM algorithm with greedy coordinate selection in Theorem 2 can be extended for importance sampling and uniform sampling as we highlight in the following (its proof follows similar lines to the proofs of Theorem 2 and Corollary 1 and hence is omitted).

**Corollary 2** *Let the conditions in Theorem 2 hold. Then, the iterates generated by the BCM algorithm enjoys the local linear convergence rate*

$$f(\bar{\sigma}) - \mathbb{E}f(\sigma^k) \leq (1 - \rho)^k \left( f(\bar{\sigma}) - f(\sigma^0) \right), \quad (27)$$

where  $\rho = \frac{\mu}{4n\|A\|_{1,1}}$  for importance sampling and  $\rho = \frac{\mu}{4n^2\|A\|_1}$  for uniform sampling.

### 3.4 Quadratic decay condition holds generically

In this section, we consider the quadratic decay condition, which is a condition on (Non-CVX), and relate it to a condition on the original problem in (CVX). In particular, we characterize sufficient conditions on (CVX) for quadratic decay to hold. We first provide some background on semidefinite programming (see for example [3], for a more detailed treatment of this topic). Consider the SDP in (CVX):

$$\begin{aligned} & \text{maximize } \langle A, X \rangle \\ & \text{subject to } X_{ii} = 1, \text{ for } i \in [n], \\ & \quad X \succeq 0, \end{aligned}$$

and its dual:

$$\begin{aligned} & \text{minimize } \langle 1, y \rangle \\ & \text{subject to } Z = \text{Diag}(y) - A, \\ & \quad Z \succeq 0, \end{aligned}$$

where  $\mathbf{1}$  is the vector of ones of appropriate size. Let  $X^*$  and  $(y^*, Z^*)$  denote the primal and dual optimal solutions, respectively, and let  $r^*$  denote the rank of  $X^*$ . Then, there exists a  $Q \in O(n)$  such that

$$\begin{aligned} X^* &= Q \text{Diag}(\lambda_1, \dots, \lambda_{r^*}, 0, \dots, 0) Q^\top, \\ Z^* &= Q \text{Diag}(0, \dots, 0, \omega_{r^*+1}, \dots, \omega_n) Q^\top. \end{aligned}$$

We say that *strict complementarity* holds if  $\lambda_i > 0$  for  $i = 1, \dots, r^*$  and  $\omega_j > 0$  for  $j = r^* + 1, \dots, n$ . Furthermore, let  $Q_1 \in \mathbb{R}^{n \times r^*}$  and  $Q_2 \in \mathbb{R}^{n \times (n-r^*)}$  respectively denote the first  $r^*$  columns and the last  $n - r^*$  columns of  $Q$  and let  $q_i$  denote the  $i$ th row of  $Q_1$ , i.e.,  $Q_1 = [q_1, q_2, \dots, q_{r^*}]^\top$ . Then,  $(y^*, Z^*)$  is *dual nondegenerate* if and only if  $\{q_1 q_1^\top, \dots, q_{r^*} q_{r^*}^\top\}$  spans  $\text{Sym}_{r^*}$ , i.e., the set of real symmetric  $r^* \times r^*$  matrices

[3, Theorem 3]. Strict complementarity and dual nondegeneracy are known to hold generically (over the set of possible cost matrices  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , i.e., they fail to hold only on a subset of measure zero of  $\mathbb{R}^{n \times n}$ ) as proven in [3, Lemma 2]. Using these definitions, we show in the next theorem that strict complementarity and dual nondegeneracy are sufficient for quadratic decay to hold at the maximizer of (Non-CVX).

**Theorem 3** *Suppose that  $\mathbf{X}^* = \boldsymbol{\sigma}\boldsymbol{\sigma}^\top$  and  $(\mathbf{y}^*, \mathbf{Z}^*) = (\text{diag}(\mathbf{A}), \mathbf{A} - \mathbf{A})$  are respectively primal and dual optimal solutions satisfying strict complementarity and dual nondegeneracy, where  $\mathbf{A} = \text{Diag}(\|g_1\|, \dots, \|g_n\|)$ . If  $r \geq \text{rank}(\mathbf{X}^*)$ , then quadratic decay is satisfied for  $f$  at all  $\bar{\boldsymbol{\sigma}}$  such that  $\bar{\boldsymbol{\sigma}}\bar{\boldsymbol{\sigma}}^\top = \mathbf{X}^*$ .*

**Proof** Suppose  $\text{rank}(\mathbf{X}^*) = r^* \leq r$ , then by strict complementarity, we have  $\text{rank}(\mathbf{Z}^*) = n - r^*$  and kernel of  $\mathbf{Z}^*$  is equal to the column space of  $\mathbf{X}^*$ , i.e.,  $\ker(\mathbf{Z}^*) = \text{col}(\mathbf{X}^*)$ . Since  $\mathbf{X}^* = \boldsymbol{\sigma}\boldsymbol{\sigma}^\top$  and  $\mathbf{Z}^* = \mathbf{A} - \mathbf{A}$ , we equivalently have  $\ker(\mathbf{A} - \mathbf{A}) = \text{col}(\boldsymbol{\sigma})$ . As  $\mathbf{Z}^*$  is feasible for the dual, then  $\mathbf{Z}^* = \mathbf{A} - \mathbf{A} \geq 0$ , and consequently  $\langle \mathbf{u}, (\mathbf{A} - \mathbf{A})\mathbf{u} \rangle \geq 0$ , for all  $\mathbf{u} \in \mathbb{R}^{n \times r}$ .

Now consider the quadratic form  $h(\mathbf{u}) := \langle \mathbf{u}, (\mathbf{A} - \mathbf{A})\mathbf{u} \rangle$  over  $\mathbf{u} \in T_\sigma \mathcal{M}_r$ . First, we show that  $h(\mathbf{u}) = 0$  if and only if  $\mathbf{u} \in \mathcal{V}_\sigma$ . The *if* direction of the proof is straightforward, i.e.,  $(\mathbf{A} - \mathbf{A})\boldsymbol{\sigma} = 0$  and  $\mathbf{u} = \boldsymbol{\sigma}\mathbf{B}$  for some skew-symmetric matrix  $\mathbf{B}$  directly imply  $h(\mathbf{u}) = 0$  for all  $\mathbf{u} \in \mathcal{V}_\sigma$ . To show the *only if* direction, let  $\mathbf{u} \in T_\sigma \mathcal{M}_r$  such that  $h(\mathbf{u}) = 0$ , or equivalently  $\text{tr}((\mathbf{A} - \mathbf{A})\mathbf{u}\mathbf{u}^\top) = 0$ . As both  $\mathbf{A} - \mathbf{A}$  and  $\mathbf{u}\mathbf{u}^\top$  are positive semidefinite matrices, this implies  $(\mathbf{A} - \mathbf{A})\mathbf{u} = 0$ . Therefore, columns of  $\mathbf{u}$  are in  $\ker(\mathbf{A} - \mathbf{A}) = \text{col}(\boldsymbol{\sigma})$ , which implies there exists  $\mathbf{B} \in \mathbb{R}^{r \times r}$  such that  $\mathbf{u} = \boldsymbol{\sigma}\mathbf{B}$  (note that it is not possible to make this claim without strict complementarity). As  $\mathbf{u} \in T_\sigma \mathcal{M}_r$ , then  $\langle \sigma_i, u_i \rangle = \langle \sigma_i, \mathbf{B}^\top \sigma_i \rangle = \langle \sigma_i \sigma_i^\top, \mathbf{B} \rangle = 0$ , for all  $i \in [n]$ . Without loss of generality, assume that the last  $r - r^*$  columns of  $\boldsymbol{\sigma}$  are equal to zero. Then, by dual nondegeneracy of the SDP, the principal submatrices of dimension  $r^* \times r^*$  of  $\{\sigma_i \sigma_i^\top\}_{i=1}^n$  spans  $\mathcal{S}^{r^*}$ . Consider the decomposition

$$\mathbf{B} = \begin{bmatrix} \mathbf{B}_{11} & \mathbf{B}_{12} \\ \mathbf{B}_{21} & \mathbf{B}_{22} \end{bmatrix},$$

where  $\mathbf{B}_{11} \in \mathbb{R}^{r^* \times r^*}$  and  $\mathbf{B}_{22} \in \mathbb{R}^{(r-r^*) \times (r-r^*)}$ . Then, the dual nondegeneracy implies that  $\mathbf{B}_{11}$  is a skew-symmetric matrix, i.e.,  $\mathbf{B}_{11}^\top = -\mathbf{B}_{11}$ . Furthermore, as the last  $r - r^*$  columns of  $\boldsymbol{\sigma}$  are equal to zero, then  $\mathbf{u} = \boldsymbol{\sigma}\mathbf{B}$  does not depend on  $\mathbf{B}_{21}$  and  $\mathbf{B}_{22}$ . Therefore, we can pick  $\mathbf{B}_{21} = -\mathbf{B}_{12}^\top$  and  $\mathbf{B}_{22} = 0$  such that  $\mathbf{B}$  is a skew-symmetric matrix and observe that  $\mathbf{u} \in \mathcal{V}_\sigma$ . The same argument can be extended for all  $\bar{\boldsymbol{\sigma}}$  such that  $\bar{\boldsymbol{\sigma}}\bar{\boldsymbol{\sigma}}^\top = \mathbf{X}^*$  using parallel transport.

To conclude the proof, we let  $\{\mathbf{u}^\ell\}_{\ell=1}^{n(r-1)}$  be an orthogonal basis to  $T_\sigma \mathcal{M}_r$  such that  $\{\mathbf{u}^\ell\}_{\ell=1}^s$  is a basis for  $\mathcal{V}_\sigma$ . Let  $\mathbf{M} \in \mathbb{R}^{n(r-1) \times n(r-1)}$  such that  $M_{ij} = \langle \mathbf{u}^i, (\mathbf{A} - \mathbf{A})\mathbf{u}^j \rangle$ . Consider the function  $\bar{h} : \mathbb{R}^{n(r-1)} \rightarrow \mathbb{R}^{n(r-1)}$  such that  $\bar{h}(\mathbf{v}) = \mathbf{v}^\top \mathbf{M} \mathbf{v}$  and observe that  $\bar{h}(\text{vec}(\mathbf{u})) = h(\mathbf{u})$ . Let  $\mathbf{L} = [\text{vec}(\mathbf{u}^1), \dots, \text{vec}(\mathbf{u}^s)]^\top \in \mathbb{R}^{s \times n(r-1)}$ , then  $\mathbf{v}^\top \mathbf{M} \mathbf{v} > 0$  for all  $\mathbf{v}$  such that  $\mathbf{L} \mathbf{v} = 0$  and  $\mathbf{v} \neq \mathbf{0}$ . Then, by Finsler's Lemma,  $\mathbf{L}_\perp^\top \mathbf{M} \mathbf{L}_\perp > 0$ , where  $\mathbf{L}_\perp$  is any basis of the right null-space of  $\mathbf{L}$ . Equivalently, there exists  $\mu > 0$  such that  $h(\mathbf{u}) \geq \mu \|\mathbf{u}\|_{\mathbb{F}}^2$  for all  $\mathbf{u} \in T_\sigma \mathcal{M}_r \setminus \mathcal{V}_\sigma$ .  $\square$

**Remark 1** Finsler's Lemma [13, Lemma C.11.2] also yields that  $\mu = \lambda_{\min}(\mathbf{L}_\perp^\top \mathbf{M} \mathbf{L}_\perp)$ .

This theorem states that quadratic decay holds for all global maxima of (Non-CVX) provided that the rank of the factorization is large enough so that the global maximum values of (CVX) and (Non-CVX) are equal to one another. For this case, the set of all global maxima is an equivalence class corresponding to a solution since strict complementarity and dual nondegeneracy imply that the primal solution of (CVX) is unique. On top of this, when  $r \geq \sqrt{2n}$ , it is known that (see [11, Theorem 2]) any local maximum is global generically (i.e., for almost all cost matrices  $\mathbf{A}$ ). As strict complementarity and dual nondegeneracy also hold generically for (CVX), then consequently, when  $r \geq \sqrt{2n}$ , quadratic decay holds for all local maxima generically as we highlight in the following corollary.

**Corollary 3** *If  $r \geq \sqrt{2n}$ , then quadratic decay holds for all local maxima generically.*

#### 4 Approximately achieving the maximum value of (CVX)

Our results in Sect. 3 show that the BCM algorithm converges with a sublinear rate to a first-order stationary solution and with a linear rate to a local maximum when initialized sufficiently close to it. In this section, we incorporate a second-order oracle to BCM in order to obtain an algorithm, which we refer to as BCM2, that returns an approximate second-order stationary point. More specifically, at the current iteration of the algorithm, if the norm of the gradient is large, we take a BCM step. Otherwise, we run a subroutine (e.g., Lanczos method) to find the leading eigenvector of the Hessian. The main motivation for designing such an algorithm is that the approximate second-order stationary solutions provide  $\mathcal{O}(1/r)$  approximation to (CVX). In particular, call  $\sigma$  an  $\varepsilon$ -approximate concave point if  $\langle \mathbf{u}, \text{Hess} f(\sigma)[\mathbf{u}] \rangle \leq \varepsilon \langle \mathbf{u}, \mathbf{u} \rangle$ , for all  $\mathbf{u} \in T_\sigma \mathcal{M}_r$ . Then, the following theorem provides an approximation ratio between the approximate concave points of (Non-CVX) and the maximum value of (CVX).

**Theorem 4** [32, Theorem 1] *Let  $\sigma \in \mathcal{M}_r$  be an  $\varepsilon$ -approximate concave point. Then, for any positive semidefinite  $\mathbf{A}$ , the following approximation ratio holds:*

$$f(\sigma) \geq \left(1 - \frac{1}{r-1}\right) \text{SDP}(\mathbf{A}) - \frac{n}{2}\varepsilon, \quad (28)$$

where  $\text{SDP}(\mathbf{A})$  is the maximum value of (CVX).

This approximation ratio follows due to a generalization of the randomized rounding approach (most famously presented by [22]) applied to an  $\varepsilon$ -approximate concave point. In fact, it can be shown that it is not possible to find a better approximation ratio (in terms of the dependence on the rank of the factorization  $r$ ) for all problems  $\mathbf{A}$ . This result is highlighted in the following theorem.

**Theorem 5** [14, Theorems 1 & 3] *Let  $\text{SDP}(\mathbf{A})$  be the maximum value of (CVX) and  $\text{SDP}_r(\mathbf{A})$  be the maximum value of (Non-CVX). Then, for all positive semidefinite matrices  $\mathbf{A}$ , the following approximation ratio holds:*

$$1 \geq \frac{\text{SDP}_r(\mathbf{A})}{\text{SDP}(\mathbf{A})} \geq \gamma(r) = \frac{2}{r} \left( \frac{\Gamma((r+1)/2)}{\Gamma(r/2)} \right)^2 = 1 - \Theta(1/r), \quad (29)$$

**Algorithm 2** BCM2

- 1: Initialize  $\sigma^0 \in \mathbb{R}^{n \times r}$  and calculate  $g_i^0 = \sum_{j \neq i} A_{ij} \sigma_j^0$ , for all  $i \in [n]$ .
- 2: **for**  $k = 0, 1, 2, \dots$  **do**
- 3:   Compute  $\|\text{grad} f(\sigma^k)\|_F^2 = 2 \sum_{i=1}^n (\|g_i^k\|^2 - \langle \sigma_i^k, g_i^k \rangle^2)$ .
- 4:   **if**  $\|\text{grad} f(\sigma^k)\|_F^2 > \varepsilon^3 / (1350 \|A\|_1)$  **then**
- 5:      $i_k \leftarrow \arg \max_{i \in [n]} (\|g_i^k\| - \langle \sigma_i^k, g_i^k \rangle)$
- 6:      $\sigma_{i_k}^{k+1} \leftarrow g_{i_k}^k / \|g_{i_k}^k\|$ .
- 7:      $g_i^{k+1} \leftarrow g_i^k - A_{i i_k} \sigma_{i_k}^k + A_{i i_k} \sigma_{i_k}^{k+1}$ , for all  $i \neq i_k$ .
- 8:   **else**
- 9:     Find a direction  $\mathbf{u}^k \in T_{\sigma^k} \mathcal{M}_r$  such that  $\langle \mathbf{u}^k, \text{Hess} f(\sigma^k)[\mathbf{u}^k] \rangle \geq \lambda_{\max}(\text{Hess} f(\sigma^k))/2$ ,  $\langle \mathbf{u}^k, \text{grad} f(\sigma^k) \rangle \geq 0$ , and  $\|\mathbf{u}^k\|_F = 1$ .
- 10:      $\sigma_i^{k+1} \leftarrow \sigma_i^k \cos(\|u_i^k\|t) + \frac{u_i^k}{\|u_i^k\|} \sin(\|u_i^k\|t)$ , for all  $i \in [n]$ , where  $t = \varepsilon / (15 \|A\|_1)$ .
- 11:      $g_i^{k+1} \leftarrow \sum_{j \neq i} A_{ij} \sigma_j^{k+1}$ , for all  $i \in [n]$ .
- 12:   **end if**
- 13: **end for**

where  $\Gamma(z) = \int_0^\infty x^{z-1} e^{-x} dx$  is the Gamma function. Furthermore, under the unique games conjecture, there is no polynomial-time algorithm that approximates  $\text{SDP}_r(A)$  with an approximation ratio greater than  $\gamma(r) + \varepsilon$  for any  $\varepsilon > 0$ .

These results provide motivation to design algorithms with second-order guarantees to solve (Non-CVX) and for this reason, we propose the BCM2 algorithm (see Algorithm 2), which can be described as follows: When the Frobenius norm of the Riemannian gradient is at least as large as  $\|\text{grad} f(\sigma^k)\|_F^2 > \varepsilon^3 / (1350 \|A\|_1)$ , we use BCM to update the current solution. Otherwise, we assume that there is a second-order oracle that returns an update direction  $\mathbf{u}^k \in T_{\sigma^k} \mathcal{M}_r$  such that  $\langle \mathbf{u}^k, \text{Hess} f(\sigma^k)[\mathbf{u}^k] \rangle \geq \lambda_{\max}(\text{Hess} f(\sigma^k))/2$ ,  $\langle \mathbf{u}^k, \text{grad} f(\sigma^k) \rangle \geq 0$ , and  $\|\mathbf{u}^k\|_F = 1$ . Notice that finding a tangent vector  $\mathbf{u}^k$  that satisfy  $\langle \mathbf{u}^k, \text{Hess} f(\sigma^k)[\mathbf{u}^k] \rangle \geq \lambda_{\max}(\text{Hess} f(\sigma^k))/2$  and  $\|\mathbf{u}^k\|_F = 1$  is an eigenpair problem and can be solved efficiently using the Lanczos method. The condition  $\langle \mathbf{u}^k, \text{grad} f(\sigma^k) \rangle \geq 0$ , on the other hand, can always be satisfied by switching the sign of  $\mathbf{u}^k$ . It is a straightforward exercise to explicitly construct such a vector and it can be found in [9, Lemma 11]. Once the update direction  $\mathbf{u}^k \in T_{\sigma^k} \mathcal{M}_r$  is obtained, we take a step towards this direction using the geodesics on the manifold. When the step size is carefully chosen, it can be shown that the objective value of the iterates generated by this procedure is a monotonically increasing sequence until the approximate second-order stationary condition is satisfied. This property is presented in the following lemma.

**Lemma 2** *Let  $\mathbf{u}^k \in T_{\sigma^k} \mathcal{M}_r$  such that  $\|\mathbf{u}^k\|_F = 1$ ,  $\langle \mathbf{u}^k, \text{grad} f(\sigma^k) \rangle \geq 0$ , and  $\langle \mathbf{u}^k, \text{Hess} f(\sigma^k)[\mathbf{u}^k] \rangle \geq \varepsilon/2$ . Consider the update rule given by the exponential map  $\sigma^{k+1} = \text{Exp}_{\sigma^k}(t\mathbf{u}^k)$ , i.e.,*

$$\sigma_i^{k+1} = \sigma_i^k \cos(\|u_i^k\|t) + \frac{u_i^k}{\|u_i^k\|} \sin(\|u_i^k\|t), \quad \text{for all } i \in [n], \tag{30}$$

where  $t = \frac{\varepsilon}{15\|A\|_1}$  is the step size. These iterates satisfy the following ascent in the function value:

$$f(\sigma^{k+1}) - f(\sigma^k) \geq \frac{\varepsilon^3}{2700\|A\|_1^2}.$$

**Proof** The Taylor expansion of  $\sigma^{k+1}$  around  $\sigma^k$  is given by

$$\begin{aligned} \sigma_i^{k+1} &= \sigma_i^k \sum_{\ell=0}^{\infty} \frac{(-1)^\ell}{(2\ell)!} (\|u_i^k\|t)^{2\ell} + u_i^k \sum_{\ell=0}^{\infty} \frac{(-1)^\ell}{(2\ell+1)!} (\|u_i^k\|t)^{2\ell+1}, \\ &= \sigma_i^k + tu_i^k - \frac{t^2}{2} \|u_i^k\|^2 \sigma_i^k - \frac{t^3}{6} \|u_i^k\|^2 u_i^k + \dots, \end{aligned}$$

and using this, we can compute the Taylor expansion of  $f(\sigma^{k+1})$  as follows

$$\begin{aligned} f(\sigma^{k+1}) &= \sum_{i=1}^n \sum_{j \neq i} A_{ij} \langle \sigma_i^{k+1}, \sigma_j^{k+1} \rangle, \\ &= \sum_{i=1}^n \sum_{j \neq i} A_{ij} \left[ \langle \sigma_i^k, \sigma_j^k \rangle + t \left( \langle \sigma_i^k, u_j^k \rangle + \langle u_i^k, \sigma_j^k \rangle \right) \right. \\ &\quad \left. + \frac{t^2}{2} \left( -\|u_j^k\|^2 \langle \sigma_i^k, \sigma_j^k \rangle + 2\langle u_i^k, u_j^k \rangle - \|u_i^k\|^2 \langle \sigma_i^k, \sigma_j^k \rangle \right) \right] - t^3 \beta, \end{aligned}$$

where  $\beta$  represents the third and higher-order terms. Using the definitions of  $f(\sigma^k)$  and its derivatives, the above equality can be written as follows

$$f(\sigma^{k+1}) = f(\sigma^k) + t \langle u^k, \text{grad } f(\sigma^k) \rangle + \frac{t^2}{2} \langle u^k, \text{Hess } f(\sigma^k)[u^k] \rangle - t^3 \beta. \tag{31}$$

Here, our aim is to upper bound the magnitude of the remainder term corresponding to the third and higher-order terms. To this end, we upper bound the higher-order terms using the Cauchy-Schwarz inequality for each term individually. This yields

$$|\beta| \leq \sum_{i=1}^n \sum_{j \neq i} |A_{ij}| \left( \sum_{\ell=3}^{\infty} \frac{t^{\ell-3}}{\ell!} (\|u_i^k\| + \|u_j^k\|)^\ell \right).$$

As  $t < 1$  and  $A$  is a symmetric matrix, we can upper bound the right hand-side of the above inequality as follows

$$|\beta| \leq \|A\|_1 \sum_{i=1}^n \left( \sum_{\ell=3}^{\infty} \frac{2^\ell}{\ell!} \|u_i^k\|^\ell \right).$$

Since  $\|u_i\| \leq 1$  for all  $i \in [n]$ , we consequently have

$$|\beta| \leq \|A\|_1 \left( \sum_{i=1}^n \|u_i^k\|^2 \right) \left( \sum_{\ell=3}^{\infty} \frac{2^\ell}{\ell!} \right) = \|A\|_1 \sum_{\ell=3}^{\infty} \frac{2^\ell}{\ell!}.$$

where the latter equality follows since  $\|u^k\|_F = 1$ . Using  $\sum_{\ell=3}^{\infty} \frac{2^\ell}{\ell!} = e^2 - 5 \leq 5/2$  above and plugging this bound back in (31), we obtain

$$f(\sigma^{k+1}) \geq f(\sigma^k) + t \langle u^k, \text{grad} f(\sigma^k) \rangle + \frac{t^2}{2} \langle u^k, \text{Hess} f(\sigma^k)[u^k] \rangle - \frac{5\|A\|_1}{2} t^3. \tag{32}$$

Since we are given that  $\langle u^k, \text{grad} f(\sigma^k) \rangle \geq 0$  and  $\langle u^k, \text{Hess} f(\sigma^k)[u^k] \rangle \geq \varepsilon/2$ , (32) yields

$$f(\sigma^{k+1}) - f(\sigma^k) \geq \frac{\varepsilon}{4} t^2 - \frac{5\|A\|_1}{2} t^3.$$

Choosing  $t = \frac{\varepsilon}{15\|A\|_1}$  maximizes the right-hand side of the above inequality and concludes the proof.  $\square$

Using this ascent lemma, we next analyze the global convergence of Algorithm 2 in Theorem 6, where we assume that we have access to a subroutine that solves the eigenpair problem to the desired accuracy. We then implement the subroutine using the Lanczos algorithm (presented in Algorithm 3) and present its convergence in Theorem 7. In particular, we have the following theorem for the former case.

**Theorem 6** Consider Algorithm 2, where BCM is used at iteration  $k$  if  $\|\text{grad} f(\sigma^k)\|_F^2 \geq \varepsilon^3/(1350\|A\|_1)$  and a second-order step (see lines 9-11 of Algorithm 2) is taken otherwise. Let  $K_{\text{BCM}}$  denote the number of BCM epochs made and let  $K_{\text{H}}$  denote the number of second-order oracle iterations made such that  $K = nK_{\text{BCM}} + K_{\text{H}}$ . Then, as soon as

$$K_{\text{BCM}} + K_{\text{H}} = \left\lceil \frac{675n\|A\|_1^2}{\varepsilon^2} \right\rceil, \tag{33}$$

Algorithm 2 is guaranteed to return a solution  $\sigma^K$  that satisfies

$$f(\sigma^K) \geq \left(1 - \frac{1}{r-1}\right) \text{SDP}(A) - \frac{n}{2}\varepsilon, \tag{34}$$

where  $\text{SDP}(A)$  is the maximum value of (CVX).

**Proof** As we have proven previously in (12), each iteration of BCM yields the following functional ascent

$$f(\boldsymbol{\sigma}^{k+1}) - f(\boldsymbol{\sigma}^k) \geq \frac{\|\text{grad} f(\boldsymbol{\sigma}^k)\|_F^2}{2n\|\mathbf{A}\|_1} \geq \frac{\varepsilon^3}{2700n\|\mathbf{A}\|_1^2}, \quad (35)$$

where the latter inequality holds since BCM is applied at iteration  $k$  of Algorithm 2 if  $\|\text{grad} f(\boldsymbol{\sigma}^k)\|_F^2 \geq \frac{\varepsilon^3}{1350\|\mathbf{A}\|_1}$ . Similarly, by Lemma 2, each iteration of the second-order oracle yields the following functional ascent

$$f(\boldsymbol{\sigma}^{k+1}) - f(\boldsymbol{\sigma}^k) \geq \frac{\varepsilon^3}{2700\|\mathbf{A}\|_1^2}. \quad (36)$$

Hence, an epoch ( $n$  iterations) of BCM yields the same amount of function value improvement as an iteration of the second-order oracle. Let

$$f^* = \left(1 - \frac{1}{r-1}\right) \text{SDP}(\mathbf{A})$$

denote the desired approximation ratio and consider the approximation gap of the solution  $\boldsymbol{\sigma}$  with respect to  $f^*$  that is given by

$$h(\boldsymbol{\sigma}) = f^* - f(\boldsymbol{\sigma}). \quad (37)$$

The aim of the algorithm is to find a solution  $\boldsymbol{\sigma}$  that satisfy  $h(\boldsymbol{\sigma}) \leq \epsilon$  for some  $\epsilon > 0$ . Consider that the BCM2 algorithm runs  $K_{\text{BCM}}$  epochs of BCM and  $K_{\text{H}}$  iterations of the second-order oracle such that a total of  $K = nK_{\text{BCM}} + K_{\text{H}}$  iterations are made. Let  $\mathcal{G} = \{0 \leq k \leq K-1 : \|\text{grad} f(\boldsymbol{\sigma}^k)\|_F^2 \geq \frac{\varepsilon^3}{1350\|\mathbf{A}\|_1}\}$  be the set of iterations at which BCM step is taken and let  $\mathcal{H} = \{0 \leq k \leq K-1\} \setminus \mathcal{G}$  be the set of iterations at which a second-order oracle step is taken. Then, the approximation gap decreases at each iteration by the following amount:

$$h(\boldsymbol{\sigma}^k) - h(\boldsymbol{\sigma}^{k+1}) \geq \frac{\varepsilon^3}{2700\|\mathbf{A}\|_1^2} \delta_k, \quad (38)$$

where, for notational simplicity, we introduced

$$\delta_k = \begin{cases} \frac{1}{n}, & \text{if } k \in \mathcal{G}, \\ 1, & \text{if } k \in \mathcal{H}. \end{cases} \quad (39)$$

By Theorem 4, we are given that any  $\varepsilon$ -approximate concave point  $\boldsymbol{\sigma}$  satisfies

$$h(\boldsymbol{\sigma}) \leq \frac{n}{2}\varepsilon. \quad (40)$$

Hence, the right-hand side of (38) can be lower bounded as follows

$$h(\sigma^k) - h(\sigma^{k+1}) \geq \frac{2\delta_k}{675n^3\|\mathbf{A}\|_1^2} h^3(\sigma^k). \tag{41}$$

Considering the reciprocal of the approximation gap, we observe that

$$\begin{aligned} \frac{1}{h^2(\sigma^{k+1})} - \frac{1}{h^2(\sigma^k)} &= \frac{(h(\sigma^k) - h(\sigma^{k+1}))(h(\sigma^k) + h(\sigma^{k+1}))}{h^2(\sigma^{k+1})h^2(\sigma^k)}, \\ &\geq \frac{2\delta_k}{675n^3\|\mathbf{A}\|_1^2} \frac{h(\sigma^k)(h(\sigma^k) + h(\sigma^{k+1}))}{h^2(\sigma^{k+1})}, \end{aligned} \tag{42}$$

where the inequality follows by (41). As the right-hand side of (41) is lower bounded by zero, we have  $h(\sigma^k) \geq h(\sigma^{k+1})$ . Thus, we can lower bound the right-hand side of (42) as follows

$$\frac{1}{h^2(\sigma^{k+1})} - \frac{1}{h^2(\sigma^k)} \geq \frac{4\delta_k}{675n^3\|\mathbf{A}\|_1^2}. \tag{43}$$

Summing (43) over  $k = 0, 1, \dots, K - 1$ , we get

$$\frac{1}{h^2(\sigma^K)} - \frac{1}{h^2(\sigma^0)} \geq \sum_{k=0}^{K-1} \frac{4\delta_k}{675n^3\|\mathbf{A}\|_1^2} = \frac{4}{675n^3\|\mathbf{A}\|_1^2} (K_{\text{BCM}} + K_{\text{H}}).$$

Given that  $\sigma^0$  is not an  $\varepsilon$ -approximate concave point (or else, there is nothing to prove), we have

$$\frac{1}{h^2(\sigma^K)} \geq \frac{4}{675n^3\|\mathbf{A}\|_1^2} (K_{\text{BCM}} + K_{\text{H}}). \tag{44}$$

Since by (40), we know that  $\frac{1}{h(\sigma)} \geq \frac{2}{n\varepsilon}$  for any  $\varepsilon$ -approximate concave point, then as soon as

$$K_{\text{BCM}} + K_{\text{H}} \geq \frac{675n\|\mathbf{A}\|_1^2}{\varepsilon^2} \tag{45}$$

iterations made, BCM2 is guaranteed to return an  $\varepsilon$ -approximate concave point, i.e., there exists a solution  $\sigma^k$  for some  $1 < k < K$  such that  $h(\sigma^k) \leq \frac{n}{2}\varepsilon$ . Since  $\{h(\sigma^k)\}_{k \geq 0}$  is a nonincreasing sequence (as we have already shown in (41)), then the final iterate of the algorithm  $\sigma^K$  is guaranteed to satisfy  $h(\sigma^K) \leq \frac{n}{2}\varepsilon$ , i.e.,  $\sigma^K$  is an  $\varepsilon$ -approximate concave point.  $\square$

In Theorem 6,  $K_{\text{BCM}} + K_{\text{H}}$  represents the total number of epochs to guarantee (34), whereas the iteration counter of the algorithm is given in terms of  $K = nK_{\text{BCM}} + K_{\text{H}}$ . This is due to the fact that, at each iteration of the BCM algorithm, a single row

of  $\sigma$  is updated and consequently  $n$  iterations of the BCM algorithm add up to an epoch. On the other hand, at each iteration of the second-order step, all entries of  $\sigma$  are updated, and hence each second-order iteration is an epoch. In terms of the computational cost, an iteration of BCM requires  $\mathcal{O}(nr)$  operations and consequently an epoch of BCM requires  $\mathcal{O}(n^2r)$  operations, whereas the second-order direction of update is typically found approximately via a few iterations of the power method or the Lanczos method (see Theorem 7 for a more rigorous treatment of this statement), which require  $\mathcal{O}(n^2r)$  operations. Therefore, an epoch of Algorithm 2 typically has a computational complexity of  $\mathcal{O}(n^2r)$ . Furthermore, by Theorem 6, we observe that in at most  $\mathcal{O}(n\|A\|_1^2/\varepsilon^2)$  epochs, Algorithm 2 returns a solution that achieves the optimal approximation ratio up to an accuracy of  $\mathcal{O}(n\varepsilon)$ . In particular, picking  $\varepsilon = 2 \text{SDP}(A)/(n(r-1))$ , we obtain the following corollary.

**Corollary 4** *Consider the setup of Theorem 6 and set  $\varepsilon = 2 \text{SDP}(A)/(n(r-1))$ . Then, as soon as*

$$K = \left\lceil \frac{675n^3(r-1)^2\|A\|_1^2}{4(\text{SDP}(A))^2} \right\rceil, \quad (46)$$

Algorithm 2 is guaranteed to return a solution  $\sigma^K$  that satisfies

$$f(\sigma^K) \geq \left(1 - \frac{2}{r-1}\right) \text{SDP}(A).$$

**Remark 2** In order to understand the total running time of BCM2, consider the following example. Let  $A$  be the adjacency matrix of a random Erdos-Rényi graph on  $n$  nodes and  $\lfloor cn \rfloor$  edges. The size of the maximum cut in this graph normalized by the number of nodes can be bounded between  $[c/2 + 0.4\sqrt{c}, c/2 + 0.6\sqrt{c}]$  with high probability as  $n$  increases, for all sufficiently large  $c$  [20]. Since the maximum value of (CVX) is within 0.878 of the maximum cut [22], we can then conclude that  $\text{SDP}(A)/n = \mathcal{O}(c)$  with high probability. We can also observe that for this graph, the degree of a node approximately follows a Poisson distribution with mean  $2c$ , which can be approximated by a normal distribution with mean  $2c$  and variance  $\sqrt{2c}$ , for large  $c$  [20]. Then, we have  $\|A\|_1 = \mathcal{O}(c \log n)$  with high probability. Therefore, for this problem, Corollary 4 states that in  $\tilde{\mathcal{O}}(nr^2)$  iterations (where tilde is used to hide the logarithmic dependences), Algorithm 2 returns a  $\mathcal{O}(1/r)$ -optimal solution with high probability. Per iteration computational cost of the algorithm is  $\mathcal{O}(nrc)$ , which results in a total running time of  $\tilde{\mathcal{O}}(n^2r^3c)$ . In comparison, Klein-Lu method (see [27, Lemma 4]) requires  $\tilde{\mathcal{O}}(n^2r^3c)$  running time and the matrix multiplicative weights method (see [5, Theorem 3]) requires  $\tilde{\mathcal{O}}(n^2r^{3.5}/c)$  running time to return a  $1/r$ -optimal solution.

**Remark 3** It has been shown in [37, Theorem 3.1] and [17, Theorem 3.5] that an exactly feasible approximately second-order stationary point to (Non-CVX) is also approximately optimal for (CVX). Our BCM2 algorithm returns such a solution and in light of these results, we can conclude that it finds a high-quality solution to (CVX)

with high probability whenever  $r \geq \sqrt{2n}$ . See Fig. 4 for an empirical validation of this result.

In the description of Algorithm 2 (see line 9), we assumed that we have access to a vector in the tangent space of the current iterate, which satisfies certain second-order conditions. In Algorithm 3, we describe an efficient subroutine to find this desired tangent vector based on the Lanczos method. In particular, the Lanczos method returns a tridiagonal real symmetric matrix whose diagonal entries are  $\{\alpha_\ell\}_{\ell \geq 1}$  and off-diagonal entries are  $\{\beta_\ell\}_{\ell \geq 2}$ , where  $\ell$  denotes the iteration counter in Algorithm 3. The entire spectrum of such a symmetric tridiagonal matrix can be efficiently computed in almost linear time in the dimension of the matrix [18]. Consequently, letting  $y$  denote the leading eigenvector of this tridiagonal matrix, we can construct the desired tangent vector in Algorithm 2 as  $\mathbf{u}^k = \sum_{\ell \geq 1} y_\ell \mathbf{u}_\ell$ . It is well-known that after  $n(r - 1)$  iterations, the Lanczos method constructs the leading eigenvector exactly (since order- $n(r - 1)$  Krylov subspace spans the entire tangent space). Furthermore, it is also possible to analyze the performance of the Lanczos method with early termination [28]. Building on these ideas, we characterize the quality of the solution returned by Algorithms 2+3 in the following theorem, whose proof can be found in ‘‘Appendix C’’.

**Theorem 7** *Suppose in Algorithm 3, we initialize  $\mathbf{u}_1$  uniformly at random over  $T_\sigma \mathcal{M}_r$ . Let*

$$\ell^* = \left\lceil \left( \frac{1}{2} + 2\sqrt{\frac{\|\mathbf{A}\|_1}{\varepsilon}} \right) \log \left( \frac{\left\lceil \frac{675n\|\mathbf{A}\|_1^2}{\varepsilon^2} \right\rceil 1.648\sqrt{n(r-1)}}{\delta} \right) \right\rceil,$$

and consider that Algorithm 3 is run for  $\min(\ell^*, n(r - 1))$  iterations at each call from Algorithm 2. Then, after  $K$  iterations (defined as in (46)), Algorithm 2 returns a solution  $\sigma^K$  that satisfies

$$f(\sigma^K) \geq \left(1 - \frac{1}{r-1}\right) \text{SDP}(\mathbf{A}) - \frac{n}{2}\varepsilon,$$

with probability at least  $1 - \delta$ .

---

**Algorithm 3** Lanczos Method

- 1: Given  $\sigma$ , define  $H[\mathbf{u}] = \text{Hess}f(\sigma)[\mathbf{u}] + 4\|\mathbf{A}\|_1\mathbf{u}$ . Initialize  $\mathbf{u}_1 \in T_\sigma \mathcal{M}_r$  such that  $\|\mathbf{u}_1\|_F = 1$ . Let  $\alpha_1 = \langle \mathbf{u}_1, H[\mathbf{u}_1] \rangle$  and  $\mathbf{r}_1 = H[\mathbf{u}_1] - \alpha_1\mathbf{u}_1$ .
  - 2: **for**  $\ell \geq 2$  **do**
  - 3:    $\beta_\ell = \|\mathbf{r}_{\ell-1}\|_F$
  - 4:    $\mathbf{u}_\ell = \mathbf{r}_{\ell-1}/\beta_\ell$  (If  $\beta_\ell = 0$ , pick  $\mathbf{u}_\ell \perp \text{span}(\mathbf{u}_1, \dots, \mathbf{u}_{\ell-1})$  arbitrarily)
  - 5:    $\alpha_\ell = \langle \mathbf{u}_\ell, H[\mathbf{u}_\ell] \rangle$
  - 6:    $\mathbf{r}_\ell = H[\mathbf{u}_\ell] - \alpha_\ell\mathbf{u}_\ell - \beta_\ell\mathbf{u}_{\ell-1}$
  - 7: **end for**
-

## 5 Numerical experiments

In this section, we evaluate the empirical performance of the BCM algorithm. In what follows,  $n$  is the dimension of the cost matrix  $A \in \mathbb{R}^{n \times n}$ , and  $r$  refers to the rank of factorization. All algorithms are implemented on Matlab and the experiments are run on a computer with 2.9 GHz processor and 16 GB memory. RGD and RTR algorithms are implemented using the Manopt package [10] with the default options and the algorithms are terminated when the maximum allowed time is reached.

In all experiments, the cost matrix is generated as  $A = (G + G^\top)/n$ , where  $G_{ij} \sim N(0, 1)$  for all  $i \neq j$ , and  $G_{ii} = 0$  for all  $i \in [n]$ . All experiments are based on 50 Monte Carlo runs over the initialization. For each run, the initial iterate  $\sigma^0 \in \mathbb{R}^{n \times r}$  is the same for all algorithms and each row of  $\sigma^0$  is generated uniformly at random on  $\mathbb{S}^{r-1}$ .

First, we compare various coordinate selection schemes for BCM (see Algorithm 1). We compare cyclic order  $i = (1, 2, \dots, n)$ , uniform random selection, random permutation order ( $i$  follows a cyclic order of a uniformly random permutation of  $[n]$ ), greedy coordinate selection, and selection by importance sampling. Figure 1 summarizes the results of our experiments on  $n \in \{200, 1000\}$  with  $r = \lceil \sqrt{2n} \rceil$ . We observe that greedy coordinate selection achieves higher function value after running each algorithm the same number of iterations; yet, due to its high per-iteration cost, cyclic, uniformly random, and random permutation selection rules perform better in terms of overall runtime complexity. Furthermore, randomized rules that do not cycle through all coordinates achieve lower function values after running each algorithm the same number of iterations. This phenomenon is observed for a number of numerical examples in different papers and unfortunately we do not have a good theoretical understanding about this behavior except for a few preliminary results [23,24,29]. It would be an interesting future direction to theoretically understand the slower convergence of randomized coordinate selection rules in practice.

Next, we evaluate the performance of these algorithms for  $n \in \{200, 1000, 5000\}$  with  $r = 2$  and  $r = \sqrt{2n}$ . In Fig. 2, empirical results illustrate the fast convergence of BCM and BCM2 (see Algorithm 2) compared to RGD and RTR, for both  $r = \lceil \sqrt{2n} \rceil$  and  $r = 2$ . The numerical results indicate our algorithms return a high-quality solution much faster than RGD and RTR regardless of the rank of the factorization is larger or smaller than the Barvinok-Pataki bound.

We next compare the final performance of different methods after convergence. That is, we run all algorithms sufficiently enough until their function value stabilize, and compare the final value obtained. In Fig. 3, we clearly observe that the final function values obtained through BCM2 (Algorithm 2 with Lanczos method) and RTR are larger than those obtained by other algorithms. We also observe that even when the problem size is large (e.g.,  $n = 20,000$ ), BCM returns a desirable solution within  $\sim 20$  seconds, whereas it takes approximately a minute for RTR to return such a solution.

Finally, we consider a random SDP with a planted solution. In particular, we consider a matrix  $X \geq 0$  such that  $\text{rank}(X) = r$  and  $X \in \mathbb{S}^n$  where  $n = \frac{r(r+1)}{2}$ . We then generate a MaxCut SDP for which  $X$  is an optimal solution, i.e., we find a cost matrix  $A$  in the normal cone of  $X$  (this requires solving an auxiliary SDP). For each  $r \in \{4, 7, 10\}$ , we generate 100 random MaxCut SDPs as described above. We perform

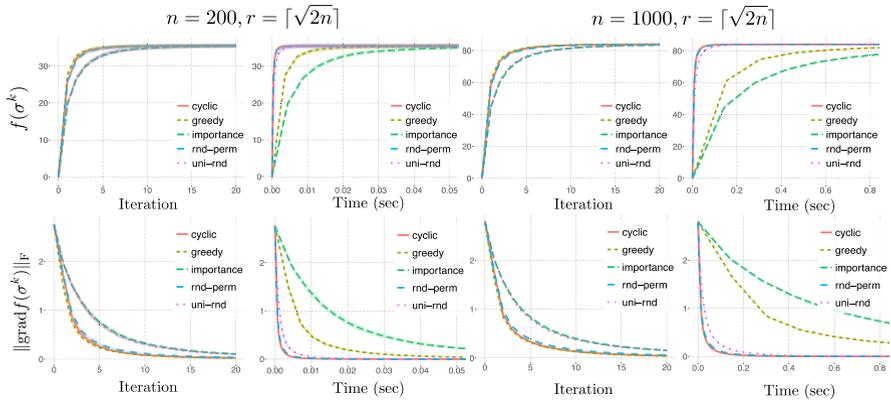


Fig. 1 Comparisons of different randomization schemes for  $n \in \{200, 1000\}$  with  $r = \lceil \sqrt{2n} \rceil$

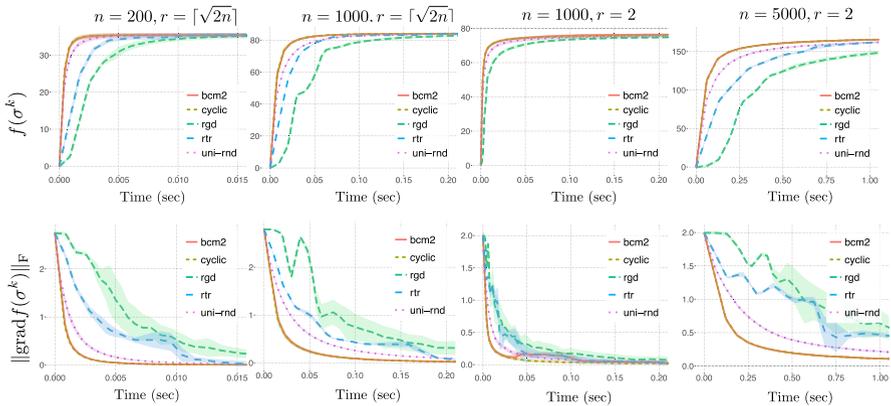


Fig. 2 Performance of BCM and BCM2 (Algorithms 1 and 2) compared to other methods. Here, RTR and RGD refer to Riemannian Trust Region and Riemannian Gradient Descent, respectively

a Burer–Monteiro factorization to these SDPs for a range of ranks in  $[r - 4, r + 4]$ . We solve the resulting non-convex problem using our BCM2 algorithm. Figure 4 shows the percentage of experiments solved correctly for each value of  $r$ . We consider a trial correct if the solution returned by BCM2 is sufficiently close to the maximizer of the SDP. Figure 4 shows that there is a sharp phase transition at the Barvinok–Pataki bound. Above this bound, the solutions returned by our BCM2 algorithm is approximately optimal to (CVX) with high probability.

## 6 Conclusion

In this paper, we studied the Burer–Monteiro approach to solve large-scale SDPs. We considered to solve this non-convex problem using the block-coordinate maximization

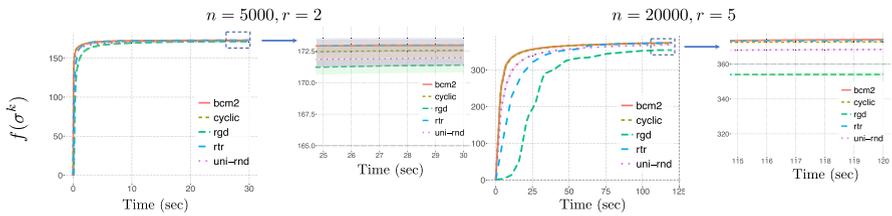


Fig. 3 Comparing the final performance of different methods after (near) convergence

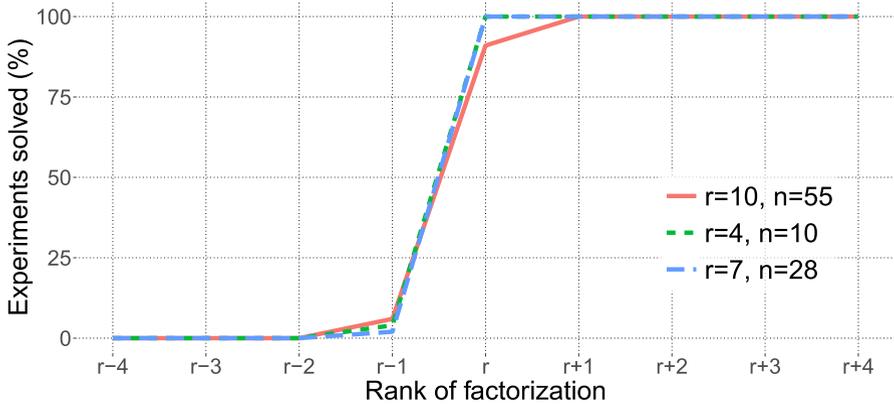


Fig. 4 Phase transition in recovering the optimal solution of (CVX) by an approximately second-order stationary solution of (Non-CVX)

algorithm that is extremely simple to implement. We proved that for various coordinate selection rules, BCM attains a global sublinear convergence rate of  $\mathcal{O}(1/\epsilon)$  to guarantee  $\mathbb{E}\|\text{grad } f(\sigma^k)\|_F^2 \leq \epsilon$ . We also showed the linear convergence of BCM around a local maximum that satisfy the quadratic decay condition. We proved that the quadratic decay condition generically holds for all local maxima provided that  $r \geq \sqrt{2n}$ . These are the first precise rate estimates for the non-convex Burer–Monteiro approach in the literature to the best of our knowledge. We then introduced a new algorithm called BCM2 based on BCM and Lanczos methods. We showed that BCM2 is guaranteed to return a solution that provides  $1 - \mathcal{O}(1/r)$  approximation to the SDP without any assumptions on the cost matrix  $A$ , where the  $r$ -dependence of this approximation is optimal under the unique games conjecture. We also presented numerical results that verify our theoretical findings and show that BCM is faster than the state-of-the-art methods. Even though in this paper, we only considered SDPs with diagonal constraints, it would be of interest to study the block-coordinate maximization approach in more generic problems.

### A Proof of Corollary 1

Similar to the proof of Theorem 1, from Proposition 1, we have

$$\begin{aligned}
 f(\boldsymbol{\sigma}^{k+1}) - f(\boldsymbol{\sigma}^k) &= 2 \left( \|g_{i_k}^k\| - \langle \sigma_{i_k}^k, g_{i_k}^k \rangle \right), \\
 &= \frac{2 \|g_{i_k}^k\| \left( \|g_{i_k}^k\| - \langle \sigma_{i_k}^k, g_{i_k}^k \rangle \right)}{\|g_{i_k}^k\|}, \\
 &\geq \frac{\|g_{i_k}^k\|^2 - \langle \sigma_{i_k}^k, g_{i_k}^k \rangle^2}{\|g_{i_k}^k\|}, \tag{47}
 \end{aligned}$$

where the inequality follows since  $\|g_{i_k}^k\| \geq \langle \sigma_{i_k}^k, g_{i_k}^k \rangle$ , for all  $\sigma_{i_k}^k \in \mathbb{R}^{n \times r}$ . Letting  $\mathbb{E}_k$  denote the expectation over  $i_k$  given  $\boldsymbol{\sigma}^k$ , we have

$$\mathbb{E}_k f(\boldsymbol{\sigma}^{k+1}) - f(\boldsymbol{\sigma}^k) \geq \sum_{i=1}^n p_i \frac{\|g_i^k\|^2 - \langle \sigma_i^k, g_i^k \rangle^2}{\|g_i^k\|}.$$

In particular, when  $p_i = \frac{1}{n}$ , for all  $i \in [n]$  (i.e., for uniform sampling case), we have

$$\mathbb{E}_k f(\boldsymbol{\sigma}^{k+1}) - f(\boldsymbol{\sigma}^k) \geq \frac{1}{n \|A\|_1} \sum_{i=1}^n \left( \|g_i^k\|^2 - \langle \sigma_i^k, g_i^k \rangle^2 \right),$$

since  $\|g_i^k\| \leq \|A\|_1$ , for all  $i \in [n]$  by (11). Therefore, we have

$$\mathbb{E}_k f(\boldsymbol{\sigma}^{k+1}) - f(\boldsymbol{\sigma}^k) \geq \frac{\|\text{grad} f(\boldsymbol{\sigma}^k)\|_{\mathbb{F}}^2}{2n \|A\|_1}. \tag{48}$$

On the other hand, when  $p_i = \frac{\|g_i^k\|}{\sum_{j=1}^n \|g_j^k\|}$  (i.e., for importance sampling case), we have

$$\mathbb{E}_k f(\boldsymbol{\sigma}^{k+1}) - f(\boldsymbol{\sigma}^k) \geq \frac{\sum_{i=1}^n \|g_i^k\|^2 - \langle \sigma_i^k, g_i^k \rangle^2}{\sum_{j=1}^n \|g_j^k\|} = \frac{\|\text{grad} f(\boldsymbol{\sigma}^k)\|_{\mathbb{F}}^2}{2 \sum_{j=1}^n \|g_j^k\|}.$$

Letting  $\|A\|_{1,1} = \sum_{i,j=1}^n |A_{ij}|$  denote the  $L_{1,1}$  norm of matrix  $A$ , we observe that  $\sum_{j=1}^n \|g_j^k\| \leq \|A\|_{1,1}$ , which in the above inequality yields

$$\mathbb{E}_k f(\boldsymbol{\sigma}^{k+1}) - f(\boldsymbol{\sigma}^k) \geq \frac{\|\text{grad} f(\boldsymbol{\sigma}^k)\|_{\mathbb{F}}^2}{2 \|A\|_{1,1}}. \tag{49}$$

In order to prove (13), which corresponds to uniform sampling case, we assume the contrary that  $\mathbb{E} \|\text{grad} f(\boldsymbol{\sigma}^k)\|_{\mathbb{F}}^2 > \epsilon$ , for all  $k \in [K - 1]$ . Then, using the boundedness

of  $f$ , we get

$$\begin{aligned} f^* - f(\sigma^0) &\geq \mathbb{E}f(\sigma^K) - f(\sigma^0) = \sum_{k=0}^{K-1} \mathbb{E} \left[ f(\sigma^{k+1}) - f(\sigma^k) \right] \\ &= \sum_{k=0}^{K-1} \mathbb{E} \left[ \mathbb{E}_k f(\sigma^{k+1}) - f(\sigma^k) \right]. \end{aligned}$$

Using the expected functional ascent of BCM in (48) above, we get

$$f^* - f(\sigma^0) \geq \sum_{k=0}^{K-1} \frac{\mathbb{E} \|\text{grad} f(\sigma^k)\|_F^2}{2n \|A\|_1} > \frac{K\epsilon}{2n \|A\|_1}, \tag{50}$$

where the last inequality follows by the assumption. Then, by contradiction, the algorithm returns a solution with  $\mathbb{E} \|\text{grad} f(\sigma^k)\|_F^2 \leq \epsilon$ , for some  $k \in [K - 1]$ , provided that

$$K \geq \frac{2n \|A\|_1 (f^* - f(\sigma^0))}{\epsilon}.$$

The proof of (14), which corresponds to importance sampling case, can be obtained by using (49) (instead of (48)) in (50), and hence is omitted.

### B Rest of the Proof of Theorem 2

In order to quantify how close  $\sigma^0$  and  $\sigma$  should be so that this convergence rate holds, we need to derive explicit bounds on the higher order terms in (21) and (23), which we do in the following. The Taylor expansion of  $\sigma^k$  around  $\sigma$  yields

$$\begin{aligned} \sigma_i^k &= \sigma_i \cos(\|u_i\|t) + \frac{u_i}{\|u_i\|} \sin(\|u_i\|t), \\ &= \sigma_i \left[ \sum_{\ell=0}^{\infty} \frac{(-1)^\ell}{(2\ell)!} (\|u_i\|t)^{2\ell} \right] + \frac{u_i}{\|u_i\|} \left[ \sum_{\ell=0}^{\infty} \frac{(-1)^\ell}{(2\ell + 1)!} (\|u_i\|t)^{2\ell+1} \right]. \end{aligned}$$

Using this expansion, we can compute  $f(\sigma^k) = \sum_{i,j=1}^n A_{ij} \langle \sigma_i^k, \sigma_j^k \rangle$ . The first three terms in the expansion are already given in (22) as follows

$$f(\sigma^k) = f(\sigma) + t^2 \sum_{i=1}^n \left( \langle u_i, v_i \rangle - \|u_i\|^2 \|g_i\| \right) + \beta_f, \tag{51}$$

where  $\beta_f$  represents the higher order terms. In order to find an upper bound on  $|\beta_f|$ , we use the Cauchy-Schwarz inequality in the higher order terms in the expansion of

$f(\sigma^k)$ , which yields the following bound

$$|\beta_f| \leq \sum_{i,j=1}^n |A_{ij}| \left( \sum_{\ell=3}^{\infty} \frac{t^\ell}{\ell!} (\|u_i\| + \|u_j\|)^\ell \right).$$

As  $\|u\|_F = 1$ , we have  $\|u_i\| \leq 1$  for all  $i \in [n]$ , which implies

$$|\beta_f| \leq \sum_{i,j=1}^n |A_{ij}| \left( \sum_{\ell=3}^{\infty} \frac{t^\ell}{\ell!} 2^\ell \right),$$

where we note that  $t$  denotes the geodesic distance between  $\sigma^k$  and  $[\bar{\sigma}]$  as highlighted in (19). Assuming that  $t \leq 1$ , we obtain the following upper bound

$$|\beta_f| \leq t^3 n \|A\|_1 \left( \sum_{\ell=3}^{\infty} \frac{2^\ell}{\ell!} \right).$$

Using the inequality  $\sum_{\ell=3}^{\infty} \frac{2^\ell}{\ell!} = e^2 - 5 \leq 5/2$  above, we get

$$|\beta_f| \leq \frac{5n \|A\|_1 t^3}{2}.$$

Plugging this value back in (51), we obtain

$$f(\sigma^k) \leq f(\sigma) + t^2 \sum_{i=1}^n \left( \langle u_i, v_i \rangle - \|u_i\|^2 \|g_i\| \right) + \frac{5n \|A\|_1 t^3}{2}. \tag{52}$$

Considering the same expansion for  $\|\text{grad} f(\sigma^k)\|_F^2 = 2 \sum_{i=1}^n (\|g_i^k\|^2 - \langle \sigma_i^k, g_i^k \rangle^2)$ , we get the following (see (21)):

$$\|\text{grad} f(\sigma^k)\|_F^2 = 2t^2 \sum_{i=1}^n \left( \|u_i\| \|g_i\| - \left\langle \frac{u_i}{\|u_i\|}, v_i \right\rangle \right)^2 + \beta_g, \tag{53}$$

where  $\beta_g$  represents the higher order terms. Upper bounding each higher order terms using the Cauchy-Schwarz inequality as follows, we obtain

$$|\beta_g| \leq 2 \sum_{i=1}^n \left[ \sum_{j,m=1}^n |A_{ij}| |A_{im}| \left( \sum_{\ell=3}^{\infty} \frac{t^\ell}{\ell!} (\|u_j\| + \|u_m\|)^\ell \right) + \sum_{j,m=1}^n |A_{ij}| |A_{im}| \left( \sum_{\substack{\ell,s=0 \\ \ell+s \geq 3}}^{\infty} \frac{t^{\ell+s}}{\ell!s!} (\|u_i\| + \|u_j\|)^{\ell+s} \right) \right].$$

Using the fact that  $\|u_i\| \leq 1$  for all  $i \in [n]$ , we get the following upper bound

$$|\beta_g| \leq 2 \sum_{i=1}^n \left[ \sum_{j,m=1}^n |A_{ij}| |A_{im}| \left( \sum_{\ell=3}^{\infty} \frac{t^\ell}{\ell!} 2^\ell \right) + \sum_{j,m=1}^n |A_{ij}| |A_{im}| \left( \sum_{\substack{\ell,s=0 \\ \ell+s \geq 3}}^{\infty} \frac{t^{\ell+s}}{\ell!s!} 2^{\ell+s} \right) \right].$$

Using the upper bound  $\sum_{j,m=1}^n |A_{ij}| |A_{im}| \leq \|A\|_1^2$  above, we obtain

$$|\beta_g| \leq 2 \|A\|_1^2 \sum_{i=1}^n \left[ \sum_{\ell=3}^{\infty} \frac{t^\ell}{\ell!} 2^\ell + \sum_{\substack{\ell,s=0 \\ \ell+s \geq 3}}^{\infty} \frac{t^{\ell+s}}{\ell!s!} 2^{\ell+s} \right].$$

Introducing a change of variables in the last sum, we get

$$\begin{aligned} |\beta_g| &\leq 2 \|A\|_1^2 \sum_{i=1}^n \left[ \sum_{\ell=3}^{\infty} \frac{t^\ell}{\ell!} 2^\ell + \sum_{\ell=3}^{\infty} \frac{t^\ell}{\ell!} 2^\ell \left( \sum_{s=0}^{\ell} \frac{\ell!}{s!(\ell-s)!} \right) \right], \\ &= 2 \|A\|_1^2 \sum_{i=1}^n \left[ \sum_{\ell=3}^{\infty} \frac{t^\ell}{\ell!} (2^\ell + 4^\ell) \right]. \end{aligned}$$

Assuming that  $t \leq 1$ , we obtain the following upper bound

$$|\beta_g| \leq 2 \|A\|_1^2 t^3 \sum_{i=1}^n \left[ \sum_{\ell=3}^{\infty} \frac{1}{\ell!} (2^\ell + 4^\ell) \right].$$

Using the inequality  $\sum_{\ell=3}^{\infty} \frac{2^\ell + 4^\ell}{\ell!} = e^2 + e^4 - 18 \leq 44$  above, we get

$$|\beta_g| \leq 88n \|A\|_1^2 t^3.$$

Plugging this value back in (53), we obtain

$$\|\text{grad} f(\sigma^k)\|_{\mathbb{F}}^2 \geq 2t^2 \sum_{i=1}^n \left( \|u_i\| \|g_i\| - \left\langle \frac{u_i}{\|u_i\|}, v_i \right\rangle \right)^2 - 88n \|A\|_1^2 t^3. \tag{54}$$

Using the same bounding technique as in (24), we get

$$\begin{aligned} \|\text{grad} f(\sigma^k)\|_{\mathbb{F}}^2 &\geq \frac{\mu}{n} \left( f(\bar{\sigma}) - f(\sigma^k) - \frac{5n \|A\|_1 t^3}{2} \right) - 88n \|A\|_1^2 t^3, \\ &= \frac{\mu}{n} \left( f(\bar{\sigma}) - f(\sigma^k) \right) - t^3 \|A\|_1 (3\mu + 88n \|A\|_1). \end{aligned}$$

Therefore, in order for (25) to hold, we need

$$t^3 \|A\|_1 (3\mu + 88n \|A\|_1) \leq \frac{\mu}{2n} (f(\bar{\sigma}) - f(\sigma^k)),$$

which can be equivalently rewritten as follows

$$t^3 \leq \frac{\mu(f(\bar{\sigma}) - f(\sigma^k))}{2n \|A\|_1 (3\mu + 88n \|A\|_1)}.$$

As  $f(\sigma^k)$  is a monotonically non-decreasing sequence, then as soon as  $\sigma^0$  is sufficiently close to  $[\bar{\sigma}]$  in the sense that

$$\text{dist}(\sigma^0, [\bar{\sigma}]) \leq \left( \frac{\mu(f(\bar{\sigma}) - f(\sigma^k))}{2n \|A\|_1 (3\mu + 88n \|A\|_1)} \right)^{1/3},$$

then the linear convergence rate presented in (26) holds.

### C Proof of Theorem 7

Before presenting the proof of Theorem 7, we first introduce the following theorem that characterizes the convergence rate of the Lanczos method with random initialization.

**Theorem 8** [28, Theorem 4.2] *Let  $A \in \mathbb{R}^{n \times n}$  be a positive semidefinite matrix,  $b \in \mathbb{R}^n$  be an arbitrary vector and  $\lambda_L^\ell(A, b)$  denote the output of the Lanczos algorithm after  $\ell$  iterations when applied to find the leading eigenvalue of  $A$  (denoted by  $\lambda_1(A)$ ) with initialization  $b$ . In particular,*

$$\lambda_L^\ell(A, b) = \max \left\{ \frac{\langle x, Ax \rangle}{\langle x, x \rangle} : 0 \neq x \in \text{span}(b, \dots, A^{\ell-1}b) \right\}.$$

*Assume that  $b$  is uniformly distributed over the set  $\{b \in \mathbb{R}^n : \|b\| = 1\}$  and let  $\epsilon \in [0, 1)$ . Then, the probability that the Lanczos algorithm does not return an  $\epsilon$ -approximation to the leading eigenvalue of  $A$  exponentially decreases as follows*

$$\mathbb{P} \left( \lambda_L^\ell(A, b) < (1 - \epsilon)\lambda_1(A) \right) \begin{cases} \leq 1.648\sqrt{n}e^{-\sqrt{\epsilon}(2\ell-1)}, & \text{if } 0 < \ell < n(r-1), \\ = 0, & \text{if } \ell \geq n(r-1). \end{cases}$$

Using this result, Theorem 7 is proven as follows. Since the tangent space  $T_\sigma \mathcal{M}_r$  has dimension  $n(r-1)$ , then we can define a symmetric matrix (where we drop the notational dependency on  $\sigma$  for simplicity)  $H \in \mathbb{R}^{n(r-1) \times n(r-1)}$  that represents the linear operator  $\text{Hess} f(\sigma)$  in the basis  $\{u^1, \dots, u^{n(r-1)}\}$  such that  $\text{span}(u^1, \dots, u^{n(r-1)}) = T_\sigma \mathcal{M}_r$ . In particular, letting  $H_{ij} = \langle u^i, \text{Hess} f(\sigma)[u^j] \rangle$  yields the desired matrix  $H$  and the Lanczos algorithm is run to find the leading eigenvalue of this matrix. Here, it is important to note that  $H$  is not a psd matrix, so it is required to shift  $H$  with a

large enough multiple of the identity matrix so that the resulting matrix is guaranteed to be positive semidefinite. In particular, by inspecting the definition of  $\text{Hess } f(\boldsymbol{\sigma})$  in (5), it is easy to observe that  $\|\text{Hess } f(\boldsymbol{\sigma})\|_{\text{op}} \leq 4\|\mathbf{A}\|_1$ . Therefore, it is sufficient to run the Lanczos algorithm to find the leading eigenvalue of  $\tilde{\mathbf{H}} = \mathbf{H} + 4\|\mathbf{A}\|_1 \mathbf{I}$ , where  $\mathbf{I}$  denotes the appropriate sized identity matrix. On the other hand, we initialize the Lanczos algorithm with a random vector  $\mathbf{u}$  of unit norm (i.e.,  $\|\mathbf{u}\|_F = 1$ ) in the tangent space  $T_{\boldsymbol{\sigma}} \mathcal{M}_r$ . Notice that  $\mathbf{u}$  can equivalently be represented as a vector  $b \in \mathbb{R}^{n(r-1)}$  in the basis  $\{\mathbf{u}^1, \dots, \mathbf{u}^{n(r-1)}\}$  as  $\mathbf{u} = \sum_{i=1}^{n(r-1)} b_i \mathbf{u}^i$  such that  $\|b\| = 1$ . Then, by Theorem 8, we have

$$\mathbb{P}\left(\lambda_L^\ell(\tilde{\mathbf{H}}, b) < (1 - \epsilon)\lambda_1(\tilde{\mathbf{H}})\right) \leq 1.648\sqrt{n(r-1)}e^{-\sqrt{\epsilon}(2\ell-1)}.$$

Letting  $\lambda_1(\mathbf{H})$  denote the leading eigenvalue of  $\mathbf{H}$ , we run the Lanczos algorithm to obtain a vector  $b^*$  such that  $\|b^*\| = 1$  and  $\langle b^*, \mathbf{H}b^* \rangle \geq \lambda_1(\mathbf{H})/2$ . Thus, we want  $\mathbb{P}\left(\lambda_L^\ell(\tilde{\mathbf{H}}, b) < 4\|\mathbf{A}\|_1 + \lambda_1(\mathbf{H})/2\right)$  to be small. Setting  $\epsilon^* = \frac{\lambda_1(\mathbf{H})}{16\|\mathbf{A}\|_1}$ , we can observe that

$$\begin{aligned} (1 - \epsilon^*)\lambda_1(\tilde{\mathbf{H}}) &= \left(1 - \frac{\lambda_1(\mathbf{H})}{16\|\mathbf{A}\|_1}\right)(4\|\mathbf{A}\|_1 + \lambda_1(\mathbf{H})), \\ &= 4\|\mathbf{A}\|_1 + \frac{3\lambda_1(\mathbf{H})}{4} - \frac{(\lambda_1(\mathbf{H}))^2}{16\|\mathbf{A}\|_1}, \\ &\geq 4\|\mathbf{A}\|_1 + \frac{\lambda_1(\mathbf{H})}{2}, \end{aligned}$$

where the inequality follows since  $\lambda_1(\mathbf{H}) \leq 4\|\mathbf{A}\|_1$ . Consequently, we have

$$\begin{aligned} &\mathbb{P}\left(\lambda_L^\ell(\tilde{\mathbf{H}}, b) < 4\|\mathbf{A}\|_1 + \lambda_1(\mathbf{H})/2\right) \\ &\leq \mathbb{P}\left(\lambda_L^\ell(\tilde{\mathbf{H}}, b) < (1 - \epsilon^*)\lambda_1(\tilde{\mathbf{H}})\right) \leq 1.648\sqrt{n(r-1)}e^{-\sqrt{\epsilon^*}(2\ell-1)}. \end{aligned}$$

By Theorem 6, we know that the Lanczos method is called at most  $\lceil 675n\|\mathbf{A}\|_1^2/\epsilon^2 \rceil$  times to search for an  $\epsilon$ -approximate concave point and for any non-desired solution we have  $\lambda_1(\mathbf{H}) \geq \epsilon$  by the definition of  $\epsilon$ -approximate concave point. Then, by using a union bound over all calls to the Lanczos method, we conclude that when the Lanczos method is run for  $\ell$  iterations, we have the following guarantee

$$\begin{aligned} &\mathbb{P}(\text{Algorithm 2+3 fails to return an } \epsilon\text{-approximate concave point}) \\ &\leq \left\lceil \frac{675n\|\mathbf{A}\|_1^2}{\epsilon^2} \right\rceil 1.648\sqrt{n(r-1)}e^{-\sqrt{\frac{\epsilon}{16\|\mathbf{A}\|_1}}(2\ell-1)}. \end{aligned}$$

In order to set this probability to some  $\delta \in (0, 1)$ , we let

$$\begin{aligned} \ell^* &= \left\lceil \left( \frac{1}{2} + 2\sqrt{\frac{\|A\|_1}{\varepsilon}} \right) \log \left( \frac{\left\lceil \frac{675n\|A\|_1^2}{\varepsilon^2} \right\rceil 1.648\sqrt{n(r-1)}}{\delta} \right) \right\rceil \\ &= \tilde{O} \left( \sqrt{\frac{\|A\|_1}{\varepsilon}} \log \left( \frac{n\sqrt{n(r-1)}}{\delta} \right) \right), \end{aligned}$$

where tilde is used to hide poly-logarithmic factors in  $\|A\|_1/\varepsilon$ . Since the Lanczos algorithm is guaranteed to return the leading eigenvalue with probability 1 in at most  $n(r-1)$  iterations, then running each Lanczos subroutine for  $\min(\ell^*, n(r-1))$  iterations, it is guaranteed that Algorithm 2+3 returns an  $\varepsilon$ -approximate concave point with probability at least  $1 - \delta$ .

## References

1. Absil, P.-A., Baker, C.G., Gallivan, K.A.: Trust-region methods on Riemannian manifolds. *Found. Comput. Math.* **7**(3), 303–330 (2007)
2. Absil, P.-A., Mahony, R., Sepulchre, R.: *Optimization Algorithms on Matrix Manifolds*. Princeton University Press, Princeton (2007)
3. Alizadeh, F., Haeblerly, J.-P.A., Overton, M.L.: Complementarity and nondegeneracy in semidefinite programming. *Math. Program.* **77**(1), 111–128 (1997)
4. Anitescu, M.: Degenerate nonlinear programming with a quadratic growth condition. *SIAM J. Optim.* **10**(4), 1116–1135 (2000)
5. Arora, S., Hazan, E., Kale, S.: Fast algorithms for approximate semidefinite programming using the multiplicative weights update method. In: *Proceedings of the 46th Annual IEEE Symposium on Foundations of Computer Science, FOCS'05*, pp. 339–348 (2005)
6. Bandeira, A.S., Boumal, N., Voroninski, V.: On the low-rank approach for semidefinite programs arising in synchronization and community detection. [arXiv:1602.04426](https://arxiv.org/abs/1602.04426) (2016)
7. Barvinok, A.I.: Problems of distance geometry and convex properties of quadratic maps. *Discrete Comput. Geom.* **13**(2), 189–202 (1995)
8. Bonnans, J.F., Ioffe, A.: Second-order sufficiency and quadratic growth for nonisolated minima. *Math. Oper. Res.* **20**(4), 801–817 (1995)
9. Boumal, N., Absil, P.-A., Cartis, C.: Global rates of convergence for nonconvex optimization on manifolds. [arXiv preprint arXiv:1605.08101](https://arxiv.org/abs/1605.08101) (2016)
10. Boumal, N., Mishra, B., Absil, P.-A., Sepulchre, R.: Manopt, a Matlab toolbox for optimization on manifolds. *J. Mach. Learn. Res.* **15**, 1455–1459 (2014)
11. Boumal, N., Voroninski, V., Bandeira, A.S.: The non-convex Burer–Monteiro approach works on smooth semidefinite programs. In: *Advances in Neural Information Processing Systems*, pp. 2757–2765 (2016)
12. Boumal, N., Voroninski, V., Bandeira, A.S.: Deterministic guarantees for Burer–Monteiro factorizations of smooth semidefinite programs. [arXiv preprint arXiv:1804.02008](https://arxiv.org/abs/1804.02008) (2018)
13. Briat, C.: *Linear Parameter-Varying and Time-Delay Systems*. Springer (2014)
14. Briët, J., de Oliveira Filho, F.M., Vallentin, F.: The positive semidefinite grothendieck problem with rank constraint. In: *Automata, Languages and Programming*, pp. 31–42 (2010)
15. Burer, S., Monteiro, R.D.C.: A nonlinear programming algorithm for solving semidefinite programs via low-rank factorization. *Math. Program.* **95**(2), 329–357 (2003)
16. Burer, S., Monteiro, R.D.C.: Local minima and convergence in low-rank semidefinite programming. *Math. Program.* **103**(3), 427–444 (2005)

17. Cifuentes, D., Moitra, A.: Polynomial time guarantees for the Burer–Monteiro method. arXiv preprint [arXiv:1912.01745](https://arxiv.org/abs/1912.01745) (2019)
18. Coakley, E.S., Rokhlin, V.: A fast divide-and-conquer algorithm for computing the spectra of real symmetric tridiagonal matrices. *Appl. Comput. Harmon. Anal.* **34**(3), 379–414 (2013)
19. Erdogdu, M.A., Deshpande, Y., Montanari, A.: Inference in graphical models via semidefinite programming hierarchies. In: *Advances in Neural Information Processing Systems*, pp. 416–424 (2017)
20. Gamarnik, D., Li, Q.: On the max-cut of sparse random graphs. arXiv preprint [arXiv:1411.1698](https://arxiv.org/abs/1411.1698) (2014)
21. Garber, D., Hazan, E.: Approximating semidefinite programs in sublinear time. In: *Proceedings of the 24th International Conference on Neural Information Processing Systems, NIPS'11*, pp. 1080–1088 (2011)
22. Goemans, M.X., Williamson, D.P.: Improved approximation algorithms for maximum cut and satisfiability problems using semidefinite programming. *J. ACM* **42**(6), 1115–1145 (1995)
23. Gurbuzbalaban, M., Ozdaglar, A., Parrilo, P.A., Vanli, N.D.: When cyclic coordinate descent outperforms randomized coordinate descent. In: Guyon, I., Luxburg, U.V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., Garnett, and (eds.) *Advances in Neural Information Processing Systems*, volume 30, pp. 6999–7007. Curran Associates, Inc. (2017)
24. Gurbuzbalaban, M., Ozdaglar, A., Vanli, N.D., Wright, S.J.: Randomness and permutations in coordinate descent methods. *Math. Program.* **181**, 03 (2018)
25. Javanmard, A., Montanari, A., Ricci-Tersenghi, F.: Phase transitions in semidefinite relaxations. *Proc. Natl. Acad. Sci.* **113**(16), E2218–E2223 (2016)
26. Journee, M., Bach, F., Absil, P.-A., Sepulchre, R.: Low-rank optimization on the cone of positive semidefinite matrices. *SIAM J. Optim.* **20**(5), 2327–2351 (2010)
27. Klein, P., Lu, H.-I.: Efficient approximation algorithms for semidefinite programs arising from MAX CUT and COLORING. In: *Proceedings of the Twenty-Eighth Annual ACM Symposium on Theory of Computing, STOC'96*, pp. 338–347. ACM, New York, NY, USA (1996)
28. Kuczynski, J., Woźniakowski, H.: Estimating the largest eigenvalues by the power and Lanczos algorithms with a random start. *SIAM J. Matrix Anal. Appl.* **13**(4), 1094–1122 (1992)
29. Lee, C.-P., Wright, S.J.: Random permutations fix a worst case for cyclic coordinate descent. *IMA J. Numer. Anal.* **39**, 07 (2016)
30. Lee, J.D., Simchowitz, M., Jordan, M.I., Recht, B.: Gradient descent only converges to minimizers. In: *29th Annual Conference on Learning Theory*, vol. 49, pp. 1246–1257. PMLR (2016)
31. Lu, Z., Xiao, L.: Randomized block coordinate non-monotone gradient method for a class of nonlinear programming. Technical Report MSR-TR-2013-66 (2013)
32. Mei, S., Misiakiewicz, T., Montanari, A., Oliveira, R.I.: Solving SDPs for synchronization and MaxCut problems via the Grothendieck inequality. arXiv preprint [arXiv:1703.08729](https://arxiv.org/abs/1703.08729) (2017)
33. Montanari, A.: A Grothendieck-type inequality for local maxima. arXiv preprint [arXiv:1603.04064](https://arxiv.org/abs/1603.04064) (2016)
34. Parrilo, P.A.: Semidefinite programming relaxations for semialgebraic problems. *Math. Program.* **96**(2), 293–320 (2003)
35. Pataki, G.: On the rank of extreme matrices in semidefinite programs and the multiplicity of optimal eigenvalues. *Math. Oper. Res.* **23**(2), 339–358 (1998)
36. Patrascu, A., Necoara, I.: Efficient random coordinate descent algorithms for large-scale structured nonconvex optimization. *J. Glob. Optim.* **61**, 05 (2013)
37. Pumar, T., Jelassi, S., Boumal, N.: Smoothed analysis of the low-rank approach for smooth semidefinite programs. In: *Advances in Neural Information Processing Systems*, pp. 2281–2290 (2018)
38. Richtárik, P., Takáč, M.: Iteration complexity of randomized block-coordinate descent methods for minimizing a composite function. *Math. Program.* **144**, 07 (2011)
39. Steurer, D.: Fast SDP algorithms for constraint satisfaction problems. In: *Proceedings of the Twenty-First Annual ACM–SIAM Symposium on Discrete Algorithms*, pp. 684–697 (2010)
40. Tropp, J.A., Yurtsever, A., Udell, M., Cevher, V.: Practical sketching algorithms for low-rank matrix approximation. *SIAM J. Matrix Anal. Appl.* **38**(4), 1454–1485 (2017)
41. Tseng, P., Yun, S.: A coordinate gradient descent method for nonsmooth separable minimization. *Math. Program.* **117**, 387–423 (2009)
42. Vandenberghe, L., Boyd, S.: Semidefinite programming. *SIAM Rev.* **38**(1), 49–95 (1996)
43. Wang, P.-W., Chang, W.-C., Kolter, J.Z.: The mixing method: coordinate descent for low-rank semidefinite programming. arXiv preprint [arXiv:1706.00476](https://arxiv.org/abs/1706.00476) (2017)

**Publisher's Note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.