

Sparse Sums of Positive Semidefinite Matrices

Marcel K. de Carli Silva* Nicholas J. A. Harvey† Cristiane M. Sato‡

Abstract

Recently there has been much interest in “sparsifying” sums of rank one matrices: modifying the coefficients such that only a few are nonzero, while approximately preserving the matrix that results from the sum. Results of this sort have found applications in many different areas, including sparsifying graphs. In this paper we consider the more general problem of sparsifying sums of positive semidefinite matrices that have arbitrary rank.

We give several algorithms for solving this problem. The first algorithm is based on the method of Batson, Spielman and Srivastava (2009). The second algorithm is based on the matrix multiplicative weights update method of Arora and Kale (2007). We also highlight an interesting connection between these two algorithms.

Our algorithms have numerous applications. We show how they can be used to construct graph sparsifiers with auxiliary constraints, sparsifiers of hypergraphs, and sparse solutions to semidefinite programs.

*Department of Combinatorics and Optimization, University of Waterloo. mksilva@uwaterloo.ca. Partially supported by an NSERC Discovery Grant of L. Tunçel.

†Department of Computer Science, University of British Columbia. nickhar@cs.ubc.ca. Supported by an NSERC Discovery Grant.

‡Department of Combinatorics and Optimization, University of Waterloo. cmsato@uwaterloo.ca. Partially supported by an NSERC Discovery Grant of N. Wormald.

1 Introduction

A *sparsifier* of a graph is a subgraph that approximately preserves some structural properties of the graph. The original work in this area studied *cut sparsifiers*, which are weighted subgraphs that approximate every cut arbitrarily well. The celebrated work of Benczúr and Karger [5, 6] proved that every undirected graph with n vertices and m edges (and potentially non-negative weights on its edges) has a subgraph with only $O(n \log n / \varepsilon^2)$ edges (and new weights on those edges) such that, for every cut, the weight of the cut in the original graph and its subgraph agree up to a multiplicative factor of $(1 \pm \varepsilon)$. Benczúr and Karger also gave a randomized algorithm to construct a cut sparsifier in $\tilde{O}(m / \varepsilon^2)$ time. Recent work has extended and improved their algorithm in various ways [10, 11, 12, 14, 15].

Spielman and Teng [39] introduced *spectral sparsifiers*, which are weighted subgraphs such that the quadratic forms defined by the Laplacians of the graph and the sparsifier agree up to a multiplicative factor of $(1 \pm \varepsilon)$. Spectral sparsifiers are also cut sparsifiers, as can be seen by evaluating these quadratic forms at $\{0, 1\}$ -vectors. They proved that every undirected graph with n vertices and m edges (and potentially non-negative weights on its edges) has a spectral sparsifier with only $n \text{polylog}(n) / \varepsilon^2$ edges (and new weights on those edges). Spielman and Srivastava [38] reduce the graph sparsification problem to the following abstract problem in matrix theory.

Problem 1. Let $v_1, \dots, v_m \in \mathbb{R}^n$ be vectors and let $B = \sum_i v_i v_i^T$. Given $\varepsilon \in (0, 1)$, find a vector $y \in \mathbb{R}^m$ with small support such that $y \geq 0$ and

$$B \preceq \sum_i y_i v_i v_i^T \preceq (1 + \varepsilon)B. \quad (1)$$

(Here the notation $X \preceq Y$ means that the matrix $Y - X$ is positive semidefinite.)

Spielman and Srivastava [38] observe that Problem 1 can be solved using known concentration bounds on operator-valued random variables, specifically Rudelson’s sampling lemma [32, 33]. This approach yields a vector y with support size $O(n \log n / \varepsilon^2)$, and therefore yields a construction of spectral sparsifiers with $O(n \log n / \varepsilon^2)$ edges. Their algorithm relies on the linear system solver of Spielman and Teng [39], which was significantly simplified by Koutis, Miller and Peng [24]. Recent work [23] has improved the space usage of Spielman and Srivastava’s algorithm.

In subsequent work, Batson, Spielman and Srivastava [4] give a deterministic algorithm that solves Problem 1 and produces a vector y with support size $O(n / \varepsilon^2)$. Consequently they obtain improved spectral sparsifiers with $O(n / \varepsilon^2)$ edges. This work led to important progress in metric embeddings [29, 34], convex geometry [40] and Banach space theory [37].

In this paper, we focus on a more general problem.

Problem 2. Let B_1, \dots, B_m be symmetric, positive semidefinite matrices of size $n \times n$ and let $B = \sum_i B_i$. Given $\varepsilon \in (0, 1)$, find a vector $y \in \mathbb{R}^m$ with small support such that $y \geq 0$ and

$$B \preceq \sum_i y_i B_i \preceq (1 + \varepsilon)B. \quad (2)$$

This problem can also be solved by known concentration bounds: Ahlswede and Winter [1] give a method for generalizing Chernoff-like bounds to operator-valued random variables, and one of their theorems [1, Theorem 19] directly yields a solution to Problem 2. (Other expositions of these results also exist [41, 16].) This approach yields a vector y with support size $O(n \log n / \varepsilon^2)$. See Section 3 for more details.

This paper gives two improved solutions to Problem 2. Our interest in this topic is motivated by several applications, such as constructing sparsifiers with certain auxiliary properties and sparsifiers for hypergraphs. We discuss these applications in Section 1.2.

1.1 Our Results

We give several efficient algorithms for solving Problem 2. Our strongest solution is:

Theorem 3. *Let B_1, \dots, B_m be symmetric, positive semidefinite matrices of size $n \times n$ and arbitrary rank. Set $B := \sum_i B_i$. For any $\varepsilon \in (0, 1)$, there is a deterministic algorithm to construct a vector $y \in \mathbb{R}^m$ with $O(n/\varepsilon^2)$ nonzero entries such that $y \geq 0$ and*

$$B \preceq \sum_i y_i B_i \preceq (1 + \varepsilon)B.$$

The algorithm runs in $O(mn^3/\varepsilon^2)$ time. Moreover, the result continues to hold if the input matrices B_1, \dots, B_m are Hermitian and positive semidefinite.

Our proof of Theorem 3 is quite simple and builds on results of Batson, Spielman and Srivastava [4]. We remark that the assumption that the B_i 's are positive semidefinite cannot be removed; see Appendix D.

We also give a second solution to Problem 2 which is quantitatively weaker, although it is based on very general machinery which might prove useful in further applications or generalizations of Problem 2. This second solution is based on the matrix multiplicative weights update method (MMWUM) of Arora and Kale [3, 22]. By a black-box application of their theorems we obtain a deterministic algorithm to construct a vector y with $O(n \log n/\varepsilon^3)$ nonzero entries. By slightly refining their analysis we can improve the number of nonzero entries to $O(n \log n/\varepsilon^2)$. We remark that Orecchia and Vishnoi [30] have used MMWUM for solving the balanced separator problem; this can be used as a subroutine in Spielman and Teng's algorithm for constructing spectral sparsifiers.

Another virtue of our second solution is that it illustrates that the surprising Batson-Spielman-Srivastava (BSS) algorithm is actually closely related to MMWUM. In particular, the algorithms underlying our two solutions are *identical*, except for the use of slightly different potential functions. We explain this connection in Section 8.

1.2 Applications

In this section, we present several applications of Problem 2. Proofs are given in Appendix A.

Sparsifiers with costs.

Corollary 4. *Let $G = (V, E)$ be a graph, let $w: E \rightarrow \mathbb{R}_+$ be a weight function, and let $c_1, \dots, c_k: E \rightarrow \mathbb{R}_+$ be cost functions, with $k = O(n)$. Let $\mathcal{L}_G(w)$ denote the Laplacian matrix for graph G with weight function w . For any real $\varepsilon \in (0, 1)$, there is a deterministic polynomial-time algorithm to find a subgraph H of G and a weight function $w_H: E(H) \rightarrow \mathbb{R}_+$ such that*

$$\begin{aligned} \mathcal{L}_G(w) \preceq \mathcal{L}_H(w_H) \preceq (1 + \varepsilon)\mathcal{L}_G(w), \\ \sum_{e \in E} w_e c_{i,e} \leq \sum_{e \in E(H)} w_{H,e} c_{i,e} \leq (1 + \varepsilon) \sum_{e \in E} w_e c_{i,e} \quad \text{for all } i \end{aligned}$$

and $|E(H)| = O(n/\varepsilon^2)$.

The inequalities $\mathcal{L}_G(w) \preceq \mathcal{L}_H(w_H) \preceq (1 + \varepsilon)\mathcal{L}_G(w)$ are equivalent to the condition that the subgraph H (with weights w_H) is a spectral sparsifier of G (with weights w). We remark that existing methods for producing sparsifiers have low probability of approximately satisfying even a single cost function (i.e., the case $k = 1$).

One potentially interesting application of sparsifiers with costs is as follows.

Corollary 5 (Rainbow Sparsifiers). *Let $G = (V, E)$ be a graph and let $w: E \rightarrow \mathbb{R}_+$ be a weight function. Let E_1, \dots, E_k be a partition of the edges, i.e., each edge is colored with one of k colors. For any real $\varepsilon \in (0, 1)$, there is a deterministic polynomial-time algorithm to find a subgraph H of G and a weight function $w_H: E(H) \rightarrow \mathbb{R}_+$ such that*

$$\begin{aligned} \mathcal{L}_G(w) &\preceq \mathcal{L}_H(w_H) \preceq (1 + \varepsilon)\mathcal{L}_G(w), \\ (1 - \varepsilon) \sum_{e \in E_i} w_e &\leq \sum_{e \in E(H) \cap E_i} w_{H,e} \leq (1 + \varepsilon) \sum_{e \in E_i} w_e \quad \text{for all } i, \end{aligned}$$

and $|E(H)| = O((n + k)/\varepsilon^2)$.

Hypergraph sparsifiers. Let $\mathcal{H} = (V, \mathcal{E})$ be a hypergraph, and let $w: \mathcal{E} \rightarrow \mathbb{R}_+$. We follow the definition of Laplacian for hypergraphs as in [31]. For each hyperedge $E \in \mathcal{E}$, define its Laplacian \mathcal{L}_E as the graph Laplacian of a graph on V whose edge set forms a clique on E . Define the Laplacian for the hypergraph \mathcal{H} with weight function w as the matrix $\mathcal{L}_{\mathcal{H}}(w) := \sum_{E \in \mathcal{E}} w_E \mathcal{L}_E$.

Corollary 6 (Spectral sparsifiers for hypergraphs). *For any real $\varepsilon \in (0, 1)$, there is a deterministic polynomial-time algorithm to find a sub-hypergraph \mathcal{G} of \mathcal{H} and a weight function $w_{\mathcal{G}}: \mathcal{E}(\mathcal{G}) \rightarrow \mathbb{R}_+$ such that*

$$\mathcal{L}_{\mathcal{H}}(w) \preceq \mathcal{L}_{\mathcal{G}}(w_{\mathcal{G}}) \preceq (1 + \varepsilon)\mathcal{L}_{\mathcal{H}}(w),$$

and $|\mathcal{E}(\mathcal{G})| = O(n/\varepsilon^2)$.

This corollary concerns spectral sparsifiers. It is also interesting to study sparsifiers that approximately preserve all cuts. There are several ways to extend the definition of ‘‘the weight of a cut’’ from ordinary graphs to hypergraphs. We consider the following two definitions, where S is any set of vertices in a hypergraph \mathcal{H} with edge weights w .

- $w(\delta_{\mathcal{H}}(S))$: This is the sum of the weights of all hyperedges that contain at least one vertex in S and at least one vertex in $\bar{S} := V \setminus S$.
- $w^*(\delta_{\mathcal{H}}(S))$: This is defined to be $\sum_{E \in \mathcal{E}} w_E \cdot |S \cap E| \cdot |\bar{S} \cap E|$.

Obviously these definitions agree in ordinary graphs.

Corollary 7 (Cut sparsifiers for hypergraphs, second definition). *For any real $\varepsilon \in (0, 1)$, there is a deterministic polynomial-time algorithm to find a sub-hypergraph \mathcal{G} of \mathcal{H} and a weight function $w_{\mathcal{G}}: \mathcal{E}(\mathcal{G}) \rightarrow \mathbb{R}_+$ such that*

$$w^*(\delta_{\mathcal{H}}(S)) \leq w_{\mathcal{G}}^*(\delta_{\mathcal{G}}(S)) \leq (1 + \varepsilon)w^*(\delta_{\mathcal{H}}(S)) \quad \text{for every } S \subseteq V,$$

and $|\mathcal{E}(\mathcal{G})| = O(n/\varepsilon^2)$.

Corollary 8 (Cut sparsifiers for hypergraphs, first definition). *Assume that \mathcal{H} is an r -uniform hypergraph. For any real $\varepsilon \in (0, 1)$, there is a deterministic polynomial-time algorithm to find a sub-hypergraph \mathcal{G} of \mathcal{H} and a weight function $w_{\mathcal{G}}: \mathcal{E}(\mathcal{G}) \rightarrow \mathbb{R}_+$ such that*

$$\frac{(r-1)}{r^2/4} w(\delta_{\mathcal{H}}(S)) \leq w_{\mathcal{G}}(\delta_{\mathcal{G}}(S)) \leq \frac{(1+\varepsilon)r^2}{4(r-1)} w(\delta_{\mathcal{H}}(S)) \quad \forall S \subseteq V,$$

and $|\mathcal{E}(\mathcal{G})| = O(n/\varepsilon^2)$. In other words, the sparsified hypergraph \mathcal{G} approximates the weight of the cuts in the hypergraph \mathcal{H} to within a factor $\Theta(r^2)$.

For the special case $r = 3$, we can achieve $(1 + \varepsilon)$ -approximate sparsification for all cuts, even under the first definition.

Corollary 9 (Cut sparsifiers for 3-uniform hypergraphs). *Assume that \mathcal{H} is a 3-uniform hypergraph. For any $\varepsilon \in (0, 1)$, there is a deterministic polynomial-time algorithm to find a sub-hypergraph \mathcal{G} of \mathcal{H} and a weight function $w_{\mathcal{G}}: \mathcal{E}(\mathcal{G}) \rightarrow \mathbb{R}_+$ such that*

$$w(\delta_{\mathcal{H}}(S)) \leq w_{\mathcal{G}}(\delta_{\mathcal{G}}(S)) \leq (1 + \varepsilon)w(\delta_{\mathcal{H}}(S)) \quad \forall S \subseteq V,$$

and $|\mathcal{E}(\mathcal{G})| = O(n/\varepsilon^2)$.

Sparse solutions to semidefinite programs.

Corollary 10. *Let A_1, \dots, A_m be symmetric, positive semidefinite matrices of size $n \times n$, and let B be a symmetric matrix of size $n \times n$. Let $c \in \mathbb{R}^m$ with $c \geq 0$. Suppose that the semidefinite program (SDP)*

$$\min \left\{ c^T z : \sum_i z_i A_i \succeq B, z \in \mathbb{R}^m, z \geq 0 \right\}$$

has a feasible solution z^ . Then, for any real $\varepsilon \in (0, 1)$, it has a feasible solution \bar{z} with at most $O(n/\varepsilon^2)$ nonzero entries and $c^T \bar{z} \leq (1 + \varepsilon)c^T z^*$.*

Several important SDPs can be cast as in Corollary 10; see, e.g., [19, 20]. Recently, Jain and Yao [21] gave a parallel approximation algorithm for SDPs in this form with B positive semidefinite.

Lovász theta number. For a graph $G = (V, E)$ on n nodes, let $t'(G)$ denote the square of the minimum radius of an Euclidean ball in \mathbb{R}^n such that there is a map from V to points in the ball such that adjacent vertices are mapped to points at distance at least 1. Also, let $\vartheta'(G)$ denote the variant of the Lovász theta number introduced in [27] and [35].

Corollary 11. *Let $G = (V, E)$ be a graph. For any real $\varepsilon \in (0, 1)$, there is a deterministic polynomial-time algorithm to find a subgraph H of G such that*

$$(1 - \varepsilon)t'(G) \leq t'(H) \leq t'(G)$$

and $|E(H)| = O(n/\varepsilon^2)$.

Corollary 12. *Let $G = (V, E)$ be a graph. For any real $\varepsilon \in (0, 1)$, there is a deterministic polynomial-time algorithm to find a supergraph H of G such that*

$$\frac{\vartheta'(G)}{1 - \varepsilon + \varepsilon\vartheta'(G)} \leq \vartheta'(H) \leq \vartheta'(G)$$

and $|E(H)| = \binom{n}{2} - O(n/\varepsilon^2)$.

Corollary 13. *Let G be a graph such that $\vartheta'(G) = o(\sqrt{n})$. For any real $\gamma > 0$, there is a supergraph H of G such that*

$$\frac{\vartheta'(G)}{1 + \gamma} \leq \vartheta'(H) \leq \vartheta'(G)$$

and $|E(H)| = \binom{n}{2} - O(n\vartheta'(G)^2/\gamma^2)$.

Corollary 14. *Let G be a graph such that $\vartheta'(G) = \Omega(\sqrt{n})$. For any real $\gamma \geq 1$, there is a supergraph H of G such that*

$$\vartheta'(H) = \Omega(\sqrt{n}/\gamma)$$

and $|E(H)| = \binom{n}{2} - O(n^2/\gamma^2)$.

Approximate Carathéodory theorems. One immediate application for Theorem 3 is an approximate Carathéodory-type theorem. A classic result of this sort is:

Theorem 15 (Althöfer [2], Lipton-Young [25]). *Let $v_1, \dots, v_m \in [0, 1]^n$ and let $\lambda \in \mathbb{R}^m$ satisfy $\lambda \geq 0$ and $\sum_i \lambda_i = 1$. Then there exists $\mu \in \mathbb{R}^m$ with $\mu \geq 0$, $\sum_i \mu_i = 1$ and only $O(\log n/\varepsilon^2)$ nonzero entries such that $\|\sum_i \lambda_i v_i - \sum_i \mu_i v_i\|_\infty \leq \varepsilon$.*

This theorem follows from simple random sampling arguments, but it has several interesting consequences, including the existence of sparse, low-regret solutions to zero-sum games. The following corollary of Theorem 3 can be viewed as a matrix generalization of Theorem 15.

Corollary 16. *Let B_1, \dots, B_m be symmetric, positive semidefinite matrices of size $n \times n$ and let $\lambda \in \mathbb{R}^m$ satisfy $\lambda \geq 0$ and $\sum_i \lambda_i = 1$. Let $B = \sum_i \lambda_i B_i$. For any $\varepsilon \in (0, 1)$, there exists $\mu \geq 0$ with $\sum_i \mu_i = 1$ such that μ has $O(n/\varepsilon^2)$ nonzero entries and*

$$(1 - \varepsilon)B \preceq \sum_i \mu_i B_i \preceq (1 + \varepsilon)B.$$

Although the support size in Theorem 15 is much smaller than in Corollary 16, the latter provides a multiplicative error bound whereas the former only provides an additive error bound. Theorem 15 can be modified to give multiplicative error bounds if we allow μ to have $O(n \log n/\varepsilon^2)$ non-zero entries. However such a result is not interesting as Carathéodory's theorem provides a μ with only $n + 1$ non-zero entries and no error (i.e., $\varepsilon = 0$). In contrast, Carathéodory's theorem is very weak in the scenario of Corollary 16 as it only provides a μ with $n(n + 1)/2 + 1$ nonzero entries.

Sparsifiers on subgraphs.

Corollary 17. *Let $G = (V, E)$ be a graph, let $w: E \rightarrow \mathbb{R}_+$ be a weight function, and let \mathcal{F} be a collection of subgraphs of G such that $\sum_{F \in \mathcal{F}} |V(F)| = O(n)$. For any real $\varepsilon \in (0, 1)$, there is a deterministic polynomial-time algorithm to find a subgraph H of G and a weight function $w_H: E(H) \rightarrow \mathbb{R}_+$ such that $|E(H)| = O(n/\varepsilon^2)$ and*

$$\begin{aligned} \mathcal{L}_G(w) &\preceq \mathcal{L}_H(w_H) \preceq (1 + \varepsilon)\mathcal{L}_G(w), \\ \mathcal{L}_F(w_F) &\preceq \mathcal{L}_{H \cap F}(w_H \upharpoonright_{E(H \cap F)}) \preceq (1 + \varepsilon)\mathcal{L}_F(w_F) \quad \text{for all } F \in \mathcal{F}, \end{aligned}$$

where $w_F := w \upharpoonright_{E(F)}$ is the restriction of w to the coordinates $E(F)$ and $H \cap F = (V(F), E(F) \cap E(H))$.

2 Preliminaries

For a non-negative integer n , we denote $[n] := \{1, \dots, n\}$. The non-negative reals are denoted by \mathbb{R}_+ . The set of $n \times n$ symmetric matrices is denoted by \mathbb{S}^n . The set of symmetric, $n \times n$ positive semidefinite (resp., positive definite) matrices is denoted by \mathbb{S}_+^n (resp., \mathbb{S}_{++}^n). Recall that $X \in \mathbb{S}^n$ is positive semidefinite if $v^T X v \geq 0$ for all $v \in \mathbb{R}^n$, and X is positive definite if X is positive semidefinite and $v^T X v = 0$ implies $v = 0$. Sometimes we denote $X \in \mathbb{S}_+^n$ by $X \succeq 0$ and the notation $X \succeq Y$ means that $X - Y \succeq 0$. For $X \in \mathbb{S}^n$ and $a, b \in \mathbb{R}$, the notation $X \in [a, b]$ means that $aI \preceq X \preceq bI$, where I is the identity matrix.

For $X \in \mathbb{S}^n$, its trace is $\text{Tr } X := \sum_{i=1}^n X_{ii}$, its largest (resp., smallest) eigenvalue is denoted by $\lambda_{\max}(X)$ (resp., $\lambda_{\min}(X)$). The vector space \mathbb{S}^n can be endowed with the trace inner product $\langle \cdot, \cdot \rangle$ defined by $\langle X, Y \rangle := \text{Tr}(XY) = \sum_{i,j} X_{ij} Y_{ij}$ for every $X, Y \in \mathbb{S}^n$. We shall repeatedly use that $\text{Tr}(XY) = \text{Tr}(YX)$ for any matrices X, Y for which the products XY and YX make sense.

Let $G = (V, E)$ be a graph. The canonical basis vectors of \mathbb{R}^V are $\{e_i : i \in V\}$, and the canonical basis vectors of \mathbb{R}^E are $\{e_{\{i,j\}} : \{i,j\} \in E\}$. The Laplacian of G is the linear transformation $\mathcal{L}_G(\cdot) : \mathbb{R}^E \rightarrow \mathbb{S}^V$ defined by $\mathcal{L}_G(w) = \sum_{\{i,j\} \in E} w_{\{i,j\}}(e_i - e_j)(e_i - e_j)^T$.

When dealing with Problem 2, we may assume that $B = I$. See [4, Proof of Theorem 1.1] for the details of the reduction.

3 Solving Problem 2 by Ahlswede-Winter

As mentioned earlier, Spielman and Srivastava [38] explain how Problem 1 can be solved by Rudelson's sampling lemma. This lemma can be easily generalized to handle matrices of arbitrary rank using the Ahlswede-Winter inequality, yielding a solution to Problem 2.

Let X be a random matrix such that $X = B_i / \text{Tr } B_i$ with probability $p_i := \text{Tr } B_i / \text{Tr } I$. Since $B_i \succeq 0$ and $\sum_i B_i = I$, the p_i 's define a probability distribution.

Theorem 18 ([1, Theorem 19]). *Let X, X_1, \dots, X_T be i.i.d. random variables with values in \mathbb{S}^n such that $X_i \in [0, 1]$ for every i and $\mathbf{E}(X) = \mu I$ with $\mu \in [0, 1]$. Let $\varepsilon \in (0, 1/2)$. Then*

$$\mathbf{P} \left(\frac{1}{\mu T} \sum_{i=1}^T X_i \notin [1 - \varepsilon, 1 + \varepsilon] \right) \leq 2n \cdot \exp \left(-T \frac{\varepsilon^2 \mu}{2 \ln 2} \right).$$

In our case, $\mathbf{E}(X) = (1/n)I$ and $X \in [0, 1]$. So $\mu = 1/n$. Thus, if $T > (2 \ln 2) \cdot \frac{\ln n + 2 \ln 2}{\varepsilon^2 \mu} = O(n \log n / \varepsilon^2)$, then $\mathbf{P} \left(\frac{1}{\mu T} \sum_{i=1}^T X_i \notin [1 - \varepsilon, 1 + \varepsilon] \right) < 1/2$. Thus, with constant probability, we obtain a solution y to Problem 2 where y has only $O(n \log n / \varepsilon^2)$ non-zero entries.

4 Solving Problem 2 by BSS

In our modification of the BSS algorithm [4], we keep a matrix A of the form $A = \sum_i y_i B_i$ with $y \geq 0$, starting with $A = 0$, and at each iteration we add another term αB_j to A . We enforce the invariant that the eigenvalues of A lie in $[\ell, u]$, where u and ℓ are parameters given by $u = u_0 + t\delta_U$ and $\ell = \ell_0 + t\delta_L$ after t iterations. This procedure is presented in Algorithm 1. The step of the algorithm which finds B_j and α can be done by exhaustive search on j and binary search on α . Instead of the binary search, one could also compare the quantities $U_{A(t-1)}(B_j)$ and $L_{A(t-1)}(B_j)$ defined below.

In the original BSS algorithm, the matrices are rank one: $B_j = v_j v_j^T$ for some vector v_j . Their Lemmas 3.3 and 3.4 give sufficient conditions on the new term $\alpha v_j v_j^T$ so that the invariant on the eigenvalues is maintained; Lemma 3.5 gives sufficient conditions on the remaining parameters so that a suitable new term $\alpha v_j v_j^T$ exists with $\alpha > 0$. In this section we generalize those lemmas to allow B_i matrices of arbitrary rank.

Let $A \in \mathbb{S}^n$. If $u \in \mathbb{R}$ with $\lambda_{\max}(A) < u$, define $\Phi^u(A) := \text{Tr}(uI - A)^{-1}$. If $\ell \in \mathbb{R}$ with $\lambda_{\min}(A) > \ell$, define $\Phi_\ell(A) := \text{Tr}(A - \ell I)^{-1}$. Note that $\Phi_\ell(A) = \sum_i 1/(\lambda_i - \ell)$ and $\Phi^u(A) = \sum_i 1/(u - \lambda_i)$, where $\lambda_1, \dots, \lambda_n$ are the eigenvalues of A .

Lemma 19 (Analog of Lemma 3.3 in [4]). *Let $A \in \mathbb{S}^n$ and $X \in \mathbb{S}_+^n$ with $X \neq 0$. Let $u \in \mathbb{R}$ and $\delta_U > 0$. Suppose $\lambda_{\max}(A) < u$. Let $u' := u + \delta_U$ and $M := u'I - A$. If*

$$\frac{1}{\alpha} \geq \frac{\langle M^{-2}, X \rangle}{\Phi^u(A) - \Phi^{u'}(A)} + \langle M^{-1}, X \rangle =: U_A(X),$$

then $\lambda_{\max}(A + \alpha X) < u'$ and $\Phi^{u'}(A + \alpha X) \leq \Phi^u(A)$.

Algorithm 1 A procedure for solving Problem 2 based on the BSS method.

procedure SparsifySumOfMatricesByBSS($B_1, \dots, B_m, \varepsilon$)

input: Matrices $B_1, \dots, B_m \in \mathbb{S}_+^n$ such that $\sum_i B_i = I$, and a parameter $\varepsilon \in (0, 1)$.

output: A vector y with $O(n/\varepsilon^2)$ nonzero entries such that $I \preceq \sum_i y_i B_i \preceq (1 + O(\varepsilon))I$.

Initially $A(0) := 0$ and $y(0) := 0$. Set parameters $u_0, \ell_0, \delta_L, \delta_U$ as in (5) and $T := 4n/\varepsilon^2$.

Define the potential functions $\Phi^u(A) := \text{Tr}(uI - A)^{-1}$ and $\Phi_\ell(A) := \text{Tr}(A - \ell I)^{-1}$.

For $t = 1, \dots, T$

Set $u_t := u_{t-1} + \delta_U$ and $\ell_t := \ell_{t-1} + \delta_L$.

Find a matrix B_j and a value $\alpha > 0$ such that $A(t-1) + \alpha B_j \in [\ell_t, u_t]$, and

$$\Phi^{u_t}(A(t-1) + \alpha B_j) \leq \Phi^{u_{t-1}}(A(t-1)) \quad \text{and} \quad \Phi_{\ell_t}(A(t-1) + \alpha B_j) \leq \Phi_{\ell_{t-1}}(A(t-1)).$$

Set $A(t) := A(t-1) + \alpha B_j$ and $y(t) := y(t-1) + \alpha e_j$.

Return $y(T)/\lambda_{\min}(A(T))$.

Proof. Clearly $M \succ 0$. Let $V := X^{1/2}$. By the Sherman-Morrison-Woodbury formula [13],

$$\begin{aligned} \Phi^{u'}(A + \alpha X) &= \text{Tr}(M - \alpha V V^T)^{-1} = \text{Tr}(M^{-1} + \alpha M^{-1} V (I - \alpha V^T M^{-1} V)^{-1} V^T M^{-1}) \\ &= \Phi^{u'}(A) + \text{Tr}(\alpha M^{-1} V (I - \alpha V^T M^{-1} V)^{-1} V^T M^{-1}). \end{aligned}$$

Since $M^{-1} \succ 0$, $X \neq 0$ and $\Phi^u(A) > \Phi^{u'}(A)$, our hypotheses imply $1/\alpha > \langle M^{-1}, X \rangle = \text{Tr}(V^T M^{-1} V) \geq \lambda_{\max}(V^T M^{-1} V) \geq 0$, so $\beta := \lambda_{\min}(I - \alpha V^T M^{-1} V) = 1 - \alpha \lambda_{\max}(V^T M^{-1} V) > 0$ and by, e.g., [18, Corollary 7.7.4],

$$0 \prec \beta I \preceq I - \alpha V^T M^{-1} V \implies 0 \prec (I - \alpha V^T M^{-1} V)^{-1} \preceq \beta^{-1} I.$$

Thus,

$$\begin{aligned} \Phi^{u'}(A + \alpha X) &\leq \Phi^{u'}(A) + \alpha \beta^{-1} \text{Tr}(V^T M^{-2} V) \\ &= \Phi^u(A) - (\Phi^u(A) - \Phi^{u'}(A)) + \alpha \beta^{-1} \langle M^{-2}, X \rangle \end{aligned}$$

To prove that $\Phi^{u'}(A + \alpha X) \leq \Phi^u(A)$, it suffices to show that $\alpha \beta^{-1} \langle M^{-2}, X \rangle \leq \Phi^u(A) - \Phi^{u'}(A)$. This is equivalent to

$$\frac{\langle M^{-2}, X \rangle}{1/\alpha - \lambda_{\max}(V^T M^{-1} V)} \leq \Phi^u(A) - \Phi^{u'}(A),$$

which follows from $1/\alpha \geq U_A(X)$ since $\lambda_{\max}(V^T M^{-1} V) \leq \text{Tr}(V^T M^{-1} V) = \langle M^{-1}, X \rangle$.

It remains to show that $\lambda_{\max}(A + \alpha X) < u'$. Suppose not. Choose $\varepsilon \in (0, \delta_U)$ such that $1/\varepsilon > \Phi^u(A)$. By continuity, for some $\alpha' \in (0, \alpha)$ we have $\lambda_{\max}(A + \alpha' X) = u' - \varepsilon$. Since $1/\alpha' \geq 1/\alpha \geq U_A(X)$, we get $\Phi^{u'}(A + \alpha' X) \geq 1/\varepsilon > \Phi^u(A) \geq \Phi^{u'}(A + \alpha' X)$, a contradiction. \square

Lemma 20 (Analog of Lemma 3.4 in [4]). *Let $A \in \mathbb{S}^n$ and $X \in \mathbb{S}_+^n$, with $n \geq 2$. Let $\ell \in \mathbb{R}$ and $\delta_L > 0$. Suppose $\lambda_{\min}(A) > \ell$ and $\Phi_\ell(A) \leq 1/\delta_L$. Let $\ell' := \ell + \delta_L$ and $N := A - \ell' I$. If*

$$0 < \frac{1}{\alpha} \leq \frac{\langle N^{-2}, X \rangle}{\Phi_{\ell'}(A) - \Phi_\ell(A)} - \langle N^{-1}, X \rangle =: L_A(X),$$

then $\lambda_{\min}(A + \alpha X) > \ell'$ and $\Phi_{\ell'}(A + \alpha X) \leq \Phi_\ell(A)$. Moreover, $N \succ 0$.

Proof. Note that $\lambda_{\min}(A) > \ell$ and $\Phi_\ell(A) \leq 1/\delta_L$ imply that $N \succ 0$, and therefore $\lambda_{\min}(A + \alpha X) > \ell'$. Let $V := X^{1/2}$. By the Sherman-Morrison-Woodbury formula,

$$\begin{aligned}\Phi_{\ell'}(A + \alpha X) &= \text{Tr}(N + \alpha VV^T)^{-1} = \text{Tr}(N^{-1} - \alpha N^{-1}V(I + \alpha V^T N^{-1}V)^{-1}V^T N^{-1}) \\ &= \Phi_{\ell'}(A) - \text{Tr}(\alpha N^{-1}V(I + \alpha V^T N^{-1}V)^{-1}V^T N^{-1}).\end{aligned}$$

For $\beta := \lambda_{\max}(I + \alpha V^T N^{-1}V)$, we have

$$0 \prec I + \alpha V^T N^{-1}V \preceq \beta I \implies 0 \prec \beta^{-1}I \preceq (I + \alpha V^T N^{-1}V)^{-1}.$$

Thus,

$$\begin{aligned}\Phi_{\ell'}(A + \alpha X) &\leq \Phi_{\ell'}(A) - \alpha\beta^{-1} \text{Tr}(V^T N^{-2}V) \\ &= \Phi_{\ell'}(A) + (\Phi_{\ell'}(A) - \Phi_\ell(A)) - \alpha\beta^{-1} \langle N^{-2}, X \rangle\end{aligned}$$

We will be done if we show that $\alpha\beta^{-1} \langle N^{-2}, X \rangle \geq \Phi_{\ell'}(A) - \Phi_\ell(A)$. This is equivalent to

$$\frac{\langle N^{-2}, X \rangle}{1/\alpha + \lambda_{\max}(V^T N^{-1}V)} \geq \Phi_{\ell'}(A) - \Phi_\ell(A)$$

which follows from $0 < 1/\alpha \leq L_A(X)$, since $\Phi_{\ell'}(A) > \Phi_\ell(A)$, $N \succ 0$, and $\lambda_{\max}(V^T N^{-1}V) \leq \text{Tr}(V^T N^{-1}V) = \langle N^{-1}, X \rangle$. \square

The next lemma can be proved by a syntactic modification of the proof of Lemma 3.5 in [4].

Lemma 21 (Analog of Lemma 3.5 in [4]). *Let $A \in \mathbb{S}^n$ with $n \geq 2$, and let $u, \ell \in \mathbb{R}$ and $\varepsilon_U, \delta_U, \varepsilon_L, \delta_L > 0$ such that $\lambda_{\max}(A) < u$, $\lambda_{\min}(A) > \ell$, $\Phi^u(A) \leq \varepsilon_U$, and $\Phi_\ell(A) \leq \varepsilon_L$. Let $B_1, \dots, B_m \in \mathbb{S}^n$ such that $\sum_i B_i = I$. If*

$$0 \leq \frac{1}{\delta_U} + \varepsilon_U \leq \frac{1}{\delta_L} - \varepsilon_L \quad (3)$$

then there exists $j \in [m]$ and $\alpha > 0$ for which $L_A(B_j) \geq 1/\alpha \geq U_A(B_j)$.

Proof. As in [4, Lemma 3.5], it suffices to show that $\sum_i L_A(B_i) \geq \sum_i U_A(B_i)$. Let $u' := u + \delta_U$, $M := u'I - A$, $\ell' := \ell + \delta_L$, and $N := A - \ell'I$. It follows from the bilinearity of $\langle \cdot, \cdot \rangle$ and the assumption $\sum_i B_i = I$ that

$$\sum_i U_A(B_i) = \frac{\text{Tr } M^{-2}}{\Phi^u(A) - \Phi^{u'}(A)} + \text{Tr } M^{-1} \quad (4a)$$

$$\sum_i L_A(B_i) = \frac{\text{Tr } N^{-2}}{\Phi_{\ell'}(A) - \Phi_\ell(A)} - \text{Tr } N^{-1} \quad (4b)$$

It is shown in [4, Lemma 3.5] that (4a) is at most (4b), completing the proof. \square

Now we set the parameters of Lemma 21 similarly as in [4]:

$$\delta_L := 1 \quad \varepsilon_L := \frac{\varepsilon}{2} \quad \ell_0 := -\frac{n}{\varepsilon_L} \quad \delta_U := \frac{2 + \varepsilon}{2 - \varepsilon} \quad \varepsilon_U := \frac{\varepsilon}{2\delta_U} \quad u_0 := \frac{n}{\varepsilon_U}. \quad (5)$$

So (3) holds with equality. If A is the matrix obtained after $T = 4n/\varepsilon^2$ iterations, then

$$\frac{\lambda_{\max}(A)}{\lambda_{\min}(A)} \leq \frac{u_0 + T\delta_U}{\ell_0 + T\delta_L} = \left(\frac{2 + \varepsilon}{2 - \varepsilon} \right)^2 \leq \frac{1 + \varepsilon}{1 - \varepsilon}$$

so $A' := A/\lambda_{\min}(A)$ satisfies $I \preceq A' \preceq (1 + \varepsilon)I/(1 - \varepsilon)$ and A' is a positive linear combination of $O(n/\varepsilon^2)$ of the matrices B_i .

It is easy to check that the previous lemmas also hold if we replace the set \mathbb{S}^n of symmetric matrices of size $n \times n$ by the set \mathbb{H}^n of Hermitian matrices of size $n \times n$.

4.1 Running Time

At each iteration, we must compute $U_A(B_j)$ and $L_A(B_j)$ for each $j \in [m]$. The functions $U_A(X)$ and $L_A(X)$ are the inner products of X with certain matrices that can be obtained from A in time $O(n^3)$. Thus, each iteration runs in time $O(n^3 + mn^2) = O(mn^2)$, and the total running time after $T = 4n/\varepsilon^2$ iterations is $O(mn^3/\varepsilon^2)$. We remark that the reduction to the case $B = I$ can be made in time $O(mn^3)$. This concludes the proof of Theorem 3.

If the matrices B_i have $O(1)$ nonzero entries, as in the graph sparsification problem, the algorithm can be made to run in time $O(n^4/\varepsilon^2 + mn/\varepsilon^2)$. We briefly sketch the details. To reduce the problem to the case that $B = I$, we first compute $(B^+)^{1/2}$, where B^+ is the Moore-Penrose pseudoinverse of B . Define the function $f(X) := (B^+)^{1/2}X(B^+)^{1/2}$ on \mathbb{S}^n .

The reduction now calls for replacing each input matrix B_i by $f(B_i)$ and the matrix B by $f(B)$. But we shall not do this. Instead, we do some preprocessing at each iteration as follows. The function $U_A(X)$ (as well as $L_A(X)$) is the inner product of X with a certain matrix V . Hence, $U_A(f(B_j)) = \langle V, f(B_j) \rangle = \langle f(V), B_j \rangle$ for every j , since f is self-adjoint. Thus, to compute $U_A(f(B_j))$ for each j , we first compute the matrix $f(V)$ in time $O(n^3)$, and now the inner product $U_A(f(B_j)) = \langle f(V), B_j \rangle$ can be computed in constant time for each j , since B_j has $O(1)$ nonzero entries. Thus, each iteration runs in time $O(n^3 + m)$ and the total running time is $O(n^4/\varepsilon^2 + mn/\varepsilon^2)$.

5 Solving Problem 2 by MMWUM

Observe that the set of all vectors y that are feasible for (2) is the feasible region of a semidefinite program (SDP). So solving Problem 2 amounts to finding a *sparse* solution to this SDP. Here “sparse” means that there are few non-zero entries in the solution y ; this differs from other notions of “low-complexity” SDP solutions, such as the low-rank solutions studied by So, Ye and Zhang [36].

It has long been known that the multiplicative weight update method can be used to construct sparse solutions for some linear programs. A prominent example is the construction of sparse, low-regret solutions to zero-sum games [9, 43, 44]. (Another example is the work of Charikar et al. [7] on approximating metrics by few tree metrics.) Building on that idea, one might imagine that Arora and Kale’s matrix multiplicative update method (MMWUM) [3] can construct sparse solutions to (2). In this section, we show that this is indeed possible: we obtain a solution y to Problem 2 with $O(n \log n/\varepsilon^3)$ nonzero entries.

5.1 Overview of MMWUM

The MMWUM is an algorithm that helps us approximately solve an SDP feasibility problem. The gist of (a slight modification of) the method is contained in the following result (its proof can be found in Appendix B):

Theorem 22. *Let T, K, n_1, \dots, n_K be positive integers. Let $C_k, A_{1,k}, \dots, A_{m,k} \in \mathbb{S}^{n_k}$ for $k \in [K]$. For each $k \in [K]$, let $\eta_k > 0$ and $0 < \beta_k \leq 1/2$. Given $X_1, \dots, X_K \in \mathbb{S}^n$, consider the system*

$$\sum_{i=1}^m y_i \langle A_{i,k}, X_k \rangle \geq \langle C_k, X_k \rangle - \eta_k \operatorname{Tr} X_k, \quad \forall k \in [K], \quad \text{and} \quad y \in \mathbb{R}_+^m. \quad (6)$$

For each $k \in [K]$, let $\{\mathcal{P}_k, \mathcal{N}_k\}$ be a partition of $[T]$, let $0 < \ell_k \leq \rho_k$, and let $W_k^{(t)} \in \mathbb{S}^n$ and $\ell_k^{(t)} \in \mathbb{R}$ for $t \in [T+1]$. Let $y^{(t)} \in \mathbb{R}^m$ for $t \in [T]$. Suppose the following properties hold:

$$W_k^{(t+1)} = \exp \left(-\frac{\beta_k}{\ell_k + \rho_k} \sum_{\tau=1}^t \left[\sum_{i=1}^m y_i^{(\tau)} A_{i,k} - C_k + \ell_k^{(\tau)} I \right] \right), \quad \forall t \in \{0, \dots, T\}, \forall k \in [K],$$

$$\begin{aligned}
y = y^{(t)} \text{ is a solution for (6) with } X_k = W_k^{(t)}, \forall k \in [K], \quad \forall t \in [T], \\
\sum_{i=1}^m y_i^{(t)} A_{i,k} - C_k \in \begin{cases} [-\ell_k, \rho_k], & \text{if } t \in \mathcal{P}_k, \\ [-\rho_k, \ell_k], & \text{if } t \in \mathcal{N}_k, \end{cases} \quad \forall t \in [T], k \in [K], \\
\ell_k^{(t)} = \ell_k, \quad \forall t \in \mathcal{P}_k, \forall k \in [K], \quad \text{and} \quad \ell_k^{(t)} = -\ell_k, \quad \forall t \in \mathcal{N}_k, \forall k \in [K].
\end{aligned}$$

Define $\bar{y} := \frac{1}{T} \sum_{t=1}^T y^{(t)}$. Then,

$$\sum_{i=1}^m \bar{y}_i A_{i,k} - C_k \succeq -\left[\beta_k \ell_k + \frac{(\rho_k + \ell_k) \ln n}{T \beta_k} + (1 + \beta_k) \eta_k \right] I, \quad \forall k \in [K]. \quad (7)$$

Take $K = 2$, set $C_1 := I$ and $C_2 := -I$, and put $A_{i,1} := B_i$ and $A_{i,2} := -B_i$ for each $i \in [m]$. Then Theorem 22 shows that finding a solution for (2) reduces to constructing an oracle that solves linear systems of the form (6) with a few extra technical properties involving the parameters ℓ_k and ρ_k , and adjusting the other parameters so that the error term on the right-hand side of (7) is $\leq \varepsilon$.

To obtain a feasible solution for (2) that is also sparse, the idea is to design an implementation of the oracle that returns a vector $y^{(t)}$ with only *one* nonzero entry at each iteration t of MMWUM, and to adjust the parameters so that, after $T = O(n \log n / \varepsilon^3)$ iterations, the smallest and largest eigenvalues of $\sum_{i=1}^m \bar{y}_i B_i$ are ε -close to 1. Since \bar{y} is the average of the $y^{(t)}$'s, the resulting \bar{y} will have at most T nonzero entries.

We set the remaining parameters as follows:

$$\begin{aligned}
\beta := \beta_1 := \beta_2 := \frac{\varepsilon}{4}, \quad T := \frac{2(\rho + \ell) \ln n}{\beta \varepsilon}, \quad \eta := \eta_1 := \eta_2 := \frac{\varepsilon}{8}, \\
\ell := \ell_1 := \ell_2 := 1, \quad \rho := \rho_1 := \rho_2 := \frac{1 + \eta}{\eta} n, \quad \mathcal{P}_1 := \mathcal{N}_2 := [T], \quad \mathcal{N}_1 := \mathcal{P}_2 := \emptyset.
\end{aligned}$$

Then the error term on the right-hand side of (7) is

$$\beta \ell + \frac{(\rho + \ell) \ln n}{T \beta} + (1 + \beta) \eta = \frac{\varepsilon}{4} + \frac{\varepsilon}{2} + \left(1 + \frac{\varepsilon}{4}\right) \frac{\varepsilon}{8} = \frac{7\varepsilon}{8} + \frac{\varepsilon^2}{32} \leq \varepsilon. \quad (8)$$

Thus, (2) follows from (7) and (8). Moreover, $T = O(n \log n / \varepsilon^3)$, as desired.

5.2 The Oracle

It remains to implement the oracle. Consider an iteration t , and let $X_1 := W_1^{(t)}$ and $X_2 := W_2^{(t)}$ be given. We must find $y := y^{(t)} \in \mathbb{R}_+^m$ with at most one nonzero entry such that

$$\sum_{i=1}^m y_i \langle X_1, B_i \rangle \geq (1 - \eta) \text{Tr } X_1, \quad \sum_{i=1}^m y_i \langle X_2, B_i \rangle \leq (1 + \eta) \text{Tr } X_2, \quad \text{and} \quad \sum_{i=1}^m y_i B_i \in [0, \rho].$$

Since y should have only one nonzero entry, it suffices to find $j \in [m]$ and $\alpha \in \mathbb{R}_+$ such that

$$\begin{aligned}
\alpha \langle X_1, B_j \rangle &\geq (1 - \eta) \text{Tr } X_1, \\
\alpha \langle X_2, B_j \rangle &\leq (1 + \eta) \text{Tr } X_2, \\
\alpha \text{Tr } B_j &\leq \rho.
\end{aligned} \quad (9)$$

Here we are using the fact that $\lambda_{\max}(B_j) \leq \text{Tr } B_j$ since $B_j \succeq 0$. We will show that such j and α exist. Due to the definition of W_1 and W_2 , the oracle can assume that X_1 is a scalar multiple of X_2^{-1} , although we will not make use of that fact.

Proposition 23. Let $B_1, \dots, B_m \in \mathbb{S}_+^n$ such that $\sum_{i=1}^m B_i = I$. Let $\eta > 0$ and $X_1, X_2 \in \mathbb{S}_{+\eta}^n$. Then, for $\rho := (1 + \eta)n/\eta$, there exist $j \in [m]$ and $\alpha \geq 0$ such that (9) holds.

Proof. By possibly dropping some B_i 's, we may assume that $B_i \neq 0$ for every $i \in [m]$. Define $p_i := \langle X_1, B_i \rangle / \text{Tr } X_1 > 0$ for every $i \in [m]$. Consider the probability space on $[m]$ where j is sampled from $[m]$ with probability p_j . The fact that $\sum_{j=1}^m p_j = 1$ follows from $\sum_{i=1}^m B_i = I$. Then $\mathbf{E}_j[p_j^{-1} \text{Tr } B_j] = \sum_{i=1}^m \text{Tr } B_i = \text{Tr } I = n$. By Markov's inequality,

$$\mathbf{P} \left(p_j^{-1} \text{Tr } B_j \leq \frac{(1 + \eta)n}{\eta} \right) = 1 - \mathbf{P} \left(p_j^{-1} \text{Tr } B_j > \frac{(1 + \eta)n}{\eta} \right) > 1 - \frac{\eta}{1 + \eta} = \frac{1}{1 + \eta}. \quad (10)$$

Next note that $\mathbf{E}_j[p_j^{-1} \langle X_2, B_j \rangle] = \sum_{i=1}^m \langle X_2, B_i \rangle = \langle X_2, I \rangle = \text{Tr } X_2$. Together with Markov's inequality, this yields

$$\mathbf{P} \left(p_j^{-1} \langle X_2, B_j \rangle \leq (1 + \eta) \text{Tr } X_2 \right) = 1 - \mathbf{P} \left(p_j^{-1} \langle X_2, B_j \rangle > (1 + \eta) \text{Tr } X_2 \right) > 1 - \frac{1}{1 + \eta}. \quad (11)$$

It follows from (10) and (11) that there exists $j \in [m]$ satisfying

$$p_j^{-1} \langle X_2, B_j \rangle \leq (1 + \eta) \text{Tr } X_2, \quad \text{and} \quad p_j^{-1} \text{Tr } B_j \leq \frac{1 + \eta}{\eta} n = \rho.$$

Set $\alpha := p_j^{-1}$ and note that

$$\alpha \langle X_1, B_j \rangle = p_j^{-1} \langle X_1, B_j \rangle = \text{Tr } X_1 \geq (1 - \eta) \text{Tr } X_1.$$

Hence, j and α satisfy (9). □

The following proposition, proven in Appendix C, shows that the parameters achieved by Proposition 23 is essentially optimal.

Proposition 24. Any oracle for satisfying (9) must have $\rho = \Omega(n/\eta)$, even if the B_i matrices have rank one, and even if X_1 is a scalar multiple of X_2^{-1} .

We also point out that a naive application of MMWUM as stated by Kale in [22] does not work. In his description of MMWUM, the parameter K is fixed as 1. So we must correspondingly adjust our input matrices to be block-diagonal, e.g., C has two blocks: I and $-I$. However, applying Theorem 22 in this manner would lead to a sparsifier with $\Omega(n^2)$ edges. The reason is that the parameter ρ needs to be $\Omega(n)$, and we must choose $\ell = \rho$ since the spectrum of $\sum_{i=1}^m y_i A_i - C$ is symmetric around zero for any y . Thus, to get the error term on the right-hand side of (7) to be $\leq \varepsilon$, we would need to take $T = \Omega(n^2)$.

6 Solving Problem 2 by a Width-Free MMWUM

The algorithm of Section 5 solves Problem 2 with only $O(n \log n / \varepsilon^3)$ nonzero entries, which is slightly worse than the $O(n \log n / \varepsilon^2)$ nonzero entries achieved by the Ahlswede-Winter method discussed in Section 3. The main reason for this discrepancy is that MMWUM requires us to bound the “width” of the oracle using the parameter ρ ; formally, the oracle must satisfy the inequality $\alpha \text{Tr } B_j \leq \rho$ in (9). In order to satisfy this width constraint, the oracle loses an extra factor of $O(1/\varepsilon)$, and this is necessary as shown in Proposition 24.

In this section, we slightly refine MMWUM to avoid its dependence on the width. This allows us to simplify our oracle and avoid losing the extra factor of $O(1/\varepsilon)$. We obtain a solution to Problem 2 with only

only $O(n \log n / \varepsilon^2)$ nonzero entries, matching the sparsity of the solutions obtained by the Ahlswede-Winter inequality.

The following theorem is our width-free variant of MMWUM. We remark that the method described in this theorem is geared towards solving Problem 2 and is not necessarily useful for all applications of MMWUM.

Theorem 25. *Let T be a positive integer. Let $B_1, \dots, B_m \in \mathbb{S}_+^n$ be nonzero. Let $\gamma, \eta, \delta_L, \delta_U > 0$. For any given $X_L, X_U \in \mathbb{S}^n$, consider the system*

$$\begin{aligned} \delta_U &\geq \frac{\exp(\gamma\alpha \operatorname{Tr} B_j) - 1}{\operatorname{Tr} B_j} \langle X_U, B_j \rangle, \\ \delta_L &\leq \frac{1 - \exp(-\gamma\alpha \operatorname{Tr} B_j)}{\operatorname{Tr} B_j} \langle X_L, B_j \rangle, \\ \alpha &\in \mathbb{R}_+, \quad j \in [m]. \end{aligned} \tag{12}$$

For each $t \in \{0, \dots, T+1\}$, let $A(t), W_L(t), W_U(t) \in \mathbb{S}^n$, let $\alpha(t) \in \mathbb{R}_+$, and let $j(t) \in [m]$. Suppose the following properties hold:

$$A(t) = \sum_{\tau=1}^t \alpha(\tau) B_{j(\tau)}, \quad \forall t \in \{0, \dots, T\},$$

$$W_U(t+1) = \exp(\gamma A(t)) \quad \text{and} \quad W_L(t+1) = \exp(-\gamma A(t)), \quad \forall t \in \{0, \dots, T\},$$

$$(\alpha, B_j) = (\alpha(t), B_{j(t)}) \text{ is a solution for (12) with } (X_U, X_L) = \left(\frac{W_U(t)}{\operatorname{Tr} W_U(t)}, \frac{W_L(t)}{\operatorname{Tr} W_L(t)} \right), \quad \forall t \in [T].$$

Then

$$\frac{A(T)}{T} \in \left[\frac{\log(1 - \delta_L)^{-1}}{\gamma} - \frac{\log n}{T\gamma}, \frac{\log(1 + \delta_U)}{\gamma} + \frac{\log n}{T\gamma} \right]. \tag{13}$$

Proof. We will use Golden-Thompson inequality:

$$\operatorname{Tr}(\exp(A + B)) \leq \operatorname{Tr}(\exp(A) \exp(B)), \quad \forall A, B \in \mathbb{S}^n. \tag{14}$$

We will also make use of the following facts. First,

$$\exp(cx) \leq 1 + \frac{\exp(c \cdot b) - 1}{b} x \quad \forall c \in \mathbb{R}, b > 0, x \in [0, b].$$

For $X \in \mathbb{S}_+^n$, we have $\lambda_{\max}(X) \leq \operatorname{Tr} X$, so $X \in [0, \operatorname{Tr} X]$, and

$$\exp(cX) \preceq I + \frac{\exp(c \cdot \operatorname{Tr} X) - 1}{\operatorname{Tr} X} X. \tag{15}$$

For each $t \in [T+1]$, define $\Phi_L(t) := \operatorname{Tr} W_L(t)$ and $\Phi_U(t) = \operatorname{Tr} W_U(t)$. For each $t \in [T]$,

$$\begin{aligned} \Phi_U(t+1) &= \operatorname{Tr} \left(\exp(\gamma A(t)) \right) = \operatorname{Tr} \left(\exp(\gamma A(t-1) + \gamma \alpha B_j) \right) \\ &\stackrel{(14)}{\leq} \operatorname{Tr} \left(\exp(\gamma A(t-1)) \exp(\gamma \alpha B_j) \right) \\ &\stackrel{(15)}{\leq} \operatorname{Tr} \left(\exp(\gamma A(t-1)) \left(\frac{\exp(\gamma \alpha \operatorname{Tr} B_j) - 1}{\operatorname{Tr} B_j} B_j + I \right) \right) \\ &= \frac{\exp(\gamma \alpha \operatorname{Tr} B_j) - 1}{\operatorname{Tr} B_j} \operatorname{Tr}(\exp(\gamma A(t-1)) B_j) + \operatorname{Tr}(\exp(\gamma A(t-1))) \\ &= \frac{\exp(\gamma \alpha \operatorname{Tr} B_j) - 1}{\operatorname{Tr} B_j} \langle W_U(t), B_j \rangle + \Phi_U(t) \\ &\stackrel{(12)}{\leq} (1 + \delta_U) \Phi_U(t), \end{aligned} \tag{16}$$

where we abbreviated $j := j(t)$ and $\alpha := \alpha(t)$.

Since $A(0) = 0$, we have that $\Phi_U(1) = \text{Tr } I = n$. Using (16), after T iterations,

$$\Phi_U(T+1) \leq (1 + \delta_U)^T n.$$

Thus,

$$\exp(\gamma \lambda_{\max}(A(T))) \leq \sum_{i=1}^n \exp(\gamma \lambda_i) = \text{Tr } W_U(T+1) = \Phi_U(T+1) \leq (1 + \delta_U)^T n,$$

where $\lambda_1, \dots, \lambda_n$ are the eigenvalues of $A(T)$. And so $\gamma \lambda_{\max}(A(T)) \leq T \log(1 + \delta_U) + \log n$, which implies the upper bound in (13). The proof of the lower bound is analogous. \square

Next we establish conditions under which we can construct an oracle for solving the system (12). The proof consists of algebraic manipulations and an averaging argument analogous to the proof of Lemma 3.5 in [4].

Theorem 26. *Let $B_1, \dots, B_m \in \mathbb{S}_+^n$ be nonzero such that $\sum_{i=1}^m B_i = I$. Let $\delta_U, \delta_L > 0$ be such that*

$$\frac{1}{\delta_L} - n \geq \frac{1}{\delta_U}. \quad (17)$$

Then, for any $X_L, X_U \in \mathbb{S}_{++}^n$ with trace one, the system (12) has a solution.

Proof. The first inequality in (12) is equivalent to

$$\frac{\text{Tr } B_j}{\exp(\gamma \alpha \text{Tr } B_j) - 1} \geq \frac{\langle X_U, B_j \rangle}{\delta_U}. \quad (18)$$

Using the identity $\frac{1}{1-1/x} = 1 + \frac{1}{x-1}$, the second inequality in (12) is equivalent to

$$\frac{\text{Tr } B_j}{\exp(\gamma \alpha \text{Tr } B_j) - 1} \leq \frac{\langle X_L, B_j \rangle}{\delta_L} - \text{Tr } B_j. \quad (19)$$

We will choose $j \in [m]$ so that

$$\frac{\langle X_L, B_j \rangle}{\delta_L} - \text{Tr } B_j \geq \frac{\langle X_U, B_j \rangle}{\delta_U} \quad (20)$$

and set α so that (18) holds with equality. Then both (18) and (19) will hold. Note that $\alpha \geq 0$ since $e^{\gamma \alpha \text{Tr } B_j} = 1 + \delta_U \text{Tr } B_j / \langle X_U, B_j \rangle > 1$ and $\gamma \text{Tr } B_j > 0$.

To see that there exists $j \in [m]$ satisfying (20), note that, by (17) and $\sum_{i=1}^m B_i = I$,

$$\sum_{i=1}^m \left[\frac{\langle X_L, B_i \rangle}{\delta_L} - \text{Tr } B_i \right] = \frac{\text{Tr } X_L}{\delta_L} - n = \frac{1}{\delta_L} - n \geq \frac{1}{\delta_U} = \frac{\text{Tr } X_U}{\delta_U} = \sum_{i=1}^m \frac{\langle X_U, B_i \rangle}{\delta_U}.$$

This concludes the proof. \square

Finally, let us show how to set the parameters to get a sparsifier. Given $\varepsilon \in (0, 1)$, set

$$\eta := \varepsilon/2, \quad \delta_U := \frac{\eta}{n}, \quad \delta_L := \frac{\eta}{(1+\eta)n}, \quad T := \frac{n \log n}{\eta^2}. \quad (21)$$

By our choice of δ_L and δ_U , we have $1/\delta_L - n = (1 + \eta)n/\eta - n = n/\eta = 1/\delta_U$, so (17) holds with equality. After we run the modified version of MMWUM given by Theorem 25, we obtain a matrix $A(T)$. Set $\bar{A} := A(T)/T$. By Theorem 25,

$$\lambda_{\max}(\bar{A}) \leq \frac{\log(1 + \delta_U)}{\gamma} + \frac{\log n}{T\gamma} \leq \left(\delta_U + \frac{\eta^2}{n}\right)/\gamma = \frac{1 + \eta}{n\gamma/\eta}.$$

We will use that $-\log(1 - x) \geq x$ for $x < 1$. Thus,

$$\lambda_{\min}(\bar{A}) \geq \frac{\log(1 - \delta_L)^{-1}}{\gamma} - \frac{\log n}{T\gamma} \geq \left(\delta_L - \frac{\eta^2}{n}\right)/\gamma = \frac{1/(1 + \eta) - \eta}{n\gamma/\eta} \geq \frac{1 - 2\eta}{n\gamma/\eta}.$$

So if we choose $\gamma = \eta/n$ then $(1 - \varepsilon)I \preceq \bar{A} \preceq (1 + \varepsilon)I$ and \bar{A} is of the form $\sum_i y_i B_i$ with $y \geq 0$ and has at most $T = O(n \log n / \varepsilon^2)$ nonzero entries.

Remark. The choice of γ is actually irrelevant here. We could choose $\gamma > 0$ arbitrarily, then define $\bar{A} = A(T) \cdot (n\gamma/\eta T)$ and the desired conclusion would hold.

7 Solving Problem 2 by Pessimistic Estimators

An anonymous reviewer for a preliminary draft of this paper raised the possibility of designing another deterministic solution to Problem 2. The proposal was to use the pessimistic estimators of Wigderson and Xiao [42] to derandomize the random sampling approach of Section 3. In this section we show that this proposal indeed works. We remark that pessimistic estimators were also used by Hofmeister and Lefmann [17] to derandomize the proof of Theorem 15.

It is known that there is a close relationship between pessimistic estimators and multiplicative weight update methods. (See, for example, the work of Young [44].) However, the two methods are not identical, and in particular the algorithm presented in this section is not identical to either of our algorithms based on MMWUM. To illustrate one difference, notice that the algorithm in Section 3 has the property that its output vector y has every component y_i equal to an integer multiple of $n/(T \cdot \text{Tr } B_i)$. The algorithm of this section also has that property as it is a derandomization of the algorithm in Section 3. However, the algorithms in Sections 4, 5 and 6 do not have that property.

Definition 27 (Definition 3.1 in [42]). *Let $\vec{X} = (X_1, \dots, X_T)$ be random variables distributed over $[m]$. Let S be an event with $\mathbf{P}(\vec{X} \in S) > 0$. We say that ϕ_0, \dots, ϕ_T , $\phi_i : [m]^i \rightarrow [0, 1]$, are pessimistic estimators for S if the following hold.*

1. *For any i and any fixed $x_1, \dots, x_i \in [m]$, we have that*

$$\mathbf{P}_{X_{i+1}, \dots, X_T} \left((x_1, \dots, x_i, X_{i+1}, \dots, X_T) \notin S \right) \leq \phi_i(x_1, \dots, x_i).$$

2. *For any i and any fixed $x_1, \dots, x_i \in [n]$:*

$$\mathbf{E}_{X_{i+1}} (\phi_{i+1}(x_1, \dots, x_i, X_{i+1})) \leq \phi_i(x_1, \dots, x_i).$$

Note that the function ϕ_0 depends on no variables and is therefore just a scalar in $[0, 1]$. A nice property of this definition is that it allows compositions very easily. That is, if we have pessimistic estimators ϕ_0, \dots, ϕ_T and ψ_0, \dots, ψ_T for events S and S' , resp., then $\phi_0 + \psi_0, \dots, \phi_T + \psi_T$ are pessimistic estimators for the event $S \cap S'$ (see Lemma 3.3 in [42]).

The key point of this method is that, if there are pessimistic estimators ϕ_0, \dots, ϕ_T , such that $\phi_0 < 1$ and each ϕ_i can be computed efficiently, then one can find $(x_1, \dots, x_T) \in S$ efficiently.

Let X_1, \dots, X_T be i.i.d. random variables with same distribution as the random variable X as defined in Section 3. Wigderson and Xiao [42] considered the event

$$S_{\geq} = \{\vec{X} : \frac{1}{T} \sum_{i=1}^T X_i \succeq (1 - \varepsilon)\mu I\}$$

and obtained¹ the following pessimistic estimators:

$$\begin{aligned} \phi_0 &= n e^{tT(1-\varepsilon)\mu} \|\mathbf{E}_X (\exp(-tX))\|^T \leq n \exp(-T\varepsilon^2\mu/(2\ln 2)); \\ \phi_i(x_1, \dots, x_i) &:= e^{tT(1-\varepsilon)\mu} \text{Tr} \left(\exp\left(-\sum_{j=1}^i tx_j\right) \right) \cdot \|\mathbf{E}_X (\exp(-tX))\|^{T-i}, \end{aligned}$$

where $t = \log \left(\frac{1-(1-\varepsilon)\mu}{(1-\mu)(1-\varepsilon)} \right)$. Similarly, for the event $S_{\leq} = \{\vec{X} : \frac{1}{T} \sum_{i=1}^T X_i \preceq (1 + \varepsilon)\mu I\}$, one can find the following pessimistic estimators

$$\begin{aligned} \psi_0 &= n e^{-t'T(1+\varepsilon)\mu} \|\mathbf{E}_X (\exp(t'X))\|^T \leq n \exp(-T\varepsilon^2\mu/(2\ln 2)); \\ \psi_i(x_1, \dots, x_i) &:= e^{-t'T(1+\varepsilon)\mu} \text{Tr} \left(\exp\left(\sum_{j=1}^i t'x_j\right) \right) \cdot \|\mathbf{E}_X (\exp(t'X))\|^{T-i}, \end{aligned}$$

where $t' = \log \left(\frac{(1+\varepsilon)(1-\mu)}{1-(1+\varepsilon)\mu} \right)$. If we choose $T > (2\ln 2) \ln(2n)/(\varepsilon^2\mu) = (2\ln 2)n \ln(2n)/\varepsilon^2$, then $\phi_0 + \psi_0 < 1$. Each ϕ_i, ψ_i can be computed efficiently and so one can find in polynomial time $(x_1, \dots, x_T) \in S_{\geq} \cap S_{\leq}$.

8 Comparing BSS and MMWUM

In this section we show a striking similarity between the algorithms presented in Sections 4 and 6. The proof of Theorem 25 defines two potential functions for each iteration t .

$$\begin{aligned} \Phi_U(t) &:= \text{Tr } W_U(t) = \text{Tr} \exp(\gamma A(t)) \\ \Phi_L(t) &:= \text{Tr } W_L(t) = \text{Tr} \exp(-\gamma A(t)) \end{aligned}$$

The proof shows that, for the algorithm of Section 6, the potentials must change as follows:

$$\begin{aligned} \Phi_U(t+1) &\leq (1 + \delta_U)\Phi_U(t) & \forall t \in \{0, \dots, T-1\} \\ \Phi_L(t+1) &\leq (1 - \delta_L)\Phi_L(t) & \forall t \in \{0, \dots, T-1\}. \end{aligned} \tag{22}$$

Instead of requiring these potentials to grow and shrink in this way, we could instead parameterize the potential functions by the iteration number t and then simply require that the potential do not grow from iteration to iteration. To formalize this alternative approach, let us define the new potential functions

$$\begin{aligned} \Psi^u(A) &:= \text{Tr} \exp(-uI + \gamma A), \\ \Psi_\ell(A) &:= \text{Tr} \exp(\ell I - \gamma A) \end{aligned}$$

and define the parameters $\Delta_U = \ln(1 + \delta_U)$ and $\Delta_L = \ln((1 - \delta_L)^{-1})$.

¹There was a factor of n in the ϕ_i that can be removed.

Algorithm 2 A procedure for solving Problem 2 based on the Width-Free MMWUM method.

procedure SparsifySumOfMatricesByMMWUM($B_1, \dots, B_m, \varepsilon$)

input: Matrices $B_1, \dots, B_m \in \mathbb{S}_+^n$ such that $\sum_i B_i = I$, and a parameter $\varepsilon \in (0, 1)$.

output: A vector y with $O(n \log n / \varepsilon^2)$ nonzero entries such that $I \preceq \sum_i y_i B_i \preceq (1 + O(\varepsilon))I$.

Initially $A(0) := 0$, and $y(0) := 0$. Set parameters

$$u_0 := 0, \quad \ell_0 := 0, \quad \Delta_U := \ln(1 + \delta_U), \quad \Delta_L := \ln((1 - \delta_L)^{-1}),$$

where δ_U, δ_L and T are as defined in (21).

Define the potential functions $\Psi^u(A) := \text{Tr} \exp(-uI + \gamma A)$ and $\Psi_\ell(A) := \text{Tr} \exp(\ell I - \gamma A)$.

For $t = 1, \dots, T$

Set $u_t := u_{t-1} + \Delta_U$ and $\ell_t := \ell_{t-1} + \Delta_L$.

Find a matrix B_j and a value $\alpha > 0$ such that

$$\Psi^{u_t}(A(t-1) + \alpha B_j) \leq \Psi^{u_{t-1}}(A(t-1)) \quad \text{and} \quad \Psi_{\ell_t}(A(t-1) + \alpha B_j) \leq \Psi_{\ell_{t-1}}(A(t-1)).$$

Set $A(t) := A(t-1) + \alpha B_j$ and $y(t) := y(t-1) + \alpha e_j$.

Return $y(T) / \lambda_{\min}(A(T))$.

Proposition 28. *The inequalities in (22) governing the algorithm's change in potentials are equivalent to inequalities in (23).*

$$\begin{aligned} \Psi^{(t+1)\Delta_U}(A(t) + \alpha B_j) &\leq \Psi^{t\Delta_U}(A(t)) \\ \Psi_{(t+1)\Delta_L}(A(t) + \alpha B_j) &\leq \Psi_{t\Delta_L}(A(t)) \end{aligned} \tag{23}$$

Proof. Obviously (22) is equivalent to

$$\begin{aligned} (1 + \delta_U)^{-(t+1)} \cdot \Phi_U(t+1) &\leq (1 + \delta_U)^{-t} \cdot \Phi_U(t) \quad \forall t \in \{0, \dots, T-1\}, \\ (1 - \delta_L)^{-(t+1)} \cdot \Phi_L(t+1) &\leq (1 - \delta_L)^{-t} \cdot \Phi_L(t) \quad \forall t \in \{0, \dots, T-1\}. \end{aligned}$$

By the definition of Φ_U and Φ_L , and by properties of the exponential function, these inequalities are equivalent to

$$\begin{aligned} \text{Tr} \exp(-(t+1)\Delta_U I + \gamma A(t+1)) &\leq \text{Tr} \exp(-t\Delta_U I + \gamma A(t)), \\ \text{Tr} \exp((t+1)\Delta_L I - \gamma A(t+1)) &\leq \text{Tr} \exp(t\Delta_L I - \gamma A(t)). \end{aligned} \tag{24}$$

Writing $A(t+1) = A(t) + \alpha B_j$, these inequalities in (24) are equivalent to (23). \square

Algorithm 2 gives pseudocode for the algorithm of Section 6, using the functions Ψ^u and Ψ_ℓ to control the change in potentials.

The main point of this section is to observe that Algorithms 1 and 2 are identical with the exception of different parameters and different potential functions. We believe that this similarity between these two algorithms is intriguing, especially since the BSS algorithm has been called “highly original” by Naor [28]. In retrospect, it would have been perhaps more natural to develop the BSS algorithm by the following logical progression of ideas: first observe that MMWUM is useful for giving sparse solutions to SDPs, then design Algorithm 2, then later realize that a clever refinement of it leads to Algorithm 1 and its improved analysis. It is remarkable that Batson, Spielman and Srivastava developed their algorithm from first principles, apparently without knowing this connection to established algorithmic techniques.

With the advantage of hindsight (i.e., the knowledge that the BSS algorithm exists), we now explain how one might be tempted to refine Algorithm 2. It is quite tempting to modify the potential functions to more strongly penalize eigenvalues which deviate from the desired range. The natural approach to do this would be to increase the derivatives of the potential function by increasing the parameter γ . However, as remarked at the end of Section 6, the algorithm is actually unaffected by varying γ ! Thus, to improve Algorithm 2, one must seek a more substantially different potential function.

Focusing on the upper potential, we consider the question: is there a function $f: \mathbb{R} \rightarrow \mathbb{R}$ with steeper derivatives than $\exp(u - x)$ and such that, for any matrices A and B , $\text{Tr } f(A + B)$ can be easily related to $\text{Tr } f(A)$? The natural candidates to try are $f(x) = -\log(u - x)$ and $f(x) = (u - x)^{-1}$ since, in both cases, $\text{Tr } f(A + B)$ can be related to $\text{Tr } f(A)$ by the Sherman-Morrison-Woodbury formula. We do not know whether the choice $f(x) = -\log(u - x)$ can be made to work. However, choosing $f(x) = (u - x)^{-1}$, one arrives at Algorithm 1, our generalization of the BSS algorithm. Of course, even after arriving at this algorithm, one must also analyze it, and this requires the delicate calculations that were accomplished by Batson, Spielman and Srivastava.

Acknowledgements

We thank Satyen Kale for helpful discussions.

References

- [1] Rudolf Ahlswede and Andreas Winter. Strong converse for identification via quantum channels. *IEEE Transactions on Information Theory*, 48(3):569–579, March 2002.
- [2] Ingo Althöfer. On sparse approximations to randomized strategies and convex combinations. *Linear Algebra and Applications*, 199:339–355, 1994.
- [3] Sanjeev Arora and Satyen Kale. A combinatorial, primal-dual approach to semidefinite programs. In *Proceedings of the 39th Annual ACM Symposium on Theory of Computing (STOC)*, 2007.
- [4] Joshua Batson, Daniel A. Spielman, and Nikhil Srivastava. Twice-Ramanujan sparsifiers. In *Proceedings of the 41st Annual ACM Symposium on Theory of Computing (STOC)*, 2009. To appear in *SIAM Journal on Scientific Computing*.
- [5] András A. Benczúr and David R. Karger. Approximate s - t min-cuts in $\tilde{O}(n^2)$ time. In *Proceedings of the 28th Annual ACM Symposium on Theory of Computing (STOC)*, 1996.
- [6] András A. Benczúr and David R. Karger. Randomized approximation schemes for cuts and flows in capacitated graphs, 2002. <http://arxiv.org/abs/cs/0207078>.
- [7] Moses Charikar, Chandra Chekuri, Ashish Goel, Sudipto Guha, and Serge A. Plotkin. Approximating a finite metric by a small number of tree metrics. In *Proceedings of the 39th Annual IEEE Symposium on Foundations of Computer Science (FOCS)*, 1998.
- [8] Marcel de Carli Silva and Levent Tunçel. Min-max theorems related to geometric representations of graphs and their SDPs, August 2011. <http://arxiv.org/abs/1010.6036>.
- [9] Yoav Freund and Robert E. Schapire. Adaptive game playing using multiplicative weights. *Games and Economic Behavior*, 29:79–103, 1999.

- [10] Wai Shing Fung, Ramesh Hariharan, Nicholas J. A. Harvey, and Debmalya Panigrahi. A general framework for graph sparsification. In *Proceedings of the 43rd Annual ACM Symposium on Theory of Computing (STOC)*, 2011.
- [11] Wai Shing Fung and Nicholas J. A. Harvey. Graph sparsification by edge-connectivity and random spanning trees, May 2010. <http://arxiv.org/abs/1005.0265>.
- [12] Ashish Goel, Michael Kapralov, and Sanjeev Khanna. Graph sparsification via refinement sampling, April 2010. <http://arxiv.org/abs/1004.4915>.
- [13] William W. Hager. Updating the inverse of a matrix. *SIAM Review*, 31(2):221–239, 1989.
- [14] Ramesh Hariharan and Debmalya Panigrahi. A general framework for graph sparsification, April 2010. <http://arxiv.org/abs/1004.4080>.
- [15] Ramesh Hariharan and Debmalya Panigrahi. A linear-time algorithm for sparsification of unweighted graphs, May 2010. <http://arxiv.org/abs/1005.0670>.
- [16] Nicholas J. A. Harvey. Lecture notes for C&O 750: Randomized algorithms, 2011. <http://www.math.uwaterloo.ca/~harvey/W11/Lecture11Notes.pdf>.
- [17] Thomas Hofmeister and Hanno Lefmann. Computing sparse approximations deterministically. *Linear Algebra and its Applications*, 240:9–19, 1996.
- [18] Roger A. Horn and Charles R. Johnson. *Matrix analysis*. Cambridge University Press, Cambridge, 1990. Corrected reprint of the 1985 original.
- [19] Garud Iyengar, David J. Phillips, and Clifford Stein. Approximation algorithms for semidefinite packing problems with applications to maxcut and graph coloring. In Michael Jünger and Volker Kaibel, editors, *Integer Programming and Combinatorial Optimization*, volume 3509 of *Lecture Notes in Computer Science*, pages 77–90. Springer Berlin / Heidelberg, 2005.
- [20] Garud Iyengar, David J. Phillips, and Clifford Stein. Approximating semidefinite packing programs. *SIAM Journal on Optimization*, 21(1):231–268, 2011.
- [21] Rahul Jain and Penghui Yao. A parallel approximation algorithm for positive semidefinite programming. In *The 52nd Annual IEEE Symposium on Foundations of Computer Science (FOCS 2011)*, 2011. (to appear).
- [22] Satyen Kale. *Efficient Algorithms using the Multiplicative Weights Update Method*. PhD thesis, Princeton University, 2007. Princeton Tech Report TR-804-07.
- [23] Jonathan A. Kelner and Alex Levin. Spectral sparsification in the semi-streaming setting. In *Proceedings of the 28th International Symposium on Theoretical Aspects of Computer Science (STACS)*, pages 440–451, 2011.
- [24] Ioannis Koutis, Gary L. Miller, and Richard Peng. Approaching optimality for solving SDD systems. In *Proceedings of the 51st Annual IEEE Symposium on Foundations of Computer Science (FOCS)*, 2010.
- [25] Richard J. Lipton and Neal E. Young. Simple strategies for large zero-sum games with applications to complexity theory. In *Proceedings of the 26th Annual ACM Symposium on Theory of Computing (STOC)*, 1994.

- [26] L. Lovász. Semidefinite programs and combinatorial optimization.
- [27] R. J. McEliece, E. R. Rodemich, and H. C. Rumsey, Jr. The Lovász bound and some generalizations. *J. Combin. Inform. System Sci.*, 3(3):134–152, 1978.
- [28] Assaf Naor. Sparse quadratic forms and their geometric applications (after batson, spielman and srivastava). In *Séminaire Bourbaki*, 2011. Exposé no. 1033.
- [29] Ilan Newman and Yuri Rabinovich. Finite volume spaces and sparsification, 2010. <http://arxiv.org/abs/1002.3541>.
- [30] Lorenzo Orecchia and Nisheeth K. Vishnoi. Towards an SDP-based approach to spectral methods: A nearly-linear time algorithm for graph partitioning and decomposition. In *Proceedings of the 22nd Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 532–545, 2011.
- [31] Juan A. Rodríguez. On the Laplacian eigenvalues and metric parameters of hypergraphs. *Linear Multilinear Algebra*, 50(1):1–14, 2002.
- [32] Mark Rudelson. Random vectors in the isotropic position. *J. of Functional Analysis*, 164(1):60–72, 1999.
- [33] Mark Rudelson and Roman Vershynin. Sampling from large matrices: An approach through geometric functional analysis. *Journal of the ACM*, 54(4), 2007.
- [34] Gideon Schechtman. Tight embedding of subspaces of L_p in ℓ_p^n for even p . *Proceedings of the AMS*. To appear.
- [35] A. Schrijver. A comparison of the Delsarte and Lovász bounds. *IEEE Trans. Inform. Theory*, 25(4):425–429, 1979.
- [36] Anthony Man-Cho So, Yinyu Ye, and Jiawei Zhang. A unified theorem on SDP rank reduction. *Mathematics of Operations Research*, 33(4):910–920, 2008.
- [37] Daniel A. Spielman and Nikhil Srivastava. An elementary proof of the restricted invertibility theorem. *Israel J. Math.* To appear.
- [38] Daniel A. Spielman and Nikhil Srivastava. Graph sparsification by effective resistances. In *Proceedings of the 40th Annual ACM Symposium on Theory of Computing (STOC)*, pages 563–568, 2008.
- [39] Daniel A. Spielman and Shang-Hua Teng. Nearly-linear time algorithms for graph partitioning, graph sparsification, and solving linear systems. In *Proceedings of the 36th Annual ACM Symposium on Theory of Computing (STOC)*, pages 81–90, 2004.
- [40] Nikhil Srivastava. On contact points of convex bodies, 2009. <http://www.cs.yale.edu/homes/srivastava/papers/contact.pdf>.
- [41] Roman Vershynin. A note on sums of independent random matrices after Ahlswede-Winter, 2008. <http://www-personal.umich.edu/~romanv/teaching/reading-group/ahlsvede-winter.pdf>.
- [42] Avi Wigderson and David Xiao. Derandomizing the Ahlswede-Winter matrix-valued chernoff bound using pessimistic estimators and applications. *Theory of Computing*, 4(3), 2008.
- [43] Neal Young. Greedy algorithms by derandomizing unknown distributions. Technical Report 1087, Department of ORIE, Cornell University, March 1994. <http://hdl.handle.net/1813/8971>.

- [44] Neal Young. Randomized rounding without solving the linear program. In *Proceedings of the 6th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 170–178, 1995.

A Proofs of the Applications

Corollary 4. *Let $G = (V, E)$ be a graph, let $w: E \rightarrow \mathbb{R}_+$ be a weight function, and let $c_1, \dots, c_k: E \rightarrow \mathbb{R}_+$ be cost functions, with $k = O(n)$. Let $\mathcal{L}_G(w)$ denote the Laplacian matrix for graph G with weight function w . For any real $\varepsilon \in (0, 1)$, there is a deterministic polynomial-time algorithm to find a subgraph H of G and a weight function $w_H: E(H) \rightarrow \mathbb{R}_+$ such that*

$$\begin{aligned} \mathcal{L}_G(w) &\preceq \mathcal{L}_H(w_H) \preceq (1 + \varepsilon)\mathcal{L}_G(w), \\ \sum_{e \in E} w_e c_{i,e} &\leq \sum_{e \in E(H)} w_{H,e} c_{i,e} \leq (1 + \varepsilon) \sum_{e \in E} w_e c_{i,e} \quad \text{for all } i \end{aligned}$$

and $|E(H)| = O(n/\varepsilon^2)$.

Proof. For every edge $e = ij \in E$, let B_e be the direct sum $w_{ij}[(e_i - e_j)(e_i - e_j)^T \oplus c_{1,e} \oplus \dots \oplus c_{k,e}]$. Let $B := \mathcal{L}_G(w) \oplus w^T c_1 \oplus \dots \oplus w^T c_k$. The result follows immediately by applying Theorem 3 to these matrices. \square

Corollary 5. *Let $G = (V, E)$ be a graph and let $w: E \rightarrow \mathbb{R}_+$ be a weight function. Let E_1, \dots, E_k be a partition of the edges, i.e., each edge is colored with one of k colors. For any real $\varepsilon \in (0, 1)$, there is a deterministic polynomial-time algorithm to find a subgraph H of G and a weight function $w_H: E(H) \rightarrow \mathbb{R}_+$ such that*

$$\begin{aligned} \mathcal{L}_G(w) &\preceq \mathcal{L}_H(w_H) \preceq (1 + \varepsilon)\mathcal{L}_G(w), \\ (1 - \varepsilon) \sum_{e \in E_i} w_e &\leq \sum_{e \in E(H) \cap E_i} w_{H,e} \leq (1 + \varepsilon) \sum_{e \in E_i} w_e \quad \text{for all } i, \end{aligned}$$

and $|E(H)| = O((n + k)/\varepsilon^2)$.

Proof. For each i , let $c_i: E \rightarrow \mathbb{R}$ be the characteristic vector of E_i . Now apply Corollary 4. \square

Corollary 6 (Spectral sparsifiers for hypergraphs). *For any real $\varepsilon \in (0, 1)$, there is a deterministic polynomial-time algorithm to find a sub-hypergraph \mathcal{G} of \mathcal{H} and a weight function $w_{\mathcal{G}}: \mathcal{E}(\mathcal{G}) \rightarrow \mathbb{R}_+$ such that*

$$\mathcal{L}_{\mathcal{H}}(w) \preceq \mathcal{L}_{\mathcal{G}}(w_{\mathcal{G}}) \preceq (1 + \varepsilon)\mathcal{L}_{\mathcal{H}}(w),$$

and $|\mathcal{E}(\mathcal{G})| = O(n/\varepsilon^2)$.

Proof. The result follows directly by applying Theorem 3 to the matrices $w_E \mathcal{L}_E$. \square

Corollary 7 (Cut sparsifiers for hypergraphs, second definition). *For any real $\varepsilon \in (0, 1)$, there is a deterministic polynomial-time algorithm to find a sub-hypergraph \mathcal{G} of \mathcal{H} and a weight function $w_{\mathcal{G}}: \mathcal{E}(\mathcal{G}) \rightarrow \mathbb{R}_+$ such that*

$$w^*(\delta_{\mathcal{H}}(S)) \leq w_{\mathcal{G}}^*(\delta_{\mathcal{G}}(S)) \leq (1 + \varepsilon)w^*(\delta_{\mathcal{H}}(S)) \quad \text{for every } S \subseteq V,$$

and $|\mathcal{E}(\mathcal{G})| = O(n/\varepsilon^2)$.

Proof. Note that $w^*(\delta_{\mathcal{H}}(S))$ is obtained by evaluating the quadratic form $x^T \mathcal{L}_{\mathcal{H}}(w)x$, where x is the characteristic vector of S . Thus the sparsifier produced by Corollary 6 satisfies the desired inequalities. \square

Corollary 8 (Cut sparsifiers for hypergraphs, first definition). *Assume that \mathcal{H} is an r -uniform hypergraph. For any real $\varepsilon \in (0, 1)$, there is a deterministic polynomial-time algorithm to find a sub-hypergraph \mathcal{G} of \mathcal{H} and a weight function $w_{\mathcal{G}}: \mathcal{E}(\mathcal{G}) \rightarrow \mathbb{R}_+$ such that*

$$\frac{(r-1)}{r^2/4}w(\delta_{\mathcal{H}}(S)) \leq w_{\mathcal{G}}(\delta_{\mathcal{G}}(S)) \leq \frac{(1+\varepsilon)r^2}{4(r-1)}w(\delta_{\mathcal{H}}(S)) \quad \forall S \subseteq V,$$

and $|\mathcal{E}(\mathcal{G})| = O(n/\varepsilon^2)$. In other words, the sparsified hypergraph \mathcal{G} approximates the weight of the cuts in the hypergraph \mathcal{H} to within a factor $\Theta(r^2)$.

Proof. For any r -uniform hypergraph \mathcal{H} , it is easy to see that

$$(r-1)w(\delta_{\mathcal{H}}(S)) \leq w^*(\delta_{\mathcal{H}}(S)) \leq \lceil r/2 \rceil \lceil r/2 \rceil w(\delta_{\mathcal{H}}(S)) \quad \forall S \subseteq V. \quad (25)$$

Thus the sparsifier produced by Corollary 6 satisfies the desired inequalities. \square

Corollary 9 (Cut sparsifiers for 3-uniform hypergraphs). *Assume that \mathcal{H} is a 3-uniform hypergraph. For any $\varepsilon \in (0, 1)$, there is a deterministic polynomial-time algorithm to find a sub-hypergraph \mathcal{G} of \mathcal{H} and a weight function $w_{\mathcal{G}}: \mathcal{E}(\mathcal{G}) \rightarrow \mathbb{R}_+$ such that*

$$w(\delta_{\mathcal{H}}(S)) \leq w_{\mathcal{G}}(\delta_{\mathcal{G}}(S)) \leq (1+\varepsilon)w(\delta_{\mathcal{H}}(S)) \quad \forall S \subseteq V,$$

and $|\mathcal{E}(\mathcal{G})| = O(n/\varepsilon^2)$.

Proof. Since $r = 3$, a consequence of (25) is that $w^*(\delta_{\mathcal{H}}(S)) = 2w(\delta_{\mathcal{H}}(S))$ for every S . Thus the sparsifier produced by Corollary 6 satisfies the desired inequalities. \square

Corollary 10. *Let A_1, \dots, A_m be symmetric, positive semidefinite matrices of size $n \times n$, and let B be a symmetric matrix of size $n \times n$. Let $c \in \mathbb{R}^m$ with $c \geq 0$. Suppose that the semidefinite program (SDP)*

$$\min \left\{ c^T z : \sum_i z_i A_i \succeq B, z \in \mathbb{R}^m, z \geq 0 \right\}$$

has a feasible solution z^* . Then, for any real $\varepsilon \in (0, 1)$, it has a feasible solution \bar{z} with at most $O(n/\varepsilon^2)$ nonzero entries and $c^T \bar{z} \leq (1+\varepsilon)c^T z^*$.

Proof. Let $B'_i := \begin{bmatrix} z_i^* A_i & \mathbf{0} \\ \mathbf{0} & c_i z_i^* \end{bmatrix}$ for every $i \in [m]$ and $B' := \begin{bmatrix} D & \mathbf{0} \\ \mathbf{0} & c^T z^* \end{bmatrix}$, where $D := \sum_i z_i^* A_i \succeq B$. Then $B'_i \succeq 0$ and $B' = \sum_i B'_i$. By applying Theorem 3, we obtain $y \in \mathbb{R}^m$ with $y \geq 0$ and $O(n/\varepsilon^2)$ nonzero entries such that $\sum_i y_i z_i^* A_i \succeq D \succeq B$ and $\sum_i c_i y_i z_i^* \leq (1+\varepsilon)c^T z^*$. Thus, we can take $\bar{z}_i = y_i z_i^*$ for every $i \in [m]$. \square

Corollary 11. *Let $G = (V, E)$ be a graph. For any real $\varepsilon \in (0, 1)$, there is a deterministic polynomial-time algorithm to find a subgraph H of G such that*

$$(1-\varepsilon)t'(G) \leq t'(H) \leq t'(G)$$

and $|E(H)| = O(n/\varepsilon^2)$.

Proof. It is straightforward to formulate $t'(G)$ as an SDP (see, e.g., [26]) so that its dual has an optimal solution and there is no duality gap. The dual can be written as:

$$\begin{aligned} \max \quad & \sum_{e \in E} z_e \\ \text{Diag}(y^*) \succeq & \mathcal{L}_G(z) \\ \sum_{v \in V} y_v = & 1 \\ z \geq & 0 \end{aligned} \tag{26}$$

The proof is now almost identical to the proof of Corollary 10. Let (z^*, y^*) be an optimal solution. Using Theorem 3, we obtain $\bar{z} \in \mathbb{R}^E$ with $\bar{z} \geq 0$ and $O(n/\varepsilon^2)$ nonzero entries such that (y^*, \bar{z}) is feasible in (26) and has objective value $\sum_{e \in E(H)} \bar{z}_e \geq (1 - \varepsilon)t(G)$, where $H = (V, E(H))$ and $E(H)$ is the support of \bar{z} . Then \bar{z} is also feasible for the SDP defined using H instead of G , which shows that $t'(H) \geq (1 - \varepsilon)t'(G)$. \square

Corollary 12. *Let $G = (V, E)$ be a graph. For any real $\varepsilon \in (0, 1)$, there is a deterministic polynomial-time algorithm to find a supergraph H of G such that*

$$\frac{\vartheta'(G)}{1 - \varepsilon + \varepsilon\vartheta'(G)} \leq \vartheta'(H) \leq \vartheta'(G)$$

and $|E(H)| = \binom{n}{2} - O(n/\varepsilon^2)$.

Proof. For a graph $G = (V, E)$, define $t(G)$ as the square of the minimum radius of a hypersphere on \mathbb{R}^n such that there is a map from V to the hypersphere such that adjacent vertices are mapped to points at distance exactly 1. Lovász [26] noted that $t(G)$ is related to the Lovász theta number $\vartheta(\bar{G})$ of the complement \bar{G} of G by the formula $2t(G) + 1/\vartheta(\bar{G}) = 1$; see [8] for a proof. By repeating the same proof for $t'(G)$, one finds that $2t'(G) + 1/\vartheta'(\bar{G}) = 1$. The result now follows from Corollary 11 via this formula. \square

Corollary 13. *Let G be a graph such that $\vartheta'(G) = o(\sqrt{n})$. For any real $\gamma > 0$, there is a supergraph H of G such that*

$$\frac{\vartheta'(G)}{1 + \gamma} \leq \vartheta'(H) \leq \vartheta'(G)$$

and $|E(H)| = \binom{n}{2} - O(n\vartheta(G)^2/\gamma^2)$.

Proof. Apply Corollary 12 with $\varepsilon := \gamma/\vartheta'(G)$. \square

Corollary 14. *Let G be a graph such that $\vartheta'(G) = \Omega(\sqrt{n})$. For any real $\gamma \geq 1$, there is a supergraph H of G such that*

$$\vartheta'(H) = \Omega(\sqrt{n}/\gamma)$$

and $|E(H)| = \binom{n}{2} - O(n^2/\gamma^2)$.

Proof. Apply Corollary 12 with $\varepsilon := \gamma/\sqrt{n}$. \square

Corollary 16. *Let B_1, \dots, B_m be symmetric, positive semidefinite matrices of size $n \times n$ and let $\lambda \in \mathbb{R}^m$ satisfy $\lambda \geq 0$ and $\sum_i \lambda_i = 1$. Let $B = \sum_i \lambda_i B_i$. For any $\varepsilon \in (0, 1)$, there exists $\mu \geq 0$ with $\sum_i \mu_i = 1$ such that μ has $O(n/\varepsilon^2)$ nonzero entries and*

$$(1 - \varepsilon)B \preceq \sum_i \mu_i B_i \preceq (1 + \varepsilon)B.$$

Proof. Let $B'_i := \begin{bmatrix} \lambda_i B_i & \mathbf{0} \\ \mathbf{0} & \lambda_i \end{bmatrix}$ for every $i \in [m]$ and $B' := \begin{bmatrix} B & \mathbf{0} \\ \mathbf{0} & 1 \end{bmatrix}$, so that $B'_i \succeq 0$ and $B' = \sum_i B'_i$. By applying Theorem 3, we obtain $y \in \mathbb{R}^m$ with $y \geq 0$ and $O(n/\varepsilon^2)$ nonzero entries such that $B' \preceq \sum_i y_i B'_i \preceq (1 + \varepsilon)B'$ or, equivalently, $B \preceq \sum_i y_i \lambda_i B_i \preceq (1 + \varepsilon)B$ and $1 \leq \sum_i y_i \lambda_i \leq 1 + \varepsilon$. Let $\mu \in \mathbb{R}^m$ be defined by $\mu_i := y_i \lambda_i / (\sum_i y_i \lambda_i)$. Then $\mu \geq 0$ and $\sum_i \mu_i = 1$, and

$$(1 - \varepsilon)B \preceq \frac{B}{1 + \varepsilon} \preceq \frac{B}{\sum_i y_i \lambda_i} \preceq \sum_i \mu_i B_i \preceq \frac{1 + \varepsilon}{\sum_i y_i \lambda_i} B \preceq (1 + \varepsilon)B.$$

This completes the proof. \square

Corollary 17. *Let $G = (V, E)$ be a graph, let $w: E \rightarrow \mathbb{R}_+$ be a weight function, and let \mathcal{F} be a collection of subgraphs of G such that $\sum_{F \in \mathcal{F}} |V(F)| = O(n)$. For any real $\varepsilon \in (0, 1)$, there is a deterministic polynomial-time algorithm to find a subgraph H of G and a weight function $w_H: E(H) \rightarrow \mathbb{R}_+$ such that $|E(H)| = O(n/\varepsilon^2)$ and*

$$\begin{aligned} \mathcal{L}_G(w) &\preceq \mathcal{L}_H(w_H) \preceq (1 + \varepsilon)\mathcal{L}_G(w), \\ \mathcal{L}_F(w_F) &\preceq \mathcal{L}_{H \cap F}(w_H \upharpoonright_{E(H \cap F)}) \preceq (1 + \varepsilon)\mathcal{L}_F(w_F) \quad \text{for all } F \in \mathcal{F}, \end{aligned}$$

where $w_F := w \upharpoonright_{E(F)}$ is the restriction of w to the coordinates $E(F)$ and $H \cap F = (V(F), E(F) \cap E(H))$.

Proof. For each edge $e \in E$, define $B_e := w_e [\mathcal{L}_G(\chi^e) \oplus \bigoplus_{F \in \mathcal{F}} \mathcal{L}_F(\chi^e \upharpoonright_{E(F)})]$, where χ^e denotes the characteristic vector of $\{e\}$ as a subset of E . Now apply Theorem 3. \square

B The MMWUM

In this section we provide some proofs about the MMWUM. These proofs are due to Kale [22]. Our set up and conclusions are slightly different and we modified the proofs accordingly. We reproduce the proofs here for the sake of completeness.

Theorem 22 can be viewed as a block-friendly version of MMWUM. First we show the version with only one block. It is basically the same as [22, Theorem 13 in Chapter 4].

Theorem 29. *Let T be a positive integer. Let $C, A_1, \dots, A_m \in \mathbb{S}^n$. Let $\eta > 0$ and $0 < \beta \leq 1/2$. For any given $X \in \mathbb{S}^n$, consider the system*

$$\sum_{i=1}^m y_i \langle A_i, X \rangle \geq \langle C, X \rangle - \eta \operatorname{Tr} X, \quad \text{and} \quad y \in \mathbb{R}_+^m. \quad (27)$$

Let $\{\mathcal{P}, \mathcal{N}\}$ be a partition of $[T]$, let $0 < \ell \leq \rho$, and let $W^{(t)} \in \mathbb{S}^n$ and $\ell^{(t)} \in \mathbb{R}$ for $t \in [T + 1]$. Let $y^{(t)} \in \mathbb{R}^m$ for $t \in [T]$. Suppose the following properties hold:

$$\begin{aligned} W^{(t+1)} &= \exp \left(-\frac{\beta}{\ell + \rho} \sum_{\tau=1}^t \left[\sum_{i=1}^m y_i^{(\tau)} A_i - C + \ell^{(\tau)} I \right] \right), \quad \forall t \in \{0, \dots, T\}, \\ y &= y^{(t)} \text{ is a solution for (27) with } X = W^{(t)}, \quad \forall t \in [T], \\ \sum_{i=1}^m y_i^{(t)} A_i - C &\in \begin{cases} [-\ell, \rho], & \text{if } t \in \mathcal{P}, \\ [-\rho, \ell], & \text{if } t \in \mathcal{N}, \end{cases} \quad \forall t \in [T], \\ \ell^{(t)} &= \ell, \quad \forall t \in \mathcal{P}, \quad \text{and} \quad \ell^{(t)} = -\ell, \quad \forall t \in \mathcal{N}. \end{aligned}$$

Define $\bar{y} := \frac{1}{T} \sum_{t=1}^T y^{(t)}$. Then

$$\sum_{i=1}^m \bar{y}_i A_i - C \succeq -\left[\beta\ell + \frac{(\rho + \ell) \ln n}{T\beta} + (1 + \beta)\eta\right] I. \quad (28)$$

The main tool for the proof of Theorem 29 is the following result:

Theorem 30 (Kale [22, Corollary 3 in Chapter 3]). *Let $0 < \beta \leq 1/2$. Let T be a positive integer. Let $\{\mathcal{P}, \mathcal{N}\}$ be a partition of $[T]$, and let $M^{(t)} \in \mathbb{S}^n$ for $t \in [T]$ and $W^{(t)} \in \mathbb{S}^n$ for $t \in [T + 1]$ with the following properties:*

$$\begin{aligned} W^{(t+1)} &= \exp\left(-\beta \sum_{\tau=1}^t M^{(\tau)}\right) \quad \forall t = 0, \dots, T, \\ 0 \preceq M^{(t)} \preceq I, \quad \forall t \in \mathcal{P}, \quad \text{and} \quad -I \preceq M^{(t)} \preceq 0, \quad \forall t \in \mathcal{N}, \end{aligned}$$

Let

$$P^{(t)} := \frac{1}{\text{Tr} W^{(t)}} W^{(t)}, \quad \forall t \in [T].$$

Then

$$(1 - \beta) \sum_{t \in \mathcal{P}} \langle M^{(t)}, P^{(t)} \rangle + (1 + \beta) \sum_{t \in \mathcal{N}} \langle M^{(t)}, P^{(t)} \rangle \leq \lambda_{\min}\left(\sum_{t=1}^T M^{(t)}\right) + \frac{\ln n}{\beta}. \quad (29)$$

Proof. Set $\Phi^{(t)} := \text{Tr}(W^{(t)})$ for $t \in [T + 1]$. Put $\beta_1 := 1 - e^{-\beta}$ and $\beta_2 := e^{\beta} - 1$. Then, for any $t \in [T]$,

$$\begin{aligned} \Phi^{(t+1)} &= \text{Tr}(W^{(t+1)}) = \text{Tr}\left(\exp\left(-\beta \sum_{\tau=1}^t M^{(\tau)}\right)\right) \\ &\leq \text{Tr}\left(\exp\left(-\beta \sum_{\tau=1}^{t-1} M^{(\tau)}\right) \exp\left(-\beta M^{(t)}\right)\right) = \text{Tr}\left(W^{(t)} \exp(-\beta M^{(t)})\right) \\ &= \langle W^{(t)}, \exp(-\beta M^{(t)}) \rangle, \end{aligned}$$

where we have used Golden-Thompson's inequality (14).

Using the fact that e^x is convex, one can prove that

$$\begin{aligned} 0 \preceq A \preceq I &\implies \exp(-\beta A) \preceq I - \beta_1 A, \\ -I \preceq A \preceq 0 &\implies \exp(-\beta A) \preceq I - \beta_2 A. \end{aligned}$$

Suppose that $t \in \mathcal{P}$. Then $\exp(-\beta M^{(t)}) \preceq I - \beta_1 M^{(t)}$, and since $W^{(t)} \succeq 0$, we get

$$\begin{aligned} \Phi^{(t+1)} &\leq \langle W^{(t)}, \exp(-\beta M^{(t)}) \rangle \leq \langle W^{(t)}, I - \beta_1 M^{(t)} \rangle \\ &= \text{Tr}(W^{(t)}) - \beta_1 \langle W^{(t)}, M^{(t)} \rangle \\ &= \text{Tr}(W^{(t)}) - \text{Tr}(W^{(t)}) \beta_1 \langle P^{(t)}, M^{(t)} \rangle \\ &= \text{Tr}(W^{(t)}) \left[1 - \beta_1 \langle P^{(t)}, M^{(t)} \rangle\right] \\ &= \Phi^{(t)} \left[1 - \beta_1 \langle P^{(t)}, M^{(t)} \rangle\right] \\ &\leq \Phi^{(t)} \exp(-\beta_1 \langle P^{(t)}, M^{(t)} \rangle). \end{aligned}$$

Similarly, if $t \in \mathcal{N}$, then

$$\Phi^{(t+1)} \leq \Phi^{(t)} \exp(-\beta_2 \langle P^{(t)}, M^{(t)} \rangle).$$

By induction on t , and using $\Phi^{(1)} = \text{Tr}(I) = n$, we get

$$\Phi^{(t+1)} \leq n \exp \left(-\beta_1 \sum_{\tau \in \mathcal{P} \cap [t]} \langle M^{(\tau)}, P^{(\tau)} \rangle - \beta_2 \sum_{\tau \in \mathcal{N} \cap [t]} \langle M^{(\tau)}, P^{(\tau)} \rangle \right), \quad \forall t \in [T].$$

For every $A \in \mathbb{S}^n$, we have $\text{Tr}(\exp(A)) = \sum_{i=1}^n e^{\lambda_i} \geq e^{\lambda_j}$ for any $j \in [n]$, where $\lambda_1, \dots, \lambda_n$ are the eigenvalues of A . Thus,

$$\begin{aligned} \Phi^{(T+1)} &= \text{Tr}(W^{(T+1)}) = \text{Tr} \left(\exp \left(-\beta \sum_{t=1}^T M^{(t)} \right) \right) \\ &\geq \exp \left(\lambda_{\max} \left(-\beta \sum_{t=1}^T M^{(t)} \right) \right) = \exp \left(-\beta \lambda_{\min} \left(\sum_{t=1}^T M^{(t)} \right) \right). \end{aligned}$$

Thus,

$$\exp \left[-\beta \lambda_{\min} \left(\sum_{t=1}^T M^{(t)} \right) \right] \leq n \exp \left[-\beta_1 \sum_{t \in \mathcal{P}} \langle M^{(t)}, P^{(t)} \rangle - \beta_2 \sum_{t \in \mathcal{N}} \langle M^{(t)}, P^{(t)} \rangle \right].$$

By taking $\ln(\cdot)$ on both sides, we get

$$-\beta \lambda_{\min} \left(\sum_{t=1}^T M^{(t)} \right) \leq \ln n - \left[\beta_1 \sum_{t \in \mathcal{P}} \langle M^{(t)}, P^{(t)} \rangle + \beta_2 \sum_{t \in \mathcal{N}} \langle M^{(t)}, P^{(t)} \rangle \right],$$

so

$$\beta_1 \sum_{t \in \mathcal{P}} \langle M^{(t)}, P^{(t)} \rangle + \beta_2 \sum_{t \in \mathcal{N}} \langle M^{(t)}, P^{(t)} \rangle \leq \beta \lambda_{\min} \left(\sum_{t=1}^T M^{(t)} \right) + \ln n,$$

and

$$\frac{\beta_1}{\beta} \sum_{t \in \mathcal{P}} \langle M^{(t)}, P^{(t)} \rangle + \frac{\beta_2}{\beta} \sum_{t \in \mathcal{N}} \langle M^{(t)}, P^{(t)} \rangle \leq \lambda_{\min} \left(\sum_{t=1}^T M^{(t)} \right) + \frac{\ln n}{\beta}.$$

Since $\sum_{t \in \mathcal{P}} \langle M^{(t)}, P^{(t)} \rangle \geq 0$ and $\sum_{t \in \mathcal{N}} \langle M^{(t)}, P^{(t)} \rangle \leq 0$, to prove (29) it suffices to show that $1 - \beta \leq \beta_1/\beta$ and $1 + \beta \geq \beta_2/\beta$. It is not hard to prove that

$$1 - e^{-x} \geq x(1 - x), \quad \forall x \in [0, +\infty) \quad \text{and} \quad e^x - 1 \leq x(1 + x), \quad \forall x \in [0, \frac{1}{2}]$$

So our choice of β_1 and β_2 ensures that $1 - \beta \leq \beta_1/\beta$ and $1 + \beta \geq \beta_2/\beta$. \square

We can now show the proof of Theorem 29.

Proof of Theorem 29. Let $M^{(t)} := \frac{1}{\ell + \rho} \left[\sum_{i=1}^m y_i^{(t)} A_i - C + \ell^{(t)} I \right]$ and $P^{(t)} := W^{(t)} / \text{Tr} W^{(t)}$ for every t . For every $t \leq T$, using (27),

$$\begin{aligned} \langle M^{(t)}, P^{(t)} \rangle &= \frac{1}{\ell + \rho} \left[\sum_{i=1}^m y_i^{(t)} \langle A_i, P^{(t)} \rangle - \langle C, P^{(t)} \rangle + \ell^{(t)} \langle I, P^{(t)} \rangle \right] \\ &= \frac{1}{(\ell + \rho) \text{Tr} W^{(t)}} \left[\sum_{i=1}^m y_i^{(t)} \langle A_i, W^{(t)} \rangle - \langle C, W^{(t)} \rangle \right] + \frac{\ell^{(t)}}{\ell + \rho} \geq -\frac{\eta}{\ell + \rho} + \frac{\ell^{(t)}}{\ell + \rho}, \end{aligned}$$

since $y^{(t)}$ is a solution for (27) with $X := W^{(t)}$. Thus, by (29),

$$\begin{aligned} \sum_{t \in \mathcal{P}} \frac{(1 - \beta)(\ell^{(t)} - \eta)}{\ell + \rho} + \sum_{t \in \mathcal{N}} \frac{(1 + \beta)(\ell^{(t)} - \eta)}{\ell + \rho} \\ \leq \frac{1}{\rho + \ell} \lambda_{\min} \left(\sum_{t=1}^T \left[\left(\sum_{i=1}^m y_i^{(t)} A_i \right) - C + \ell^{(t)} I \right] \right) + \frac{\ln n}{\beta}. \end{aligned}$$

Multiply through by $\ell + \rho$ and move $\ell^{(t)} I$ out of $\lambda_{\min}(\cdot)$:

$$\begin{aligned} \sum_{t \in \mathcal{P}} (1 - \beta)\ell^{(t)} + \sum_{t \in \mathcal{N}} (1 + \beta)\ell^{(t)} - T(1 + \beta)\eta \\ \leq \lambda_{\min} \left(\sum_{t=1}^T \left[\left(\sum_{i=1}^m y_i^{(t)} A_i \right) - C \right] \right) + \left(\sum_{t=1}^T \ell^{(t)} \right) + \frac{(\rho + \ell) \ln n}{\beta}. \end{aligned}$$

Thus,

$$\sum_{t \in \mathcal{P}} -\beta\ell^{(t)} + \sum_{t \in \mathcal{N}} \beta\ell^{(t)} \leq \lambda_{\min} \left(\sum_{t=1}^T \left[\left(\sum_{i=1}^m y_i^{(t)} A_i \right) - C \right] \right) + \frac{(\rho + \ell) \ln n}{\beta} + T(1 + \beta)\eta.$$

Next note that $\sum_{t \in \mathcal{P}} -\ell^{(t)} + \sum_{t \in \mathcal{N}} \ell^{(t)} = \sum_{t \in \mathcal{P}} -\ell + \sum_{t \in \mathcal{N}} -\ell = -T\ell$, so

$$0 \leq \lambda_{\min} \left(\sum_{t=1}^T \left[\left(\sum_{i=1}^m y_i^{(t)} A_i \right) - C \right] \right) + \beta T\ell + \frac{(\rho + \ell) \ln n}{\beta} + T(1 + \beta)\eta.$$

and

$$0 \leq \lambda_{\min} \left(\frac{1}{T} \sum_{t=1}^T \left[\left(\sum_{i=1}^m y_i^{(t)} A_i \right) - C \right] \right) + \beta\ell + \frac{(\rho + \ell) \ln n}{T\beta} + (1 + \beta)\eta.$$

Thus,

$$\sum_{i=1}^m \bar{y}_i A_i - C = \frac{1}{T} \sum_{t=1}^T \left[\left(\sum_{i=1}^m y_i^{(t)} A_i \right) - C \right] \succeq - \left[\beta\ell + \frac{(\rho + \ell) \ln n}{T\beta} + (1 + \beta)\eta \right] I.$$

□

Theorem 22 can be easily proved from Theorem 29. First, we apply Theorem 29 separately for each block. In each iteration, $y^{(t)}$ is a solution for (27) for all blocks simultaneously, and so the conclusion in (28) holds for all blocks with same \bar{y} . This new algorithm can be seen as equivalent to running K copies of MMWUM, each with different input data, with the caveat that all copies run for the same number of iterations and the vector $y^{(t)}$ returned from the oracle is the same for all copies at each iteration t .

C Optimality of MMWUM Oracle

Proposition 24. *Any oracle for satisfying (9) must have $\rho = \Omega(n/\eta)$, even if the B_i matrices have rank one, and even if X_1 is a scalar multiple of X_2^{-1} .*

Proof. Let $k = n/3$, let I_k be the identity of size $k \times k$, and let $e_j \in \mathbb{R}^k$ be the j th standard basis vector. Let $\zeta = 3\eta$ and define

$$X_1 = \text{diag}(1, \zeta^3, \zeta) \otimes I_k, \quad X_2 = \text{diag}(1, 1/\zeta^3, 1/\zeta) \otimes I_k,$$

where \otimes denotes tensor product. For $j = 1, \dots, k$, define

$$v_{1,j} = [1/\sqrt{2}, -1/\sqrt{2}, 0] \otimes e_j, \quad v_{2,j} = [1/\sqrt{2}, 1/\sqrt{2}, 0] \otimes e_j, \quad v_{3,j} = [0, 0, 1] \otimes e_j.$$

Let $B_{i,j} = v_{i,j}v_{i,j}^T$. Note that $\sum_{i,j} B_{i,j} = I$.

The oracle cannot choose a matrix $B_{i,j}$ with $i \in \{1, 2\}$, since satisfying (9) would lead to a contradiction:

$$\begin{aligned} \frac{\langle X_2, B_i \rangle}{\text{Tr}(X_2)(1+\eta)} &\leq \frac{1}{\alpha} \leq \frac{\langle X_1, B_i \rangle}{\text{Tr}(X_1)(1-\eta)} \\ \implies 1 + 3\eta = 1 + \zeta &< \frac{\langle X_2, B_i \rangle / \text{Tr} X_2}{\langle X_1, B_i \rangle / \text{Tr} X_1} \leq \frac{1+\eta}{1-\eta} < 1 + 3\eta, \end{aligned}$$

for sufficiently small η .

So the oracle must choose a matrix $B_{i,j}$ with $i = 3$. In this case,

$$\begin{aligned} \frac{\text{Tr} B_{i,j}}{\rho} &\leq \frac{1}{\alpha} \leq \frac{\langle X_1, B_{i,j} \rangle}{\text{Tr}(X_1)(1-\eta)} \\ \implies \frac{n}{9\eta} = \frac{n}{3\zeta} &\leq \frac{(1 + \zeta^3 + \zeta)k}{\zeta} = \frac{\text{Tr}(B_{i,j}) \text{Tr}(X_1)}{\langle X_1, B_{i,j} \rangle} \leq \frac{\rho}{1-\eta}. \end{aligned}$$

This shows that $\rho = \Omega(n/\eta)$. □

D The positive semidefiniteness assumption

Proposition 31. *For every positive integer n , there exist matrices $B_1, \dots, B_m \in \mathbb{S}^n$ with $m = \Omega(n^2)$ such that $B := \sum_i B_i$ is positive definite and with the following property: for every $\varepsilon \in (0, 1)$ and $y \in \mathbb{R}^m$ such that $(1 - \varepsilon)B \preceq \sum_i y_i B_i$, all entries of y are nonzero.*

Proof. Let $\mathcal{P} := \{(i, j) : i, j \in [n], i < j\}$. For $(i, j) \in \mathcal{P}$, let $E_{ij} := e_i e_j^T + e_j e_i^T$. Let J denote the matrix of all ones. Then $2I + \sum_{(i,j) \in \mathcal{P}} E_{ij} = I + J =: B \succ 0$. Let $\varepsilon \in (0, 1)$ and suppose that $(1 - \varepsilon)B \preceq 2tI + \sum_{(i,j) \in \mathcal{P}} z_{ij} E_{ij}$ for some $t \in \mathbb{R}$ and $z \in \mathbb{R}^{\mathcal{P}}$. By taking the inner product with E_{ab} on both sides, we see that $0 < 2(1 - \varepsilon) \leq z_{ab}$ for every $(a, b) \in \mathcal{P}$. Similarly, we find that $0 < 2n(1 - \varepsilon) \leq 2nt$. □