# A Matrix Hyperbolic Cosine Algorithm and Applications

Anastasios Zouzias

Department of Computer Science
University of Toronto, Canada

**Abstract.** In this paper, we generalize Spencer's hyperbolic cosine algorithm to the matrix-valued setting. We apply the proposed algorithm to several problems by analyzing its computational efficiency under two special cases of matrices; one in which the matrices have a group structure and an other in which they have rank-one. As an application of the former case, we present a deterministic algorithm that, given the multiplication table of a finite group of size $n$, it constructs an expanding Cayley graph of logarithmic degree in near-optimal $\mathcal{O}(n^2 \log^3 n)$ time. For the latter case, we present a fast deterministic algorithm for spectral sparsification of positive semi-definite matrices, which implies an improved deterministic algorithm for spectral graph sparsification of dense graphs. In addition, we give an elementary connection between spectral sparsification of positive semi-definite matrices and element-wise matrix sparsification. As a consequence, we obtain improved element-wise sparsification algorithms for diagonally dominant-like matrices.

## 1 Introduction

A non-trivial generalization of Chernoff bound type inequalities for matrix-valued random variables was introduced by Ahlswede and Winter [2]. In parallel, Vershynin and Rudelson introduced similar matrix-valued concentration inequalities using different machinery [30,31]. Following these two seminal papers, many variants have been proposed in the literature [29]; see [42] for more. Such inequalities, similarly to their real-valued ancestors, provide powerful tools to analyze probabilistic constructions and the performance of randomized algorithms. There is a rapidly growing line of research exploiting the power of these inequalities including new proofs of probabilistic constructions of expander graphs [3,24,26], matrix approximation by element-wise sparsification [14], graph approximation via edge sparsification [38], analysis of algorithms for matrix completion and decomposition of low rank matrices [29,25], semi-definite relaxation and rounding of quadratic maximization problems [32].

In many settings, it is desirable to convert the above probabilistic proofs into *efficient* deterministic procedures. That is, to derandomize the proofs. Wigderson and Xiao presented an efficient derandomization of the matrix Chernoff bound by generalizing Raghavan's method of pessimistic estimators to the matrix-valued setting [44]. In this paper, we generalize Spencer's hyperbolic cosine algorithm to the matrix-valued setting [33]. In an earlier, preliminary version of our paper [45] the generalization of Spencer's hyperbolic cosine algorithm was also based on the method of pessimistic estimators. However, here we present a proof which is based on a simple averaging argument. Next, we carefully analyze two special cases of matrices; one in which the matrices have a group structure and the other in which they have rank-one. We apply our main result to the following problems: deterministically constructing Alon-Roichman expanding Cayley graphs, approximating graphs via edge sparsification and approximating matrices via element-wise sparsification.

The Alon-Roichman theorem asserts that Cayley graphs obtained by choosing a logarithmic number of group elements independently and uniformly at random are expanders [3]. The original proof of Alon and Roichman is based on Wigner's trace method, whereas recent proofs rely on matrix-valued deviation bounds [24]. Wigderson and Xiao's derandomization of the matrix Chernoff bound implies a deterministic $\mathcal{O}(n^4 \log n)$ time algorithm for constructing Alon-Roichman graphs. Independently, Arora and Kale generalized the multiplicative weights update (MWU) method to the matrix-valued setting and, among other interesting implications, they improved the running time to $\mathcal{O}(n^3 \text{polylog}\,(n))$ [22]. Here we further improve the running time to $\mathcal{O}(n^2 \log^3 n)$ by exploiting the group structure of the problem. In addition, our algorithm is combinatorial in the sense that it only requires counting the number of all closed (even) paths of size at

most $\mathcal{O}(\log n)$ in Cayley graphs. All previous algorithms involve numerical matrix computations such as eigenvalue decompositions and matrix exponentiation.

The second problem that we study is the graph sparsification problem. This problem poses the question whether any dense graph can be approximated by a sparse graph under different notions of approximation. Given any undirected graph, the most well-studied notions of approximation by a sparse graph include approximating, *all* pairwise distances up to an additive error [28], every cut to an arbitrarily small multiplicative error [8] and every eigenvalue of the difference of their Laplacian matrices to an arbitrarily small relative error [37]; the resulting graphs are usually called *graph spanners*, *cut sparsifiers* and *spectral sparsifiers*, respectively. Given that the notion of spectral sparsification is stronger than cut sparsification, so we focus on spectral sparsifiers. An efficient randomized algorithm to construct an $(1 + \varepsilon)$-spectral sparsifier with $\mathcal{O}(n \log n / \varepsilon^2)$ edges was given in [38]. Furthermore, an $(1 + \varepsilon)$-spectral sparsifier with $\mathcal{O}(n / \varepsilon^2)$ edges can be computed in $\mathcal{O}(mn^3 / \varepsilon^2)$ deterministic time [7]. The latter result is a direct corollary of the spectral sparsification of positive semi-definite (psd) matrices problem as defined in [39]; see also [27] for more applications. Here we present a fast deterministic algorithm for spectral sparsification of psd matrices and, as a consequence, we obtain an improved deterministic spectral graph sparsification algorithm for the case of dense graphs.

The last problem that we analyze is the element-wise matrix sparsification problem. This problem was first introduced by Achlioptas and McSherry in [1]. They described sampling-based algorithms that select a small number of entries from an input matrix $A$, forming a sparse matrix $\widetilde{A}$, which is close to $A$ in the operator norm sense. The motivation to study this problem lies on the need to speed up several matrix computations including approximate eigenvector computations [1] and semi-definite programming solvers [4,13]. Recently, there are many follow-up results on this problem [5,14]. To the best of our knowledge, all known algorithms for this problem are randomized (see Table 1 of [14]). In this paper we present the first deterministic algorithm and strong sparsification bounds for self-adjoint matrices that have an approximate diagonally dominant[1] property. Diagonally dominant matrices arise in many applications such as the solution of certain elliptic differential equations via the finite element method [11], several optimization problems in computer vision [23] and computer graphics [21], to name a few.

*Organization of the Paper.* The paper is organized as follows. In § 2, we present the matrix hyperbolic cosine algorithm (Algorithm 1). We apply the matrix hyperbolic cosine algorithm to derive improved deterministic algorithms for the construction of Alon-Roichman expanding Cayley graphs in § 3, spectral sparsification of psd matrices in § 4 and element-wise matrix sparsification in § 5. Due to space constraints, almost all proofs have been deferred to the Appendix.

**Our Results**

The main contribution of this paper is a generalization of Spencer's hyperbolic cosine algorithm to the matrix-valued setting [33], [36, Lecture 4], see Algorithm 1. As mentioned in the introduction, our main result has connections with a recent derandomization of matrix concentration inequalities [44]. We should highlight a few advantages of our result compared to [44]. First, our construction does not rely on composing two separate estimators (or potential functions) to achieve operator norm bounds and does not require knowledge of the sampling probabilities of the matrix samples as in [44]. In addition, the algorithm of [44] requires computations of matrix expectations with matrix exponentials which are computationally expensive, see [44, Footnote 6, p. 63]. In this paper, we demostrate that overcoming these limitations leads to faster and in some cases simpler algorithms.

Next, we demonstrate the usefulness of the main result by analyzing its computational efficiency under two special cases of matrices. We begin by presenting the following result

**Theorem 1 (Restatement of Theorem 5).** *There is a deterministic algorithm that, given the multiplication table of a group $G$ of size $n$, constructs an Alon-Roichman expanding Cayley graph of logarithmic*

---

[1] A self-adjoint matrix $A$ of size $n$ is called *diagonally dominant* if $|A_{ii}| \geq \sum_{j \neq i} |A_{ij}|$ for every $i \in [n]$.

*degree in $\mathcal{O}(n^2 \log^3 n)$ time. Moreover, the algorithm performs only group algebra operations that correspond to counting closed paths in Cayley graphs.*

To the best of our knowledge, the above theorem improves the running time of all previously known deterministic constructions of Alon-Roichman Cayley graphs [6,44,22]. Moreover, notice that the running time of the above algorithm is optimal up-to poly-logarithmic factors since the size of the multiplication table of a finite group of size $n$ is $\mathcal{O}(n^2)$.

In addition, we study the computational efficiency of the matrix hyperbolic cosine algorithm on the case of matrix samples with rank-one. The motivation for studying this special setting is its connection with problems such as graph approximation via edge sparsification as was shown in [7,39] and matrix approximation via element-wise sparsification as we will see later in this paper. The main result for this setting can be summarized in the following theorem (see § 4), which improves the $\mathcal{O}(mn^3/\varepsilon^2)$ running time of [39] when, say, $m = \Omega(n^2)$ and $\varepsilon$ is a constant.

**Theorem 2.** *Suppose $0 < \varepsilon < 1$ and $A = \sum_{i=1}^m v_i \otimes v_i$ are given, with column vectors $v_i \in \mathbb{R}^n$. Then there are non-negative real weights $\{s_i\}_{i \leq m}$, at most $\lceil n/\varepsilon^2 \rceil$ of which are non-zero, such that*

$$(1 - \varepsilon)^3 A \preceq \widetilde{A} \preceq (1 + \varepsilon)^3 A,$$

*where $\widetilde{A} = \sum_{i=1}^m s_i v_i \otimes v_i$. Moreover, there is a deterministic algorithm which computes the weights $s_i$ in[2] $\widetilde{\mathcal{O}}(mn^2 \log^3 n/\varepsilon^2 + n^4 \log n/\varepsilon^4)$ time.*

First, as we have already mentioned the graph sparsification problem can be reduced to spectral sparsification of positive semi-definite matrix. Hence as a corollary of the above theorem (proof omitted, see [39] for details), we obtain a fast deterministic algorithm for sparsifying dense graphs, which improves the currently best known $\mathcal{O}(n^5/\varepsilon^2)$ running time for this problem.

**Corollary 1.** *Given a weighted dense graph $H = (V, E)$ on $n$ vertices with positive weights and $0 < \varepsilon < 1$, there is a deterministic algorithm that returns an $(1 + \varepsilon)$-spectral sparsifier with $\mathcal{O}(n/\varepsilon^2)$ edges in $\widetilde{\mathcal{O}}(n^4 \log n/\varepsilon^2 \max\{\log^2 n, 1/\varepsilon^2\})$ time.*

Second, we give an elementary connection between element-wise matrix sparsification and spectral sparsification of psd matrices. A direct application of this connection implies strong sparsification bounds for self-adjoint matrices that are close to being *diagonally dominant*. More precisely, we give two element-wise sparsification algorithms for self-adjoint and diagonally dominant-like matrices; in its randomized and the other in its derandomized version (see Table 1 of [14] for comparison). Here, for the sake of presentation, we state our results for diagonally dominant matrices, although the results hold under a more general setting (see § 5 for details).

**Theorem 3.** *Let $A$ be any self-adjoint and diagonally dominant matrix of size $n$ and $0 < \varepsilon < 1$. Assume for normalization that $\|A\| = 1$.*

(a) *There is a randomized linear time algorithm that outputs a matrix $\widetilde{A} \in \mathbb{R}^{n \times n}$ with at most $\mathcal{O}(n \log n/\varepsilon^2)$ non-zero entries such that, with probability at least $1 - 1/n$, $\left\|A - \widetilde{A}\right\| \leq \varepsilon$.*

(b) *There is a deterministic $\widetilde{\mathcal{O}}(\varepsilon^{-2} \boldsymbol{nnz}(A) \, n^2 \log n \max\{\log^2 n, 1/\varepsilon^2\})$ time algorithm that outputs a matrix $\widetilde{A} \in \mathbb{R}^{n \times n}$ with at most $\mathcal{O}(n/\varepsilon^2)$ non-zero entries such that $\left\|A - \widetilde{A}\right\| \leq \varepsilon$.*

*Preliminaries.* The next discussion reviews several definitions and facts from linear algebra; for more details, see [9]. By $[n]$ to be the set $\{1, 2, \ldots, n\}$. We denote by $\mathcal{S}^{n \times n}$ the set of symmetric matrices of size $n$. Let $x \in \mathbb{R}^n$, we denote by $\mathbf{diag}(x)$ the diagonal matrix containing $x_1, x_2, \ldots, x_n$. For a square matrix $M$, we also write $\mathbf{diag}(M)$ to denote the diagonal matrix that contains the diagonal entries of $M$. Let $A$ be an $m \times n$ matrix. $A^{(j)}$ will denote the $j$-th column of $A$ and $A_{(i)}$ the $i$-th row of $A$. We denote

---

[2] The $\widetilde{\mathcal{O}}(\cdot)$ notation hides $\log \log n$ and $\log \log(1/\varepsilon)$ factors throughout the paper.

$\|A\| = \max\{\|Ax\| \mid \|x\| = 1\}$, $\|A\|_\infty = \max_{i \in [m]} \sum_{j \in [n]} |A_{ij}|$ and by $\|A\|_{\mathrm{F}} = \sqrt{\sum_{i,j} A_{ij}^2}$ the Frobenius norm of $A$. Also $\mathbf{sr}(A) := \|A\|_{\mathrm{F}}^2 / \|A\|^2$ is the *stable rank* of $A$ and by $\mathbf{nnz}(A)$ the number of its non-zero entries. The trace of a square matrix $B$ is denoted as $\mathbf{tr}(B)$. We write $\mathbf{J}_n$ for the all-ones square matrices of size $n$. For two self-adjoint matrices $X, Y$, we say that $Y \succeq X$ if and only if $Y - X$ is a positive semi-definite (psd) matrix. Let $x \in \mathbb{R}^n$, then $x \otimes x$ is the $n \times n$ matrix such that $(x \otimes x)_{i,j} = x_i x_j$. Given any matrix $A$, its *dilation* is defined as $\mathcal{D}(A) = \begin{bmatrix} \mathbf{0} & A \\ A^\top & \mathbf{0} \end{bmatrix}$. It is easy to see that $\lambda_{\max}(\mathcal{D}(A)) = \|A\|$, see e.g. [40, Theorem 4.2].

*Functions of Matrices.* Here we review some basic facts about the matrix exponential and the hyperbolic cosine function, for more details see [19]. All proofs of this section have been deferred to the appendix. The matrix exponential of a self-adjoint matrix $A$ is defined as $\mathbf{exp}[A] = \mathbf{I} + \sum_{k=1}^\infty \frac{A^k}{k!}$. Let $A = Q \Lambda Q^\top$ be the eigendecomposition of $A$. It is easy to see that $\mathbf{exp}[A] = Q \mathbf{exp}[\Lambda] Q^\top$. For any real square matrices $A$ and $B$ of the same size that commute, i.e., $AB = BA$, we have that $\mathbf{exp}[A + B] = \mathbf{exp}[A]\mathbf{exp}[B]$. In general, when $A$ and $B$ do not commute, the following estimate is known for self-adjoint matrices.

**Lemma 1.** *[16,41] For any self-adjoint matrices $A$ and $B$, $\mathbf{tr}(\mathbf{exp}[A + B]) \leq \mathbf{tr}(\mathbf{exp}[A]\,\mathbf{exp}[B])$.*

We will also need the following fact about matrix exponential for rank one matrices.

**Lemma 2.** *Let $x$ be a non-zero vector in $\mathbb{R}^n$. Then $\mathbf{exp}[x \otimes x] = \mathbf{I}_n + \frac{e^{\|x\|^2} - 1}{\|x\|^2} x \otimes x$. Similarly, $\mathbf{exp}[-x \otimes x] = \mathbf{I}_n - \frac{1 - e^{-\|x\|^2}}{\|x\|^2} x \otimes x$.*

Let us define the *matrix hyperbolic cosine* function of a self-adjoint matrix $A$ as $\mathbf{cosh}[A] := (\mathbf{exp}[A] + \mathbf{exp}[-A])/2$. Next, we state a few properties of the matrix hyperbolic cosine.

**Lemma 3.** *Let $A$ be a self-adjoint matrix. Then $\mathbf{tr}(\mathbf{exp}[\mathcal{D}(A)]) = 2\mathbf{tr}(\mathbf{cosh}[A])$.*

**Lemma 4.** *Let $A$ be a self-adjoint matrix and $P$ be a projector matrix that commutes with $A$, i.e., $PA = AP$. Then $\mathbf{cosh}[PA] = P\mathbf{cosh}[A] + \mathbf{I} - P$.*

**Lemma 5.** *[43, Lemma 2.2] For any positive semi-definite self-adjoint matrix $A$ of size $n$ and any two self-adjoint matrices $B, C$ of size $n$, $B \preceq C$ implies $\mathbf{tr}(AB) \leq \mathbf{tr}(AC)$.*

## 2 Balancing Matrices: a matrix hyperbolic cosine algorithm

We briefly describe Spencer's balancing vectors game and then generalize it to the matrix-valued setting [36, Lecture 4]. Let a two-player perfect information game between Alice and Bob. The game consists of $n$ rounds. On the $i$-th round, Alice sends a vector $v_i$ with $\|v_i\|_\infty \leq 1$ to Bob, and Bob has to decide on a sign $s_i \in \{\pm 1\}$ knowing only his previous choices of signs and $\{v_k\}_{k<i}$. At the end of the game, Bob pays Alice $\|\sum_{i=1}^n s_i v_i\|_\infty$. We call the latter quantity, the *value* of the game.

It has been shown in [35] that, in the above limited online variant, Spencer's six standard deviations bound [34] does not hold and the best value that we can hope for is $\Omega(\sqrt{n \ln n})$. Such a bound is easy to obtain by picking the signs $\{s_i\}$ uniformly at random. Indeed, a direct application of Azuma's inequality to each coordinate of the random vector $\sum_{i=1}^n s_i v_i$ together with a union bound over all the coordinates gives a bound of $\mathcal{O}(\sqrt{n \ln n})$.

Now, we generalize the balancing vectors game to the matrix-valued setting. That is, Alice now sends to Bob a sequence $\{M_i\}$ of self-adjoint matrices of size $n$ with[3] $\|M_i\| \leq 1$, and Bob has to pick a sequence of signs $\{s_i\}$ so that, at the end of the game, the quantity $\|\sum_{i=1}^n s_i M_i\|$ is as small as possible. Notice that the balancing vectors game is a restriction of the balancing matrices game in which Alice is allowed to send only diagonal matrices with entries bounded in absolute value by one. Similarly to the balancing vectors game, using matrix-valued concentration inequalities, one can prove that Bob has a randomized strategy that achieves at most $\mathcal{O}(\sqrt{n \ln n})$ w.p. at least $1/2$. Indeed,

---

[3] A curious reader may ask him/her-self why the operator norm is the right choice. It turns out the the operator norm is the correct matrix-norm analog of the $\ell_\infty$ vector-norm, viewed as the *infinity* Schatten norm on the space of matrices.

**Lemma 6.** *Let $M_i \in \mathcal{S}^{n \times n}$, $\|M_i\| \leq 1$, $1 \leq i \leq n$. Pick $s_i^* \in \{\pm 1\}$ uniformly at random for every $i \in [n]$. Then $\|\sum_{i=1}^{n} s_i^* M_i\| = \mathcal{O}(\sqrt{n \ln n})$ w.p. at least $1/2$.*

Now, let's assume that Bob wants to achieve the above probabilistic guarantees using a *deterministic* strategy. Is it possible? We answer this question in the affirmative by generalizing Spencer's hyperbolic cosine algorithm (and its proof) to the matrix-valued setting. We call the resulting algorithm *matrix hyperbolic cosine* (Algorithm 1). It is clear that this simple greedy algorithm implies a deterministic strategy for Bob that achieves the probabilistic guarantees of Lemma 6 (set $f_j \sim s_j M_j$, $t = n$ and $\varepsilon = \mathcal{O}(\sqrt{\ln n/n})$ and notice that $\gamma, \rho^2$ are at most one).

Algorithm 1 requires an extra assumption on its random matrices compared to Spencer's original algorithm. That is, we assume that our random matrices have uniformly bounded their "matrix variance", denoted by $\rho^2$. This requirement is motivated by the fact that in the applications that are studied in this paper such an assumption translates bounds that depend quadratically on the matrix dimensions to bounds that depend linearly on the dimensions.

We will need the following technical lemma for proving the main result of this section, which is a Bernstein type argument generalized to the matrix-valued setting [42].

**Lemma 7.** *Let $f : [m] \to \mathcal{S}^{n \times n}$ with $\|f(i)\| \leq \gamma$ for all $i \in [m]$. Let $X$ be a random variable over $[m]$ such that $\mathbb{E} f(X) = \mathbf{0}$ and $\left\| \mathbb{E} f(X)^2 \right\| \leq \rho^2$. Then, for any $\theta > 0$, $\left\| \mathbb{E}[\mathbf{exp}\left[\mathcal{D}\left(\theta f(X)\right)\right]] \right\| \leq \exp\left(\rho^2 (e^{\theta \gamma} - 1 - \theta \gamma)/\gamma^2\right)$. In particular, for any $0 < \varepsilon < 1$, setting $\theta = \varepsilon/\gamma$ implies that $\mathbb{E}[\mathbf{exp}\left[\mathcal{D}\left(\varepsilon f(X)/\gamma\right)\right]] \preceq e^{\varepsilon^2 \rho^2/\gamma^2} \mathbf{I}_{2n}$.*

Now we are ready to prove the correctness of the matrix hyperbolic cosine algorithm.

---

**Algorithm 1** Matrix Hyperbolic Cosine

---

1: **procedure** MATRIX-HYPERBOLIC($\{f_j\}$, $\varepsilon$, $t$)        $\triangleright$ $f_j : [m] \to \mathcal{S}^{n \times n}$ as in Theorem 4, $0 < \varepsilon < 1$.
2:    Set $\theta = \varepsilon/\gamma$
3:    **for** $i = 1$ to $t$ **do**
4:        Compute $x_i^* \in [m]$: $x_i^* = \arg\min_{k \in [m]} \mathbf{tr}\left(\mathbf{cosh}\left[\theta \sum_{j=1}^{i-1} f_j(x_j^*) + \theta f_i(k)\right]\right)$
5:    **end for**
6:    **Output:** $t$ indices $x_1^*, x_2^*, \ldots, x_t^*$ such that $\left\| \frac{1}{t} \sum_{j=1}^{t} f_j(x_j^*) \right\| \leq \frac{\gamma \ln(2n)}{t\varepsilon} + \frac{\varepsilon \rho^2}{\gamma}$
7: **end procedure**

---

**Theorem 4.** *Let $f_j : [m] \to \mathcal{S}^{n \times n}$ with $\|f_j(i)\| \leq \gamma$ for all $i \in [m]$ and $j = 1, 2, \ldots$. Suppose that there exists independent random variables $X_1, X_2, \ldots$ over $[m]$ such that $\mathbb{E} f_j(X_j) = \mathbf{0}$ and $\left\| \mathbb{E} f_j(X_j)^2 \right\| \leq \rho^2$. Algorithm 1 with input $\{f_j\}, \varepsilon, t$ outputs a set of indices $\{x_j^*\}_{j \in [t]}$ over $[m]$ such that $\left\| \frac{1}{t} \sum_{j=1}^{t} f_j(x_j^*) \right\| \leq \frac{\gamma \ln(2n)}{t\varepsilon} + \frac{\varepsilon \rho^2}{\gamma}$.*

We conclude with an open question related to Spencer's six standard deviation bound [34]. Does Spencer's six standard deviation bound holds under the matrix setting? More formally, given any sequence of $n$ self-adjoint matrices $\{M_i\}$ with $\|M_i\| \leq 1$, does there exist a set of signs $\{s_i\}$ so that $\|\sum_{i=1}^{n} s_i M_i\| = \mathcal{O}(\sqrt{n})$?

## 3   Alon-Roichman Expanding Cayley Graphs

We start by describing expander graphs. Given a connected undirected $d$-regular graph $H = (V, E)$ on $n$ vertices, let $A$ be its adjacency matrix, i.e., $A_{ij} = w_{ij}$ where $w_{ij}$ is the number of edges between vertices $i$ and $j$. Moreover, let $\widehat{A} = \frac{1}{d}A$ be its normalized adjacency matrix. We allow self-loops and multiple edges. Let $\lambda_1(\widehat{A}), \ldots, \lambda_n(\widehat{A})$ be its eigenvalues in decreasing order. We have that $\lambda_1(\widehat{A}) = 1$ with corresponding eigenvector $\mathbf{1}/\sqrt{n}$, where $\mathbf{1}$ is the all-one vector. The graph $H$ is called a spectral expander if $\lambda(\widehat{A}) := \max_{2 \leq j}\{|\lambda_j(\widehat{A})|\} \leq \varepsilon$ for some positive constant $\varepsilon < 1$.

Denote by $m_k = m_k(H) := \mathbf{tr}\left(A^k\right)$. By definition, $m_k$ is equal to the number of self-returning walks of length $k$ of the graph $H$. A graph-spectrum-based invariant, recently proposed by Estrada is defined as $EE(A) := \mathbf{tr}\left(\exp[A]\right)$ [15], which also equals to $\sum_{k=0}^{\infty} m_k/k!$. For $\theta > 0$, we define the *even $\theta$-Estrada index* by $EE_{\text{even}}(A, \theta) := \sum_{k=0}^{\infty} m_{2k}(\theta A)/(2k)!$.

Now let $G$ be any finite group of order $n$ with identity element $\mathtt{id}$. Let $S$ be a multi-set of elements of $G$, we denote by $S \sqcup S^{-1}$ the symmetric closure of $S$, namely the number of occurrences of $s$ and $s^{-1}$ in $S \sqcup S^{-1}$ equals the number of occurrences of $s \in S$. Let $R$ be the right regular representation[4], i.e., $(R(g_1)\phi)(g_2) = \phi(g_1 g_2)$ for every $\phi : G \to \mathbb{R}$ and $g_1, g_2 \in G$. The Cayley graph $\mathrm{Cay}\,(G; S)$ on a group $G$ with respect to the multi-set $S \subset G$ is the graph whose vertex set is $G$, and where $g_1$ and $g_2$ are connected by an edge if there exists $s \in S$ such that $g_2 = g_1 s$ (allowing multiple edges for multiple elements in $S$). In this section we prove the correctness of the following greedy algorithm for constructing expanding Cayley graphs.

**Theorem 5.** *Algorithm 2, given the multiplication table of a finite group $G$ of size $n$ and $0 < \varepsilon < 1$, outputs a (symmetric) multi-set $S \subset G$ of size $\mathcal{O}(\log n/\varepsilon^2)$ such that $\lambda(\mathrm{Cay}\,(G; S)) \leq \varepsilon$ in $\mathcal{O}(n^2 \log^3 n/\varepsilon^3)$ time. Moreover, the algorithm performs only group algebra operations that correspond to counting closed paths in Cayley graphs.*

---

**Algorithm 2** Expander Cayley Graph via even Estrada Index Minimization

---
1: **procedure** GREEDYESTRADAMIN($G, \varepsilon$)                  ▷ Multiplication table of $G$, $0 < \varepsilon < 1$
2:     Set $S^{(0)} = \emptyset$ and $t = \mathcal{O}(\log n/\varepsilon^2)$
3:     **for** $i = 1, \ldots t$ **do**
4:         Let $g_* \in G$ that (approximately) min. the even $\varepsilon/2$-Estrada index of $\mathrm{Cay}\left(G; S^{(i-1)} \cup g \cup g^{-1}\right)$ over all $g \in G$
                                                                    ▷ Use Lemma 9
5:         Set $S^{(i)} = S^{(i-1)} \cup g_* \cup g_*^{-1}$
6:     **end for**
7:     **Output:** A multi-set $S := S^{(t)}$ of size $2t$ such that $\lambda(\mathrm{Cay}\,(G; S)) \leq \varepsilon$
8: **end procedure**

---

Let $\widehat{A}$ be the normalized adjacency matrix of $\mathrm{Cay}\left(G; S \sqcup S^{-1}\right)$ for some $S \subset G$. It is not hard to see that $\widehat{A} = \frac{1}{2|S|}\sum_{s \in S}(R(s) + R(s^{-1}))$. We want to bound $\lambda(A)$. Notice that $\lambda(A) = \|(\mathbf{I} - \mathbf{J}/n)A\|$. Since we want to analyze the second-largest eigenvalue (in absolute value), we consider $(\mathbf{I} - \mathbf{J}/n)A = \frac{1}{|S|}\sum_{s \in S}(R(s) + R(s^{-1}))/2 - \mathbf{J}/n$. Based on the above calculation, we define our matrix-valued function as

$$f(g) := (R(g) + R(g^{-1}))/2 - \mathbf{J}/n \tag{1}$$

for every $g \in G$. The following lemma connects the potential function that is used in Theorem 4 and the even Estrada index.

**Lemma 8.** *Let $S \subset G$ and $A$ be the adjacency matrix of $\mathrm{Cay}\left(G; S \sqcup S^{-1}\right)$. For any $\theta > 0$, $\boldsymbol{tr}\left(\boldsymbol{cosh}\left[\theta \sum_{s \in S} f(s)\right]\right) = EE_{even}(A, \theta/2) + 1 - \cosh(\theta|S|)$.*

The following lemma indicates that it is possible to efficiently compute the (even) Estrada index for Cayley graphs with small generating set.

**Lemma 9.** *Let $S \subset G$, $\theta, \delta > 0$, and $A$ be the adjacency matrix of $\mathrm{Cay}\,(G; S)$. There is an algorithm that, given $S$, computes an additive $\delta$ approximation to $EE(\theta A)$ or $EE_{even}(A, \theta)$ in $\mathcal{O}(n|S| \max\{\log(n/\delta), 2e^2|S|\theta\})$ time.*

*Proof.* (of Theorem 5) By Lemma 8, minimizing the even $\varepsilon/2$-Estrada index in the $i$-th iteration is equivalent to minimizing $\boldsymbol{tr}\left(\boldsymbol{cosh}\left[\theta \sum_{s \in S^{(i-1)}} f(s) + \theta f(g)\right]\right)$ over all $g \in G$ with $\theta = \varepsilon$. Notice that $f(g) \in \mathcal{S}^{n \times n}$ for $g \in G$, $\mathbb{E}_{g \in_R G} f(g) = \mathbf{0}_n$ since $\sum_{g \in G} R(g) = \mathbf{J}$. It is easy to see that $\|f(g)\| \leq 2$ and moreover a calculation implies that $\left\|\mathbb{E}_{g \in_R G} f(g)^2\right\| \leq 2$ as well. Theorem 4 implies that we get a multi-set $S$ of size $t$ such that $\lambda(\mathrm{Cay}\left(G; S \sqcup S^{-1}\right)) = \left\|\frac{1}{|S|}\sum_{s \in S} f(s)\right\| \leq \varepsilon$. The moreover part follows from Lemma 9 with $\delta = \frac{e^{\varepsilon^2}}{n^c}$ for a sufficient large constant $c > 0$. Indeed, in total we incur (following the proof of Theorem 4) at most an additive $\ln(\delta n e^{\varepsilon^2 t})/\varepsilon$ error which is bounded by $\varepsilon$.

---
[4] In other words, represent each group algebra element with a permutation matrix of size $n$ that preserves the group structure. This is always possible due to Cayley's theorem.

## 4  Fast Isotropic Sparsification and Spectral Sparsification

Let $A$ be an $m \times n$ matrix with $m \gg n$ whose columns are in isotropic position, i.e., $A^\top A = \mathbf{I}_n$. For $0 < \varepsilon < 1$, consider the problem of finding a small subset of (rescaled) rows of $A$ forming a matrix $\widetilde{A}$ such that $\left\| \widetilde{A}^\top \widetilde{A} - \mathbf{I} \right\| \le \varepsilon$. The matrix Bernstein inequality (see [42]) tells us that there exists such a set with size $\mathcal{O}(n \log n/\varepsilon^2)$. Indeed, set $f(i) = A_{(i)} \otimes A_{(i)}/p_i - \mathbf{I}_n$ where $p_i = \left\| A_{(i)} \right\|^2 / \|A\|_{\mathrm{F}}^2$. A calculation shows that $\gamma$ and $\rho^2$ are $\mathcal{O}(n)$. Moreover, Algorithm 1 implies an $\mathcal{O}(mn^4 \log n/\varepsilon^2)$ time algorithm for finding such a set. The running time of Algorithm 1 for rank-one matrix samples can be improved to $\mathcal{O}(mn^3 \text{polylog}\,(n)/\varepsilon^2)$ by exploiting their rank-one structure. More precisely, using fast algorithms for computing all the eigenvalues of matrices after rank-one updates [18]. Next we show that we can further improve the running time by a more careful analysis.

We show how to improve the running time of Algorithm 1 to $\mathcal{O}(\frac{mn^2}{\varepsilon^2}\text{polylog}\,(n, \frac{1}{\varepsilon}))$ utilizing results from numerical linear algebra including the Fast Multipole Method [12] (FMM) and ideas from [18]. The main idea behind the improvement is that the trace is invariant under any change of basis. At each iteration, we perform a change of basis so that the matrix corresponding to the previous choices of the algorithm is diagonal. Now, Step 4 of Algorithm 1 corresponds to computing all the eigenvalues of $m$ different eigensystems with special structure, i.e., diagonal plus a rank-one matrix. Such eigensystem can be solved in $\mathcal{O}(n\text{polylog}\,(n))$ time using the FMM as was observed in [18]. However, the problem now, is that at each iteration we have to represent all the vectors $A_{(i)}$ in the new basis, which may cost $\mathcal{O}(mn^2)$. The key observation is that the change of basis matrix at each iteration is a Cauchy matrix (see Appendix). It is known that matrix-vector multiplication with Cauchy matrices can be performed efficiently and numerically stable using FMM. Therefore, at each iteration, we can perform the change of basis in $\mathcal{O}(mn\text{polylog}\,(n))$ and $m$ eigenvalue computations in $\mathcal{O}(mn\text{polylog}\,(n))$ time. The next theorem states that the resulting algorithm runs in $\mathcal{O}(mn^2\text{polylog}\,(n))$ time (see Appendix for proof).

**Theorem 6.** *Let $A$ be an $m \times n$ matrix with $A^\top A = \mathbf{I}_n$, $m \ge n$ and $0 < \varepsilon < 1$. Algorithm 3 returns at most $t = \mathcal{O}(n \ln n/\varepsilon^2)$ indices $x_1^*, x_2^*, \ldots x_t^*$ over $[m]$ with corresponding scalars $s_1, s_2, \ldots, s_t$ using $\widetilde{\mathcal{O}}(mn^2 \log^3 n/\varepsilon^2)$ operations such that*

$$\left\| \sum_{i=1}^t s_i A_{(x_i^*)} \otimes A_{(x_i^*)} - \mathbf{I}_n \right\| \le \varepsilon. \tag{2}$$

---

**Algorithm 3** Fast Isotropic Sparsification

1: **procedure** $\textsc{Isotrop}(A, \varepsilon)$      $\triangleright$ $A \in \mathbb{R}^{m \times n}$, $\sum_{k=1}^m A_{(k)} \otimes A_{(k)} = \mathbf{I}_n$ and $0 < \varepsilon < 1$
2:      Set $\theta = \varepsilon/n$, $t = \mathcal{O}(n \ln n/\varepsilon^2)$, and $A_{(k)} \leftarrow A_{(k)}/\sqrt{p_k}$ for every $k \in [m]$, where $p_k = \left\| A_{(k)} \right\|^2 / n$
3:      Set $\Lambda_{\{0\}} = \mathbf{0}_n$ and $Z = \sqrt{\theta}\, A$
4:      **for** $i = 1$ to $t$ **do**
5:          $x_i^* = \arg\min_{k \in [m]} \mathbf{tr}\left(\exp\left[\Lambda_{\{i-1\}} + Z_{(k)} \otimes Z_{(k)}\right] \mathrm{e}^{-\theta i} + \exp\left[-\Lambda_{\{i-1\}} - Z_{(k)} \otimes Z_{(k)}\right] \mathrm{e}^{\theta i}\right)$    $\triangleright$ Apply $m$ times Lemma 12
6:          $[\Lambda_{\{i\}}, U_{\{i\}}] = \mathbf{eigs}(\Lambda_{\{i-1\}} + Z_{(x_i^*)} \otimes Z_{(x_i^*)})$        $\triangleright$ **eigs** computes eigensystem
7:          $Z = Z U_{\{i\}}$        $\triangleright$ Apply fast matrix-vector multiplication
8:      **end for**
9:      **Output:** $t$ indices $x_1^*, x_2^*, \ldots, x_t^*$, $x_i^* \in [m]$ s.t. $\left\| \sum_{k=1}^t \frac{A_{(x_k^*)} \otimes A_{(x_k^*)}}{t p_{x_k^*}} - \mathbf{I}_n \right\| \le \varepsilon$
10: **end procedure**

---

Next, we show that Algorithm 3 can be used as a bootstrapping procedure to improve the time complexity of [39, Theorem 3.1], see also [7, Theorem 3.1]. Such an improvement implies faster algorithms for constructing graph sparsifiers and, as we will see in § 5, element-wise sparsification of matrices.

**Theorem 7.** *Suppose $0 < \varepsilon < 1$ and $A = \sum_{i=1}^m v_i \otimes v_i$ are given, with column vectors $v_i \in \mathbb{R}^n$ and $m \ge n$. Then there are non-negative weights $\{s_i\}_{i \le m}$, at most $\lceil n/\varepsilon^2 \rceil$ of which are non-zero, such that*

$$(1 - \varepsilon)^3 A \preceq \widetilde{A} \preceq (1 + \varepsilon)^3 A, \tag{3}$$

where $\widetilde{A} = \sum_{i=1}^{m} s_i v_i \otimes v_i$. Moreover, there is an algorithm that computes the weights $\{s_i\}_{i \leq m}$ in deterministic $\widetilde{\mathcal{O}}(mn^2 \log^3 n/\varepsilon^2 + n^4 \log n/\varepsilon^4)$ time.

# 5 Element-wise Matrix Sparsification

A deterministic algorithm for the element-wise matrix sparsification problem can be obtained by derandomizing a recent result whose analysis is based on the matrix Bernstein inequality [14].

**Theorem 8.** *Let $A$ be an $n \times n$ matrix and $0 < \varepsilon < 1$. There is a deterministic polynomial time algorithm that, given $A$ and $\varepsilon$, outputs a matrix $\widetilde{A} \in \mathbb{R}^{n \times n}$ with at most $28n \ln(\sqrt{2n}) \boldsymbol{sr}(A) / \varepsilon^2$ non-zero entries such that $\left\| A - \widetilde{A} \right\| \leq \varepsilon \|A\|$.*

Next, we give two improved element-wise sparsification algorithms for self-adjoint and diagonally dominant-like matrices; one of them is randomized and the other is its derandomized version. Both algorithms share a crucial difference with all previously known algorithms for this problem; during their execution they may densify the diagonal entries of the input matrices. On the one hand, there are at most $n$ diagonal entries, so this does not affect asymptotically their sparsity guarantees. On the other hand, as we will see later this twist turns out to give strong sparsification bounds.

Recall that the results of [38,7] imply an element-wise sparsification algorithm that works only for Laplacian matrices. It is easy to verify that Laplacian matrices are also diagonally dominant. Here we extend these results to a wider class of matrices (with a weaker notion of approximation). The diagonally dominant assumption is too restrictive and we will show that our sparsification algorithms work for a wider class of matrices. To accommodate this, we say that a matrix $A$ is $\theta$-symmetric diagonally dominant (abbreviate by $\theta$-SDD) if $A$ is self-adjoint and the inequality $\|A\|_\infty \leq \sqrt{\theta} \|A\|$ holds.

By definition, any diagonally dominant matrix is also a 4-SDD matrix. On the other extreme, every self-adjoint matrix of size $n$ is $n$-SDD since the inequality $\|A\|_\infty \leq \sqrt{n} \|A\|$ is always valid. The following elementary lemma gives a connection between element-wise matrix sparsification and spectral sparsification as defined in [39].

**Lemma 10.** *Let $A$ be a self-adjoint matrix of size $n$ and $R = \boldsymbol{diag}(R_1, R_2, \ldots, R_n)$ where $R_i = \sum_{j \neq i} |A_{ij}|$. Then there is a matrix $C$ of size $n \times m$ with $m \leq \binom{n}{2}$ such that*

$$A = CC^\top + \boldsymbol{diag}(A) - R. \tag{4}$$

*Moreover, each column of $C$ is indexed by the ordered pairs $(i,j)$, $i < j$ and equals to $C^{(i,j)} = \sqrt{|A_{ij}|} e_i + \boldsymbol{sgn}(A_{ij}) \sqrt{|A_{ij}|} e_j$ for every $i < j$, $i,j \in [n]$.*

*Remark 1.* In the special case where $A$ is the Laplacian matrix of some graph, the above decomposition is precisely the vertex-edge decomposition of the Laplacian matrix, since in this case $\boldsymbol{diag}(A) = R$.

Using the above lemma, we give a randomized and a deterministic algorithm for sparsifying $\theta$-SDD matrices. First we present the randomized algorithm.

**Theorem 9.** *Let $A$ be a $\theta$-SDD matrix of size $n$ and $0 < \varepsilon < 1$. There is a randomized linear time algorithm that, given $A$, $\|A\|$ and $\varepsilon$, outputs a matrix $\widetilde{A} \in \mathbb{R}^{n \times n}$ with at most $\mathcal{O}(n\theta \log n/\varepsilon^2)$ non-zero entries such that w.p. at least $1 - 1/n$, $\left\| A - \widetilde{A} \right\| \leq \varepsilon \|A\|$.*

Next we state the derandomized algorithm of the above result.

**Theorem 10.** *Let $A$ be a $\theta$-SDD matrix of size $n$ and $0 < \varepsilon < 1/2$. There is an algorithm that, given $A$ and $\varepsilon$, outputs a matrix $\widetilde{A} \in \mathbb{R}^{n \times n}$ with at most $\mathcal{O}(n\theta/\varepsilon^2)$ non-zero entries such that $\left\| A - \widetilde{A} \right\| \leq \varepsilon \|A\|$. Moreover, the algorithm computes $\widetilde{A}$ in deterministic $\widetilde{\mathcal{O}}(\boldsymbol{nnz}(A) n^2 \theta \log^3 n/\varepsilon^2 + n^4 \theta^2 \log n/\varepsilon^4)$ time.*

*Remark 2.* The results of [7,39] imply a deterministic $\mathcal{O}(\boldsymbol{nnz}(A) \theta n^3/\varepsilon^2)$ time algorithm that outputs a matrix $\widetilde{A}$ with at most $\lceil 19(1 + \sqrt{\theta})^2/\varepsilon^2 \rceil n$ non-zero entries such that $\left\| \widetilde{A} - A \right\| \leq \varepsilon \|A\|$.

## Acknowledgements

## References

1. D. Achlioptas and F. McSherry. Fast Computation of Low-rank Matrix Approximations. *SIAM J. Comput.*, 54(2):9, 2007.
2. R. Ahlswede and A. Winter. Strong Converse for Identification via Quantum Channels. *IEEE Transactions on Information Theory*, 48(3):569–579, 2002.
3. N. Alon and Y. Roichman. Random Cayley Graphs and Expanders. *Random Struct. Algorithms*, 5:271–284, 1994.
4. S. Arora, E. Hazan, and S. Kale. Fast Algorithms for Approximate Semidefinite Programming using the Multiplicative Weights Update Method. In *Proceedings of the Symposium on Foundations of Computer Science (FOCS)*, pages 339–348, 2005.
5. S. Arora, E. Hazan, and S. Kale. A Fast Random Sampling Algorithm for Sparsifying Matrices. In *Proceedings of the International Workshop on Randomization and Approximation Techniques (RANDOM)*, pages 272–279, 2006.
6. S. Arora and S. Kale. A Combinatorial, Primal-Dual Approach to Semidefinite Programs. In *Proceedings of the Symposium on Theory of Computing (STOC)*, pages 227–236, 2007.
7. J. D. Batson, D. A. Spielman, and N. Srivastava. Twice-ramanujan sparsifiers. In *Proceedings of the Symposium on Theory of Computing (STOC)*, pages 255–262, 2009.
8. A. A. Benczúr and D. R. Karger. Approximating $s$-$t$ Minimum Cuts in $\widetilde{\mathcal{O}}(n^2)$ Time. In *Proceedings of the Symposium on Theory of Computing (STOC)*, 1996.
9. R. Bhatia. *Matrix Analysis*, volume 169. Graduate Texts in Mathematics, Springer, First edition, 1996.
10. D. Bini and V. Y. Pan. *Polynomial and Matrix Computations: Fundamental Algorithms*, volume 1. Birkhauser Verlag, 1994.
11. E. G. Boman, B. Hendrickson, and S. Vavasis. Solving Elliptic Finite Element Systems in Near-Linear Time with Support Preconditioners. *SIAM J. on Numerical Analysis*, 46(6):3264–3284, 2008.
12. J. Carrier, L. Greengard, and V. Rokhlin. A Fast Adaptive Multipole Algorithm for Particle Simulations. *SIAM J. on Scientific and Statistical Computing*, 9(4):669–686, 1988.
13. A. 'd Aspremont. Subsampling Algorithms for Semidefinite Programming. *Stochastic Systems*, pages 274–305, 2011.
14. P. Drineas and A. Zouzias. A note on Element-wise Matrix Sparsification via a Matrix-valued Bernstein Inequality. *Information Processing Letters*, 111(8), 2011.
15. E. Estrada and J. A. Rodríguez-Velázquez. Subgraph Centrality in Complex Networks. *Phys. Rev. E*, 71, May 2005.
16. S. Golden. Lower Bounds for the Helmholtz Function. *Phys. Rev.*, 137(4B):B1127–B1128, 1965.
17. G. Golub. Some Modified Eigenvalue Problems. In *Conference on Applications of Numerical Analysis*, volume 228 of *Lecture Notes in Mathematics*, pages 56–56. 1971.
18. M. Gu and S. C. Eisenstat. A Stable and Efficient Algorithm for the Rank-One Modification of the Symmetric Eigenproblem. *SIAM J. Matrix Anal. Appl.*, 15, 1994.
19. N. J. Higham. *Functions of Matrices: Theory and Computation*. Society for Industrial and Applied Mathematics (SIAM), 2008.
20. R. A. Horn and C. R. Johnson. *Matrix Analysis*. Cambridge University Press, 1990.
21. P. Joshi, M. Meyer, T. DeRose, B. Green, and T. Sanocki. Harmonic Coordinates for Character Articulation. *ACM Trans. Graph.*, 26, 2007.
22. S. Kale. *Efficient Algorithms Using the Multiplicative Weights Update Method*. PhD in Computer Science, Princeton University, 2007.
23. I. Koutis, G. L. Miller, and D. Tolliver. Combinatorial Preconditioners and Multilevel Solvers for Problems in Computer Vision and Image Processing. In *Proceedings of the Symposium on Advances in Visual Computing (ISVC)*, pages 1067–1078, 2009.
24. Z. Landau and A. Russell. Random Cayley Graphs are Expanders: a simplified proof of the Alon-Roichman theorem. *The Electronic J. of Combinatorics*, 11(1), 2004.

25. A. Magen and A. Zouzias. Low Rank Matrix-Valued Chernoff Bounds and Approximate Matrix Multiplication. In *Proceedings of the ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 1422–1436, 2011.
26. A. Naor. On the Banach Space Valued Azuma Inequality and Small set Isoperimetry of Alon-Roichman Graphs. To Appear in Combinatorics, Probability and Computing. Available at arxiv:1009.5695, September 2010.
27. A. Naor. Sparse Quadratic Forms and their Geometric Applications (after Batson, Spielman and Srivastava). Available at arxiv:1101.4324, January 2011.
28. D. Peleg and A. A. Schäffer. Graph Spanners. *J. of Graph Theory*, 13(1):99–116, 1989.
29. B. Recht. A Simpler Approach to Matrix Completion. *J. of Machine Learning Research*, pages 3413–3430, December 2011.
30. M. Rudelson. Random Vectors in the Isotropic Position. *J. Funct. Anal.*, 164(1):60–72, 1999.
31. M. Rudelson and R. Vershynin. Sampling from Large Matrices: An Approach through Geometric Functional Analysis. *J. ACM*, 54(4):21, 2007.
32. A. Man-Cho So. Moment Inequalities for sums of Random Matrices and their Applications in Optimization. *Mathematical Programming*, pages 1–27, 2009.
33. J. Spencer. Balancing Games. *J. Comb. Theory, Ser. B*, 23(1):68–74, 1977.
34. J. Spencer. Six Standard Deviations Suffice. *Transactions of The American Mathematical Society*, 289:679–679, 1985.
35. J. Spencer. Balancing Vectors in the max Norm. *Combinatorica*, 6:55–65, 1986.
36. J. Spencer. *Ten Lectures on the Probabilistic Method*. Society for Industrial and Applied Mathematics (SIAM), Second edition, 1994.
37. D. A. Spielman. Algorithms, Graph Theory, and Linear Equations in Laplacian Matrices. In *Proceedings of the International Congress of Mathematicians*, volume IV, pages 2698–2722, 2010.
38. D. A. Spielman and N. Srivastava. Graph Sparsification by Effective Resistances. In *Proceedings of the Symposium on Theory of Computing (STOC)*, 2008.
39. N. Srivastava. *Spectral Sparsification and Restricted Invertibility*. PhD in Computer Science, Yale University, 2010.
40. G. W. Stewart and J. G. Sun. *Matrix Perturbation Theory (Computer Science and Scientific Computing)*. Academic Press, 1990.
41. C. J. Thompson. Inequality with Applications in Statistical Mechanics. *J. of Mathematical Physics*, 6(11):1812–1813, 1965.
42. J. A. Tropp. User-Friendly Tail Bounds for Sums of Random Matrices. *Foundations of Computational Mathematics*, pages 1–46, 2011.
43. K. Tsuda, G. Rätsch, and M. K. Warmuth. Matrix Exponentiated Gradient Updates for on-line Learning and Bregman Projections. *JMLR*, 6:995–1018, 2005.
44. A. Wigderson and D. Xiao. Derandomizing the Ahlswede-Winter Matrix-valued Chernoff Bound using Pessimistic Estimators, and Applications. *Theory of Computing*, 4(1):53–76, 2008.
45. A. Zouzias. A Matrix Hyperbolic Cosine Algorithm and Applications. Ver. 1. Available at arxiv:1103.2793, March 2011.

# Appendix

## Fast Multiplication with Cauchy Matrices and Special Eigensystems

We start by defining the so-called Cauchy (generalized Hilbert) matrices. An $m \times n$ matrix $C$ defined by

$$C_{i,j} := \frac{1}{t_i - s_j}, \quad i \in [m], j \in [n],$$

where $t = (t_1, \ldots, t_m)$, $t \in \mathbb{R}^m$ and $s = (s_1, \ldots, s_n)$, $s \in \mathbb{R}^n$ and $t_i \neq s_j$ for all $i \in [m]$ and $j \in [n]$ is called *Cauchy*. Given a vector $x \in \mathbb{R}^n$, the naive algorithm for computing the matrix-vector product $Cx$ requires $\mathcal{O}(mn)$ operations. It is not clear if it is possible to perform this computation in less than $\mathcal{O}(mn)$ operations. Surprisingly enough, it is possible to compute this product with $\mathcal{O}((m + n) \log^2(m + n))$ operations. This computation can be done by two different approaches. The first one is based on fast polynomial multiplication, polynomial interpolation and polynomial evaluation at distinct points [10, Algorithm 1, p. 130]. The main drawback of this approach is its numerical instability. The second approach is based on the so-called Fast

Multipole Method (FMM) introduced in [12]. This method returns an approximate solution to the matrix-vector product for any given error parameter[5]. Ignoring numerical issues that are beyond the scope of this work, we summarize our discussion to the following

**Lemma 11.** *[10,12] Let $x \in \mathbb{R}^n$ and $C$ be a Cauchy matrix defined as above with $t \in \mathbb{R}^m, s \in \mathbb{R}^n$. There is an algorithm that, given vectors $s, t, x$, computes the product $Cx$ using $\mathcal{O}((m+n)\log^2(m+n))$ operations.*

Given a self-adjoint matrix $B = \Sigma + \rho u \otimes u$, where $\Sigma = \mathbf{diag}(\sigma_1, \ldots, \sigma_n)$, $\rho > 0$ and $u \in \mathbb{R}^n$ is a unit vector, our goal is to efficiently compute all the eigenvalues of $B$. It is well-known that the eigenvalues of $B$ are the roots of a special function, known as secular function [17] and are interlaced with $\{\sigma_i\}_{i \leq n}$. In addition, evaluating the secular function requires $\mathcal{O}(n)$ operations implying that a standard (Newton) root-finding procedure requires $\mathcal{O}(n)$ operations per each eigenvalue. Hence, $\mathcal{O}(n^2)$ operations are required for all eigenvalues. In their seminal paper [18], Gu and Eisenstat showed that it is possible to encode the updates of the root-finding procedure for *all* eigenvalues as matrix-vector multiplication with an $n \times n$ Cauchy matrix. Based on this observation, they showed how to use the Fast Multipole Method for approximately computing all the eigenvalues of this special type of eigenvalue problem.

**Lemma 12.** *[18] Let $b \in \mathbb{N}$, $\rho > 0$, $\Sigma = \mathbf{diag}(\sigma_1, \sigma_2, \ldots, \sigma_n)$ and $u \in \mathbb{R}^n$ be a unit vector. There is an algorithm that given $\Sigma, \rho, u$ computes all the eigenvalues of $B = \Sigma + \rho u \otimes u$ within an additive error $2^{-b}\|B\|$ in $\mathcal{O}(n\log^2 n \log b)$ operations.*

## Omitted Proofs

*Proof.* (of Lemma 2) The proof is immediate by the definition of the matrix exponential. Notice that $(x \otimes x)^k = \|x\|^{2(k-1)} x \otimes x$ for $k \geq 1$.

$$\mathbf{exp}[x \otimes x] = \mathbf{I} + \sum_{k=1}^{\infty} \frac{(x \otimes x)^k}{k!} = \mathbf{I} + \sum_{k=1}^{\infty} \frac{\|x\|^{2(k-1)} x \otimes x}{k!} = \mathbf{I} + \frac{e^{\|x\|^2} - 1}{\|x\|^2} x \otimes x.$$

Similar considerations give that $\mathbf{exp}[-x \otimes x] = \mathbf{I} - \frac{1 - e^{-\|x\|^2}}{\|x\|^2} x \otimes x$.

*Proof.* (of Lemma 3) Set $B := \mathcal{D}(A) = \begin{bmatrix} \mathbf{0} & A \\ A^\top & \mathbf{0} \end{bmatrix}$. Notice that for any integer $k \geq 1$, $B^{2k} = \begin{bmatrix} A^{2k} & \mathbf{0} \\ \mathbf{0} & A^{2k} \end{bmatrix}$ and $B^{2k+1} = \begin{bmatrix} \mathbf{0} & A^{2k+1} \\ A^{2k+1} & \mathbf{0} \end{bmatrix}$. Since the odd powers of $B$ are trace-less, it follows that

$$\mathbf{tr}(\mathbf{exp}[B]) = \mathbf{tr}\left(\mathbf{I}_{2n} + \sum_{k=1}^{\infty} \frac{B^{2k}}{(2k)!} + \sum_{k=0}^{\infty} \frac{B^{2k+1}}{(2k+1)!}\right) = \mathbf{tr}\left(\mathbf{I}_{2n} + \sum_{k=1}^{\infty} \frac{B^{2k}}{(2k)!}\right)$$

$$= 2\mathbf{tr}\left(\mathbf{I}_n + \sum_{k=1}^{\infty} \frac{A^{2k}}{(2k)!}\right) = \mathbf{tr}(\mathbf{exp}[A] + \mathbf{exp}[-A]) = 2\mathbf{tr}(\mathbf{cosh}[A]).$$

*Proof.* (of Lemma 4) By the definition of $\mathbf{cosh}[\cdot]$, it suffices to show that $\mathbf{exp}[PA] = P\mathbf{exp}[A] + \mathbf{I} - P$,

$$\mathbf{exp}[PA] = \mathbf{I} + \sum_{k=1}^{\infty} \frac{(PA)^k}{k!} = \mathbf{I} + P \sum_{k=1}^{\infty} \frac{A^k}{k!} = P\mathbf{exp}[A] + \mathbf{I} - P.$$

---

[5] That is, given an $n \times n$ Cauchy matrix, a vector $x \in \mathbb{R}^n$ and $0 < \varepsilon < 1$, returns a vector $y \in \mathbb{R}^n$ so that $\|y - Cx\|_\infty \leq \varepsilon$ in time $\mathcal{O}(n\log^2(1/\varepsilon))$. In an actual implementation, setting $\varepsilon$ to be a small constant relative to the machine's (numerical) precision suffices; see [18, § 3] for a more careful implementation and discussion on numerical issues.

*Proof.* (of Lemma 6) We wish to apply matrix Azuma's inequality, see [42, Theorem 7.1]. For every $j \in [n]$, define the matrix-valued difference sequence $f_j : [2] \to \mathcal{S}^{n \times n}$ as $f_j(k) = (2(k-1)-1)M_j$ with $\|f_j(\cdot)\| \leq 1$. Let $X$ be a uniform random variable over the set $[2]$. Then $\mathbb{E}_X f_j(X) = \mathbf{0}_n$. Set $\varepsilon = \sqrt{10 \ln(4n)/n}$. Matrix-valued Azuma's inequality tells us that w.p. at least $1/2$, a random set of signs $\{s_j\}_{j \in [n]}$ satisfies $\left\| \frac{1}{n} \sum_{j=1}^n s_j M_j \right\| \leq \varepsilon$. Rescale the last inequality to conclude.

*Proof.* (of Theorem 4) Using the notation of Algorithm 1, for every $i = 1, 2, \ldots, t$, define recursively $W(i) := \theta \sum_{j=1}^i f_j(x_j^*)$ and the potential function $\Phi^{(i)} := 2\mathbf{tr}\left(\mathbf{cosh}\left[W(i)\right]\right)$. For all steps $i = 1, 2, \ldots, t$, we will prove that

$$\Phi^{(i)} \leq \Phi^{(i-1)} \exp\left(\varepsilon^2 \rho^2 / \gamma^2\right). \tag{5}$$

Assume that the algorithm has fixed the first $(i-1)$ indices $x_1^*, \ldots, x_{(i-1)}^*$. An averaging argument applied on the expression of the argmin of Step 4 gives that

$$\begin{aligned}
\mathbb{E}_{X_i} 2\mathbf{tr}\left(\mathbf{cosh}\left[\theta W(i-1) + \theta f_i(X_i)\right]\right) &= \mathbb{E}_{X_i} \mathbf{tr}\left(\exp\left[\theta \mathcal{D}\left(W(i-1)\right) + \theta \mathcal{D}\left(f_i(X_i)\right)\right]\right) \\
&\leq \mathbf{tr}\left(\exp\left[\mathcal{D}\left(\theta W(i-1)\right)\right] \mathbb{E}_{X_i} \exp\left[\mathcal{D}\left(\theta f_i(X_i)\right)\right]\right) \\
&\leq \mathbf{tr}\left(\exp\left[\mathcal{D}\left(\theta W(i-1)\right)\right] \mathbf{I}_{2n}\right) \exp\left(\varepsilon^2 \rho^2 / \gamma^2\right) \\
&= \Phi^{(i-1)} \exp\left(\varepsilon^2 \rho^2 / \gamma^2\right)
\end{aligned}$$

where in the first inequality we used Lemma 3 and linearity of dilation, in the second inequality we used the Golden-Thompson inequality (Lemma 1) and linearity of trace, in the third inequality we used Lemma 5 together with Lemma 7 and in the last equality we used again Lemma 3. Since the algorithm seeks the minimum of the expression in Step 4, it follows that $\Phi^{(i)} \leq \mathbb{E}_{X_i} 2\mathbf{tr}\left(\exp\left[\theta \mathcal{D}\left(W(i-1)\right) + \theta \mathcal{D}\left(f_i(X_i)\right)\right]\right)$ which proves Ineq. (5). Apply $t$ times Ineq. (5) to conclude that $\Phi^{(t)} \leq \Phi^{(0)} \exp\left(t \frac{\varepsilon^2 \rho^2}{\gamma^2}\right)$. Recall that $\Phi^{(0)} = 2\mathbf{tr}\left(\mathbf{cosh}\left[\mathbf{0}_n\right]\right) = 2\mathbf{tr}\left(\mathbf{I}_n\right) = 2n$. On the other hand, we can lower bound $\Phi^{(t)}$

$$\Phi^{(t)} = 2\mathbf{tr}\left(\mathbf{cosh}\left[\theta \sum_{j=1}^t f_j(x_j^*)\right]\right) \geq \exp\left(\left\|\theta \sum_{j=1}^t f_j(x_j^*)\right\|\right).$$

The last inequality follows since $2\mathbf{tr}\left(\mathbf{cosh}\left[C\right]\right) = 2 \sum_{i=1}^n \cosh(\lambda_i(C)) \geq 2\cosh\left(\lambda_{\max}(C)\right) + 2\cosh\left(\lambda_{\min}(C)\right) \geq \exp(\|C\|)$ for any matrix $C \in \mathcal{S}^{n \times n}$. Take logarithms on both sides and divide by $\theta$, we conclude that $\left\|\sum_{j=1}^t f_j(x_j^*)\right\| \leq \frac{\ln(2n)}{\theta} + t\frac{\varepsilon^2 \rho^2}{\theta \gamma^2}$. Rescale by $t$ the last inequality to conclude the proof.

*Proof.* (of Lemma 8) For notational convenience, set $P := \mathbf{I}_n - \mathbf{J}_n/n$ and $B := \frac{\theta}{2} \sum_{s \in S}(R(s) + R(s)^{-1})$. Since $\mathbf{J}R(g) = R(g)\mathbf{J} = \mathbf{J}$, we have that $\mathbf{tr}\left(\mathbf{cosh}\left[\theta \sum_{s \in S} f(s)\right]\right) = \mathbf{tr}\left(\mathbf{cosh}\left[PB\right]\right)$. Now using Lemma 4, it follows $\mathbf{tr}\left(\mathbf{cosh}\left[PB\right]\right) = \mathbf{tr}\left(P\mathbf{cosh}\left[B\right] + \mathbf{I} - P\right) = \mathbf{tr}\left(\mathbf{cosh}\left[B\right]\right) + \mathbf{tr}\left(-\frac{\mathbf{J}}{n}\mathbf{cosh}\left[B\right] + \mathbf{I} - P\right)$ Notice that $\mathbf{J}/n$ is a projector matrix, hence applying Lemmas 2,4 we get that

$$\mathbf{tr}\left(-\frac{\mathbf{J}}{n}\mathbf{cosh}\left[B\right] + \mathbf{I} - P\right) = \mathbf{tr}\left(-\mathbf{cosh}\left[\mathbf{J}/nB\right] + P + \mathbf{I} - P\right) = 1 - \cosh(\theta|S|).$$

*Proof.* (of Lemma 9) We will prove the Lemma for $EE(A, \theta)$, the other case is similar. Let $h := \theta \sum_{s \in S} s$ be a group algebra element of $G$, i.e, $h \in \mathbb{R}[G]$. Define $\mathbf{exp}\left[h\right] := \mathtt{id} + \sum_{k=1}^\infty \frac{h^{\star k}}{k!}$ and $T_l(h) := \mathtt{id} + \sum_{k=1}^l \frac{h^{\star k}}{k!}$ (where $h^{\star k}$ is the $k$-folded convolution/multiplication over $\mathbb{R}[G]$) the exponential operator and its $l$ truncated Taylor series, respectively. Notice that $\theta A = \theta \sum_{s \in S} R(s) = R(h)$, so $EE(A, \theta) = \mathbf{tr}\left(\exp\left[R(h)\right]\right) = \mathbf{tr}\left(R(\mathbf{exp}\left[h\right])\right)$. We will show that the quantity $\mathbf{tr}\left(R(T_l(h))\right)$ is a $\delta$ approximation for $EE(A, \theta)$ when $l \geq \max\{\log(n/\delta), 2\mathrm{e}^2|S|\theta\}$.

Compute the sum of $T_l(h)$ by summing each term one by one and keeping track of all the coefficients of the group algebra elements. The main observation is that at each step there are at most $n$ such coefficients since we

are working over $\mathbb{R}[G]$. For $k > 1$, compute the $k$-th term of the sum by $(\sum_{s \in S} c_s s)^k / k! = (\sum_{s \in S} c_s s)^{k-1} / (k-1)! \cdot \sum_{s \in S} (c_s/k)s$. Assume that we have computed the first term of the above product, which is some group algebra element denote it by $\sum_{g \in G} \beta_g g$ for some $\beta_g \in \mathbb{R}$. Hence, at the next iteration, we have to compute the product/convolution of $\sum_{g \in G} \beta_g g$ with $\theta/k \sum_{s \in S} s$, which can be done in $\mathcal{O}(n|S|)$ time. Since the sum has $l$ terms, in total we require $\mathcal{O}(n|S|l)$ operations. Now, we show that it is a $\delta$ approximation. We need the following fact (see [19, Theorem 10.1, p. 234])

**Fact 11** *For any $B \in \mathbb{R}^{n \times n}$, let $T_l(B) := \sum_{k=0}^{l} \frac{B^k}{k!}$. Then, $\|\boldsymbol{exp}\,[B] - T_l(B)\| \leq \frac{\|B\|^{l+1}}{(l+1)!} e^{\|B\|}$.*

Notice that $\|\theta A\| = \left\|\sum_{s \in S} \theta R(s)\right\| \leq \theta|S|$ by triangle inequality and the fact that $\|R(g)\| = 1$ for any $g \in G$. Applying Fact 11 on $\theta A$ we get that

$$\|\mathbf{exp}\,[\theta A] - T_l(\theta A)\| \leq \frac{(\theta|S|)^{l+1}}{(l+1)!} e^{\theta|S|} \leq \left(\frac{e\theta|S|}{l+1}\right)^{l+1} e^{\theta|S|}$$

$$= \left(\frac{e^{1+(\theta|S|)/(l+1)}\theta|S|}{l+1}\right)^{l+1} \leq \frac{1}{2^{l+1}} \leq \frac{\delta}{n}.$$

where we used the inequality $(l+1)! \geq (\frac{l+1}{e})^{l+1}$ and the assumption that $l \geq \max\{\log(n/\delta), 2e^2\theta|S|\}$.

**Lemma 13.** *Assume that the first $(i-1)$ indices, $i < t$ have been fixed by Algorithm 3. Let $\Phi_k^{(i)}$ be the value of the potential function when the index $k$ has been selected at the next iteration of the algorithm. Similarly, let $\widetilde{\Phi}_k^{(i)}$ be the (approximate) value of the potential function computed using Lemma 12 within an additive error $\delta > 0$ for all eigenvalues. Then,*

$$e^{-\delta}\Phi_k^{(i)} \leq \widetilde{\Phi}_k^{(i)} \leq e^{\delta}\Phi_k^{(i)}$$

*Proof.* Let $\tau_1, \tau_2, \ldots, \tau_n$ be the eigenvalues of $\Lambda_{\{i-1\}} + Z_{(k)} \otimes Z_{(k)}$. Let $\widetilde{\tau}_1, \widetilde{\tau}_2, \ldots, \widetilde{\tau}_n$ be the approximate eigenvalues of the latter matrix when computed via Lemma 12 within an additive error $\delta > 0$, i.e, $|\widetilde{\tau}_j - \tau_j| \leq \delta$ for all $j \in [n]$.

First notice that, by Step 5 of Algorithm 3, $\Phi_k^{(i)} = 2\sum_{j=1}^{n} \cosh(\tau_j - \lambda i)$. Similarly, $\widetilde{\Phi}_k^{(i)} := 2\sum_{j=1}^{n} \cosh(\widetilde{\tau}_j - \lambda i)$. By the definition of the hyperbolic cosine, we get that

$$\sum_{j=1}^{n} \cosh(\widetilde{\tau}_j - \lambda i) = \sum_{j=1}^{n} \cosh(\tau_j - \lambda i + \widetilde{\tau}_j - \tau_j)$$

$$= \frac{1}{2}\sum_{j=1}^{n}\left[\exp(\tau_j - \lambda i)\exp(\widetilde{\tau}_j - \tau_j) + \exp(-\tau_j + \lambda i)\exp(-\widetilde{\tau}_j + \tau_j)\right].$$

To derive the upper bound notice that $\sum_{j=1}^{n} \cosh(\widetilde{\tau}_j - \lambda i) \leq \sum_{j=1}^{n} \cosh(\tau_j - \lambda i)\max_{j \in [n]}\{\exp(\widetilde{\tau}_j - \tau_j), \exp(-\widetilde{\tau}_j + \tau_j)\}$ and the maximum is upper bounded by $\exp(\delta)$. Similarly, for the lower bound.

*Proof.* (of Theorem 6) The proof consists of three steps: $(a)$ we show that Algorithm 3 is a reformulation of Algorithm 1; $(b)$ we prove that in Step 5 of Algorithm 3 it is enough to compute the values of the potential function within a sufficiently small multiplicative error using Lemma 12, and $(c)$ we give the advertised bound on the running time of Algorithm 3.

Set $p_i = \|A_{(i)}\|^2 / \|A\|_F^2$, $f(i) = A_{(i)} \otimes A_{(i)}/p_i - \mathbf{I}_n$ and $s_i = 1/p_i$ for every $i \in [m]$. Observe that $\|A\|_F^2 = \mathbf{tr}\,(A^\top A) = \mathbf{tr}\,(\mathbf{I}_n) = n$. Let $X$ be a random variable distributed over $[m]$ with probability $p_i$. Notice that $\mathbb{E}\,f(X) = \mathbf{0}_n$ and $\gamma = n$, since $\|f(i)\| = \left\|nA_{(i)} \otimes A_{(i)}/\|A_{(i)}\|^2 - \mathbf{I}_n\right\| \leq n$ for every $i \in [m]$. Moreover, a direct calculation shows that $\mathbb{E}\,f(X)^2 = \mathbb{E}\,(A_{(X)} \otimes A_{(X)}/p_X)^2 - \mathbf{I}_n = n\sum_{i=1}^{m} A_{(i)} \otimes A_{(i)} - \mathbf{I}_n = (n-1)\mathbf{I}_n$, hence $\rho^2 \leq n$. Algorithm 1 with $t = \mathcal{O}(n \ln n/\varepsilon^2)$ returns indices $x_1^*, x_2^*, \ldots, x_t^*$ so that $\left\|\frac{1}{t}\sum_{j=1}^{t} f_j(x_j^*)\right\| \leq \frac{\gamma \ln(2n)}{t\varepsilon} + \varepsilon\rho^2/\gamma \leq 2\varepsilon$. We next prove by induction that the same set of indices are also returned by Algorithm 3.

For ease of presentation, rescale every row of the input matrix $A$, i.e., set $\widehat{A}_{(k)} = A_{(k)}\sqrt{\theta/p_k}$ for every $k \in [m]$ (see Steps 2 and 3 of Algorithm 3). For sake of the analysis, let us define the following sequence of self-adjoint matrices of size $n$

$$T_{\{0\}} := \mathbf{0}_n,$$
$$T_{\{i\}} := T_{\{i-1\}} + \widehat{A}_{(x_i^*)} \otimes \widehat{A}_{(x_i^*)} \text{ for } i \in [t]$$

with eigenvalue decompositions $T_{\{i\}} = Q_{\{i\}}\Lambda_{\{i\}}Q_{\{i\}}^\top$, where $\Lambda_{\{i\}}$ are diagonal matrices containing the eigenvalues and the columns of $Q_{\{i\}}$ contain the corresponding eigenvectors. Set $Q_{\{0\}} = \mathbf{I}$ and $\Lambda_{\{0\}} = \mathbf{0}$. Notice that for every $k \in [m]$, by the eigenvalue decomposition of $T_{\{i-1\}}$, $T_{\{i-1\}} + \widehat{A}_{(k)} \otimes \widehat{A}_{(k)} = Q_{\{i-1\}}\left(\Lambda_{\{i-1\}} + Q_{\{i-1\}}^\top\widehat{A}_{(k)} \otimes Q_{\{i-1\}}^\top\widehat{A}_{(k)}\right)Q_{\{i-1\}}^\top$. Observe that the above matrix (left hand side) and $\Lambda_{\{i-1\}} + Q_{\{i-1\}}^\top\widehat{A}_{(k)} \otimes Q_{\{i-1\}}^\top\widehat{A}_{(k)}$ have the same eigenvalues, since they are similar matrices. Let $\Lambda_{\{i-1\}} + Q_{\{i-1\}}^\top\widehat{A}_{(x_i^*)} \otimes Q_{\{i-1\}}^\top\widehat{A}_{(x_i^*)} = U_{\{i\}}\Lambda_{\{i\}}U_{\{i\}}^\top$ be its eigenvalue decomposition[6]. Then

$$T_{\{i-1\}} + \widehat{A}_{(x_i^*)} \otimes \widehat{A}_{(x_i^*)} = Q_{\{i-1\}}\left(\Lambda_{\{i-1\}} + Q_{\{i-1\}}^\top\widehat{A}_{(x_i^*)} \otimes Q_{\{i-1\}}^\top\widehat{A}_{(x_i^*)}\right)Q_{\{i-1\}}^\top$$
$$= Q_{\{i-1\}}U_{\{i\}}\Lambda_{\{i\}}U_{\{i\}}^\top Q_{\{i-1\}}^\top.$$

It follows that $Q_{\{i\}} = Q_{\{i-1\}}U_{\{i\}}$ for every $i \geq 1$, so $Q_{\{i\}} = U_{\{1\}}U_{\{2\}}\ldots U_{\{i\}}$. The base case of the induction is immediate. Now assume that Algorithm 3 has returned the same indices as Algorithm 1 up to the $(i-1)$-th iteration. It suffices to prove that at the $i$-th iteration Algorithm 3 will return the index $x_i^*$.

We start with the expression in Step 4 of Algorithm 1 and prove that it's equivalent (up to a fixed multiplicative constant factor) with the expression in Step 5 of Algorithm 3. Indeed, for any $k \in [m]$, (let $C := \theta\sum_{j=1}^{i-1}f(x_j^*)$)

$$2\mathbf{tr}\left(\mathbf{cosh}\left[C + \theta f(k)\right]\right) = \mathbf{tr}\left(\mathbf{exp}\left[C + \theta f(k)\right] + \mathbf{exp}\left[-C - \theta f(k)\right]\right)$$
$$= \mathbf{tr}\left(\mathbf{exp}\left[T_{\{i-1\}} + \widehat{A}_{(k)} \otimes \widehat{A}_{(k)}\right]\mathrm{e}^{-\theta i} + \mathbf{exp}\left[-T_{\{i-1\}} - \widehat{A}_{(k)} \otimes \widehat{A}_{(k)}\right]\mathrm{e}^{\theta i}\right)$$

where we used the definition of $\mathbf{cosh}\left[\cdot\right]$, $f(i)$ and $T_{\{i-1\}}$ and the fact that the matrices commute. In light of Algorithm 3 and the induction hypothesis, observe that the $m \times n$ matrix $Z$ at the start of the $i$-th iteration of Algorithm 3 is equal to $\widehat{A}U_{\{1\}}U_{\{2\}}\ldots U_{\{i-1\}} = \widehat{A}Q_{\{i-1\}}$. Now, multiply the latter expression that appears inside the trace with $Q_{\{i-1\}}^\top$ from the left and $Q_{\{i-1\}}$ from the right, it follows that ((let $C := \theta\sum_{j=1}^{i-1}f(x_j^*)$))

$$2\mathbf{tr}\left(\mathbf{cosh}\left[C + \theta f(k)\right]\right) = \mathbf{tr}\left(\mathbf{exp}\left[\Lambda_{\{i-1\}} + Z_{(k)} \otimes Z_{(k)}\right]\mathrm{e}^{-\theta i} + \mathbf{exp}\left[-\Lambda_{\{i-1\}} - Z_{(k)} \otimes Z_{(k)}\right]\mathrm{e}^{\theta i}\right)$$

using that $Q_{\{i-1\}}$ are the eigenvectors of $T_{\{i-1\}}$ and the cyclic property of trace. This concludes part $(a)$.

Next we discuss how to deal with the technicality that arises from the approximate computation of the arg min expression in Step 5 of Algorithm 3. First, let's assume that we have approximately (by invoking Lemma 12) minimized the potential function in Step 5 of Algorithm 3; denote this sequence of potential function values by $\widetilde{\Phi}^{(1)},\ldots,\widetilde{\Phi}^{(t)}$. Next, we sufficiently bound the parameter $b$ of Lemma 12 so that the above approximation will not incur a significant multiplicative error.

Recall that at every iteration, by Ineq. (5) there exists an index over $[m]$ such that the current value of the potential function increases by at most a multiplicative factor $\exp\left(\varepsilon^2\rho^2/\gamma^2\right)$. Lemma 13 tells us that at every iteration of Algorithm 3 we increase the value of the potential function (by not selecting the optimal index over $[m]$) by at most an *extra* multiplicative factor $\mathrm{e}^{2\delta}$, where $\delta$ is the additive error when computing the eigenvalues of the matrix in Step 5 via Lemma 12. Therefore, it follows that $\widetilde{\Phi}^{(t)} \leq \exp(2\delta t)\Phi^{(t)}$.

Observe that, at the $i$-th iteration we apply Lemma12 on a matrix $\sum_{j=1}^i \widehat{A}_{(x_j)} \otimes \widehat{A}_{(x_j)}$ for some indices $x_j \in [m]$ and moreover $\left\|\sum_{j=1}^i \widehat{A}_{(x_j)} \otimes \widehat{A}_{(x_j)}\right\| =$

---

[6] by its definition, $T_{\{i\}}$ has the same eigenvalues with $\Lambda_{\{i-1\}} + Q_{\{i-1\}}^\top\widehat{A}_{(x_i^*)} \otimes Q_{\{i-1\}}^\top\widehat{A}_{(x_i^*)}$.

$\left\| \theta \sum_{j=1}^{i} A_{(x_j)} \otimes A_{(x_j)} / p_{x_j} \right\| = \left\| \theta \sum_{j=1}^{i} f(x_j) - \theta i \mathbf{I} \right\|$. Triangle inequality tells us that $\left\| \sum_{j=1}^{i} \widehat{A}_{(x_j)} \otimes \widehat{A}_{(x_j)} \right\|$ is at most $2\gamma\theta t$ for every $i \in [t]$. It follows that $\delta$ is at most $2^{-b+1}\theta t\gamma$ where $b$ is specified in Lemma 12. The above discussion suggests that by setting $b = \mathcal{O}(\log(\theta\gamma t)) = \mathcal{O}(\log(n \log n/\varepsilon^3))$ we can guarantee that the potential function $\widetilde{\Phi}^{(t)} \le 2n \exp\left(3t\varepsilon^2\right)$. This concludes part $(b)$.

Finally, we conclude the proof by analyzing the running time of Algorithm 3. Steps 2 and 3 can be done in $\mathcal{O}(mn)$ time. Step 5 requires $\widetilde{\mathcal{O}}(mn \log^2 n)$ operations by invoking $m$ times Lemma 12. Steps 6 can be done in $\mathcal{O}(n^2)$ time and Step 7 requires $\widetilde{\mathcal{O}}(mn \log^2 n)$ operations by invoking Lemma 11. In total, since the number of iterations is $\mathcal{O}(n \log n/\varepsilon^2)$, the algorithm requires $\widetilde{\mathcal{O}}(mn^2 \log^3 n/\varepsilon^2)$ operations.

*Proof.* (of Theorem 7) Assume without loss of generality that $A$ has full rank. Define $u_i = A^{-1/2}v_i$ and notice that $\sum_{i=1}^{m} u_i \otimes u_i = \mathbf{I}_n$. Run Algorithm 3 with input $\{u_i\}_{i \in [m]}$ and $\varepsilon$ which returns $\{\tau_i\}_{i \le m}$, at most $t = \mathcal{O}(n \log n/\varepsilon^2)$ of which are non-zero such that

$$\left\| \sum_{i=1}^{m} \tau_i u_i \otimes u_i - \mathbf{I}_n \right\| \le \varepsilon. \tag{6}$$

Define $\widehat{A} = A^{1/2}\left(\sum_{i=1}^{m} \tau_i u_i \otimes u_i\right) A^{1/2} = \sum_{i=1}^{m} \tau_i v_i \otimes v_i$. Eqn. (6) is equivalent to $(1-\varepsilon)\mathbf{I}_n \preceq \sum_{i=1}^{m} \tau_i u_i \otimes u_i \preceq (1+\varepsilon)\mathbf{I}_n$. Conjugating the latter expression by $A^{1/2}$, see [20, § 7.7], we get that $(1-\varepsilon)A \preceq \widehat{A} \preceq (1+\varepsilon)A$. Apply [39, Theorem 3.1] on $\widehat{A}$ which outputs a matrix $\widetilde{A} = \sum_{i=1}^{m} s_i v_i \otimes v_i$ with non-negative weights $\{s_i\}_{i \in [m]}$ at most $\lceil n/\varepsilon^2 \rceil$ of which are non-zero, such that $(1-\varepsilon)^2\widehat{A} \preceq \widetilde{A} \preceq (1+\varepsilon)^2\widehat{A}$. Using the positive semi-definite partial ordering, we conclude that $(1-\varepsilon)^3 A \preceq \widetilde{A} \preceq (1+\varepsilon)^3 A$.

*Proof.* (of Theorem 8) By homogeneity, assume that $\|A\| = 1$. Following the proof of [14], we can assume that w.l.o.g. all non-zero entries of $A$ have magnitude at least $\varepsilon/(2n)$ in absolute value, otherwise we can zero-out these entries and incur at most an error of $\varepsilon/2$ (see [14, § 4.1]).

Consider the bijection $\pi$ between the sets $[n^2]$ and $[n] \times [n]$ defined by $\pi(l) \mapsto (\lceil l/n \rceil, (l-1) \mod n + 1)$ for every $l \in [n^2]$. Let $E_{ij} \in \mathbb{R}^{n \times n}$ be the all zeros matrix having one only in the $(i,j)$ entry. Set $h(l) = \mathcal{D}\left(\frac{A_{\pi(l)}}{p_l}E_{\pi(l)} - A\right)$ where $p_l = A^2_{\pi(l)}/\|A\|^2_F$ for every $l \in [n^2]$. Observe that $h(\cdot) \in \mathcal{S}^{2n \times 2n}$. Let $X$ be a random variable over $[n^2]$ with distribution $p_l$, $l \in [n^2]$. The same analysis as in Lemmas 2 and 3 of [14] together with properties of the dilation map imply that $\|h(l)\| \le 4n\mathbf{sr}(A)/\varepsilon$ for every $l \in [n^2]$, $\mathbb{E}\, h(X) = \mathbf{0}_{2n}$, and $\left\|\mathbb{E}\, h(X)^2\right\| \le n\mathbf{sr}(A)$.

Run Algorithm 1 with $h(\cdot)$ as above. Algorithm 1 returns at most $t = 28n \ln(\sqrt{2}n)\mathbf{sr}(A)/\varepsilon^2$ indices $x_1^*, x_2^*, \dots x_t^*$ over $[n^2]$ using $\mathcal{O}(n^6\mathbf{sr}(A) \log n/\varepsilon^2)$ operations such that

$$\left\| \frac{1}{t} \sum_{l=1}^{t} h(x_l^*) \right\| \le \varepsilon/2. \tag{7}$$

Set $\widetilde{A} := \frac{1}{t} \sum_{l=1}^{t} A_{\pi(x_l^*)}/p_{x_l^*} E_{\pi(x_l^*)}$. Observe that $\widetilde{A}$ has at most $t$ non-zero entries. Now, by the definition of $h(\cdot)$ and properties of the dilation map, it follows that Ineq. (7) is equivalent to $\left\| \mathcal{D}\left(\widetilde{A} - A\right) \right\| = \left\| \widetilde{A} - A \right\| \le \varepsilon/2$.

*Proof.* (of Lemma 10) The key identity is $CC^\top := \sum_{l,k \in [n],\ l<k} C^{(l,k)} \otimes C^{(l,k)}$. Let $l, k \in [n]$ with $l < k$, it follows that

$$C^{(l,k)} \otimes C^{(l,k)} = \left(\sqrt{|A_{lk}|}e_l + \mathbf{sgn}(A_{lk})\sqrt{|A_{lk}|}e_k\right)\left(\sqrt{|A_{lk}|}e_l + \mathbf{sgn}(A_{lk})\sqrt{|A_{lk}|}e_k\right)^\top$$
$$= |A_{lk}|e_l \otimes e_l + A_{lk}e_k \otimes e_k + A_{lk}e_k \otimes e_l + |A_{lk}|e_k \otimes e_k.$$

Therefore

$$CC^\top = \sum_{l,k \in [n]:\ l<k} [|A_{lk}|e_l \otimes e_l + A_{lk}e_k \otimes e_k + A_{lk}e_k \otimes e_l + |A_{lk}|e_k \otimes e_k]. \tag{8}$$

15

Let's first prove the equality for the off-diagonal entries of Eqn (4). Let $l < k$ and $l, k \in [n]$. By construction, the only term of the sum that contributes to the $(i, j)$ and $(j, i)$ entry of the right hand side of Eqn. (8) is the term $C^{(i,j)} \otimes C^{(i,j)}$. Moreover, this term equals $|A_{ij}|e_i \otimes e_i + A_{ij}e_i \otimes e_j + A_{ij}e_j \otimes e_i + |A_{ij}|e_j \otimes e_j$. Since $A_{ij} = A_{ji}$ this proves that the off-diagonal entries are equal.

For the diagonal entries of Eqn. (4), it suffices to prove that $(CC^\top)_{ii} = R_i$. First observe that the last two terms of the sum in the right hand side of (8) do not contribute to any diagonal entry. Second, the first two terms contribute only when $l = i$ or $k = i$. In the case where $l = i$, the contribution of the sum equals to $\sum_{i<k} |A_{ik}|$. On the other case ($k = i$), the contribution of the sum is equal to $\sum_{l<i} |A_{li}|$. However, $A$ is self-adjoint so $A_{li} = A_{il}$ for every $l < i$. It follows that the total contribution is $\sum_{i<k} |A_{ik}| + \sum_{l<i} |A_{il}| = \sum_{j \neq i} |A_{ij}| = R_i$.

*Proof.* ( of Theorem 9) In one pass over the input matrix $A$ normalize the entries of $A$ by $\|A\|$, so assume without loss of generality that $\|A\| = 1$. Let $C$ be the $n \times m$ matrix guaranteed by Lemma 10, where $m = \binom{n}{2}$, each column of $C$ is indexed by the ordered pairs $(i, j)$, $i < j$ and $A = CC^\top + \mathbf{diag}(A) - R$. By definition of $C$ and the hypothesis, we have that $\|CC^\top\| = \|A - \mathbf{diag}(A) + R\| \leq \|A\| + \|A\|_\infty \leq 2\sqrt{\theta}$ and $\|C\|_F^2 = 2\sum_{i,j} |A_{ij}| \leq 2n \|A\|_\infty \leq 2n\sqrt{\theta}$.

Consider the bijection between the sets $[m]$ and $\{(i, j) \mid i < j, \; i, j \in [n]\}$ defined by $\pi(l) \mapsto (\lceil l/n \rceil, (l-1) \bmod n + 1)$. For each $l \in [m]$, set $p_l = \|C^{\pi(l)}\|^2 / \|C\|_F^2$ and define $f(l) := C^{\pi(l)} \otimes C^{\pi(l)}/p_l - CC^\top$. Let $X$ be a real-valued random variable over $[m]$ with distribution $p_l$. It is easy to verify that $\mathbb{E} f(X) = \mathbf{0}_n$, $\|f(l)\| \leq 2\|C\|_F^2$ for every $l \in [m]$. A direct calculation gives that $\|\mathbb{E} f(X)^2\| \leq 2\|C\|_F^2 \|CC^\top\|$. Matrix Bernstein inequality (see [42]) with $f(\cdot)$ as above ($\gamma = 4n\sqrt{\theta}$ and $\rho^2 = 8n\theta$) tells us that if we sample $t = 38n\theta \ln(\sqrt{2}n)/\varepsilon^2$ indices $x_1^*, x_2^*, \ldots, x_t^*$ over $[m]$ then with probability at least $1 - 1/n$, $\left\|\frac{1}{t}\sum_{j=1}^t f(x_j^*)\right\| \leq \varepsilon$. Now, set $\widetilde{C} \in \mathbb{R}^{n \times t}$ where the $j$-th column of $\widetilde{C}^{(j)}$ equals $\frac{1}{\sqrt{t}} C^{\pi(x_j^*)}$. It follows that $\left\|\frac{1}{t}\sum_{j=1}^t f(x_j^*)\right\| = \left\|\frac{1}{t}\sum_{j=1}^t C^{\pi(x_j^*)} \otimes C^{\pi(x_j^*)} - CC^\top\right\| = \left\|\widetilde{C}\widetilde{C}^\top - CC^\top\right\|$. Define $\widetilde{A} = \widetilde{C}\widetilde{C}^\top + \mathbf{diag}(A) - R$. First notice that $\left\|\widetilde{A} - A\right\| = \left\|\widetilde{C}\widetilde{C}^\top - CC^\top\right\| \leq \varepsilon$. It suffices to bound the number of non-zeros of $\widetilde{A}$. To do so, view the matrix-product $\widetilde{C}\widetilde{C}^\top$ as a sum of rank-one outer-products over all columns of $\widetilde{C}$. By the special structure of the entries of $\widetilde{C}$, every outer-product term of the sum contributes to at most four non-zero entries, two of which are off-diagonal. Since $\widetilde{C}$ has at most $t$ columns, $\widetilde{A}$ has at most $n + 2t$ non-zero entries; $n$ for the diagonal entries and $2t$ for the off-diagonal.

*Proof.* (of Theorem 10) Let $C$ be the $n \times m$ matrix such that $A = CC^\top + \mathbf{diag}(A) - R$ and $m \leq \mathbf{nnz}(A)$ guaranteed by Lemma 10. Apply Theorem 7 on the matrix $CC^\top$ and $\varepsilon$ which outputs, in deterministic $\widetilde{\mathcal{O}}(\mathbf{nnz}(A) n^2\theta \log^3 n/\varepsilon^2 + n^4\theta^2 \log n/\varepsilon^4)$ time, an $n \times \lceil n/\varepsilon^2 \rceil$ matrix $\widetilde{C}$ such that $(1-\varepsilon)^3 CC^\top \preceq \widetilde{C}\widetilde{C}^\top \preceq (1+\varepsilon)^3 CC^\top$. By Weyl's inequality [20, Theorem 4.3.1] and the fact that $\varepsilon < 1/2$, it follows that $\left\|CC^\top - \widetilde{C}\widetilde{C}^\top\right\| \leq 5\varepsilon \|CC^\top\|$. Define $\widetilde{A} := \widetilde{C}\widetilde{C}^\top + \mathbf{diag}(A) - R$. First we argue that the number of non-zero entries of $\widetilde{A}$ is at most $n + \lceil 2n/\varepsilon^2 \rceil$. Recall that every column of $\widetilde{C}$ is a rescaled column of $C$. Now, think the matrix-product $\widetilde{C}\widetilde{C}^\top$ as a sum of rank-one outer-products over all columns of $\widetilde{C}$. By the special structure of the entries of $\widetilde{C}$, every outer-product term of the sum contributes to at most four non-zero entries, two of which are off-diagonal. Since $\widetilde{C}$ has at most $\lceil n/\varepsilon^2 \rceil$ columns, $\widetilde{A}$ has at most $n + \lceil 2n/\varepsilon^2 \rceil$ non-zero entries; $n$ for the diagonal entries and $\lceil 2n/\varepsilon^2 \rceil$ for the off-diagonal. Moreover, $\widetilde{A}$ is close to $A$ in the operator norm sense. Indeed,

$$\left\|A - \widetilde{A}\right\| = \left\|CC^\top - \widetilde{C}\widetilde{C}^\top\right\| \leq 5\varepsilon \|CC^\top\| = 5\varepsilon \|A - \mathbf{diag}(A) + R\|$$
$$\leq 5\varepsilon(\|A\| + \|A\|_\infty) \leq 10\varepsilon\sqrt{\theta} \|A\|$$

where we used the definition of $\widetilde{A}$, Eqn. (4), triangle inequality, the assumption that $A$ is $\theta$-SDD and the fact that $\theta \geq 1$. Repeating the proof with $\varepsilon' = \frac{\varepsilon}{10\sqrt{\theta}}$ and elementary manipulations conclude the proof.