

Solving Laplacian Linear Equations

Rasmus Kyng, Scribe: Hongjie Chen

Lecture 9 — Wednesday, April 22nd

1 Solving Linear Equations Approximately

Given a Laplacian L of a connected graph and a demand vector $\mathbf{d} \perp \mathbf{1}$, we want to find \mathbf{x}^* solving the linear equation $L\mathbf{x}^* = \mathbf{d}$. We are going to focus on fast algorithms for finding approximate (but highly accurate) solutions.

This means we need a notion of an approximate solution. Since our definition is not special to Laplacians, we state it more generally for positive semi-definite matrices.

Definition 1.1. Given PSD matrix M and $\mathbf{d} \in \ker(M)^\perp$, let $M\mathbf{x}^* = \mathbf{d}$. We say that $\tilde{\mathbf{x}}$ is an ϵ -approximate solution to the linear equation $M\mathbf{x} = \mathbf{d}$ if

$$\|\tilde{\mathbf{x}} - \mathbf{x}^*\|_M^2 \leq \epsilon \|\mathbf{x}^*\|_M^2.$$

Remark 1.2. The requirement $\mathbf{d} \in \ker(M)^\perp$ can be removed, but this is not important for us.

Theorem 1.3 (Spielman and Teng (2004) [ST04]). *Given a Laplacian L of a weighted undirected graph $G = (V, E, \mathbf{w})$ with $|E| = m$ and $|V| = n$ and a demand vector $\mathbf{d} \in \mathbb{R}^V$, we can find $\tilde{\mathbf{x}}$ that is an ϵ -approximate solution to $L\mathbf{x} = \mathbf{d}$, using an algorithm that takes time $O(m \log^c n \log(1/\epsilon))$ for some fixed constant c and succeeds with probability $1 - 1/n^{10}$.*

In the original algorithm of Spielman and Teng, the exponent on the log in the running time was $c \approx 70$.

Today, we are going to see a simpler algorithm. But first, we'll look at one of the key tools behind all algorithms for solving Laplacian linear equations quickly.

2 Preconditioning and Approximate Gaussian Elimination

Recall our definition of two positive semi-definite matrices being approximately equal.

Definition 2.1 (Spectral approximation). Given $\mathbf{A}, \mathbf{B} \in S_+^n$, we say that

$$\mathbf{A} \approx_K \mathbf{B} \text{ if and only if } \frac{1}{1+K} \mathbf{A} \preceq \mathbf{B} \preceq (1+K) \mathbf{A}.$$

Suppose we have a positive definite matrix $M \in S_{++}^n$ and want to solve a linear equation $M\mathbf{x} = \mathbf{d}$. We can do this using gradient descent or accelerated gradient descent, as we covered in Graded Homework 1. But if we have access to an easy-to-invert matrix that happens to also be a good spectral approximation of M , then we can use this to speed up the (accelerated) gradient descent

algorithm. An example of this would be that we have a factorization $\mathcal{L}\mathcal{L}^\top \approx_K \mathbf{M}$, where \mathcal{L} is lower triangular and sparse, which means we can invert it quickly.

The following lemma, which you will prove in Problem Set 6, makes this preconditioning precise.

Lemma 2.2. *Given a matrix $\mathbf{M} \in S_{++}^n$, a vector \mathbf{d} and a decomposition $\mathbf{M} \approx_K \mathcal{L}\mathcal{L}^\top$, we can find $\tilde{\mathbf{x}}$ that ϵ -approximately solves $\mathbf{M}\mathbf{x} = \mathbf{d}$, using $O((1+K)\log(K/\epsilon)(T_{\text{matvec}} + T_{\text{sol}} + n))$ time.*

- T_{matvec} denotes the time required to compute $\mathbf{M}\mathbf{z}$ given a vector \mathbf{z} , i.e. a “matrix-vector multiplication”.
- T_{sol} denotes the time required to compute $\mathcal{L}^{-1}\mathbf{z}$ or $(\mathcal{L}^\top)^{-1}\mathbf{z}$ given a vector \mathbf{z} .

Dealing with pseudo-inverses. When our matrices have a null space, preconditioning becomes slightly more complicated, but as long as it is easy to project to the complement of the null space, there’s no real issue. The following describes precisely what we need (but you can ignore the null-space issue when first reading these notes without losing anything significant).

Lemma 2.3. *Given a matrix $\mathbf{M} \in S_+^n$, a vector $\mathbf{d} \in \ker(\mathbf{M})^\perp$ and a decomposition $\mathbf{M} \approx_K \mathcal{L}\mathcal{D}\mathcal{L}^\top$, where \mathcal{L} is invertible, we can find $\tilde{\mathbf{x}}$ that ϵ -approximately solves $\mathbf{M}\mathbf{x} = \mathbf{d}$, using $O((1+K)\log(K/\epsilon)(T_{\text{matvec}} + T_{\text{sol}} + T_{\text{proj}} + n))$ time.*

- T_{matvec} denotes the time required to compute $\mathbf{M}\mathbf{z}$ given a vector \mathbf{z} , i.e. a “matrix-vector multiplication”.
- T_{sol} denotes the time required to compute $\mathcal{L}^{-1}\mathbf{z}$ and $(\mathcal{L}^\top)^{-1}\mathbf{z}$ and $\mathcal{D}^+\mathbf{z}$ given a vector \mathbf{z} .
- T_{proj} denotes the time required to compute $\mathbf{\Pi}_\mathbf{M}\mathbf{z}$ given a vector \mathbf{z} .

Theorem 2.4 (Kying and Sachdeva (2015) [KS16]). *Given a Laplacian \mathbf{L} of a weighted undirected graph $G = (V, E, \mathbf{w})$ with $|E| = m$ and $|V| = n$, we can find a decomposition $\mathcal{L}\mathcal{L}^\top \approx_{0.5} \mathbf{L}$, such that \mathcal{L} has number of non-zeroes $\text{nnz}(\mathcal{L}) = O(m \log^3 n)$, with probability at least $1 - 3/n^5$. in time $O(m \log^3 n)$.*

We can combine Theorem 2.4 with Lemma 2.3 to get a fast algorithm for solving Laplacian linear equations.

Corollary 2.5. *Given a Laplacian \mathbf{L} of a weighted undirected graph $G = (V, E, \mathbf{w})$ with $|E| = m$ and $|V| = n$ and a demand vector $\mathbf{d} \in \mathbb{R}^V$, we can find $\tilde{\mathbf{x}}$ that is an ϵ -approximate solution to $\mathbf{L}\mathbf{x} = \mathbf{d}$, using an algorithm that takes time $O(m \log^3 n \log(1/\epsilon))$ and succeeds with probability $1 - 1/n^{10}$.*

Proof sketch. First we need to get a factorization that conforms to Lemma 2.3. The decomposition $\mathcal{L}\mathcal{L}^\top$ provided by Theorem 2.4 can be rewritten as $\mathcal{L}\mathcal{L}^\top = \tilde{\mathcal{L}}\mathcal{D}(\tilde{\mathcal{L}})^\top$ where $\tilde{\mathcal{L}}$ is equal to \mathcal{L} except $\mathcal{L}(n, n) = 1$ and we let \mathcal{D} be the identity matrix, except $\mathcal{D}(n, n) = 0$. This ensures $\mathcal{D}^+ = \mathcal{D}$ and that $\tilde{\mathcal{L}}$ is invertible and lower triangular with $O(m \log^3 n \log(1/\epsilon))$ non-zeros. We note that the inverse of an invertible lower or upper triangular matrix with N non-zeros can be applied in time $O(N)$ given an adjacency list representation of the matrix. Finally, as $\ker(\mathcal{L}\mathcal{L}^\top) = \text{span}\{\mathbf{1}\}$, we have $\mathbf{\Pi}_{\tilde{\mathcal{L}}\mathcal{D}(\tilde{\mathcal{L}})^\top} = \mathbf{I} - \frac{1}{n}\mathbf{1}\mathbf{1}^\top$, and this projection matrix can be applied in $O(n)$ time. Altogether, this means that $T_{\text{matvec}} + T_{\text{sol}} + T_{\text{proj}} = O(n)$, which suffices to complete the proof. \square

3 Approximate Gaussian Elimination Algorithm

Recall *Gaussian Elimination / Cholesky decomposition* of a graph Laplacian \mathbf{L} . We will use $\mathbf{A}(:, i)$ to denote the the i th column of a matrix \mathbf{A} . We can write the algorithm as

Algorithm 1: Gaussian Elimination / Cholesky Decomposition

Input: Graph Laplacian \mathbf{L}

Output: Lower triangular \mathcal{L} s.t. $\mathcal{L}\mathcal{L}^\top = \mathbf{L}$

Let $\mathbf{S}_0 = \mathbf{L}$;

for $i = 1$ *to* $i = n - 1$ **do**

$l_i = \frac{1}{\sqrt{\mathbf{S}_{i-1}(i,i)}} \mathbf{S}_{i-1}(:, i);$
 $\mathbf{S}_i = \mathbf{S}_{i-1} - l_i l_i^\top.$

$l_n = \mathbf{0}_{n \times 1};$

return $\mathcal{L} = [l_1 \cdots l_n];$

We want to introduce some notation that will help us describe and analyze a faster version of Gaussian elimination – one that uses sampling to create a sparse approximation of the decomposition.

Consider a Laplacian \mathbf{S} of a graph H and a vertex v of H . We define $\text{STAR}(v, \mathbf{S})$ to be the Laplacian of the subgraph of H consisting of edges incident on v . We define

$$\text{CLIQUE}(v, \mathbf{S}) = \text{STAR}(v, \mathbf{S}) - \frac{1}{\mathbf{S}(v, v)} \mathbf{S}(:, v) \mathbf{S}(:, v)^\top$$

For example, suppose

$$\mathbf{L} = \begin{pmatrix} W & -\mathbf{a}^\top \\ -\mathbf{a} & \text{diag}(\mathbf{a}) + \mathbf{L}_{-1} \end{pmatrix}$$

Then

$$\text{STAR}(1, \mathbf{L}) = \begin{pmatrix} W & -\mathbf{a}^\top \\ -\mathbf{a} & \text{diag}(\mathbf{a}) \end{pmatrix} \text{ and } \text{CLIQUE}(1, \mathbf{L}) = \begin{pmatrix} 0 & \mathbf{0} \\ \mathbf{0} & \text{diag}(\mathbf{a}) - \frac{1}{W} \mathbf{a} \mathbf{a}^\top \end{pmatrix}$$

which is illustrated in Figure 1.

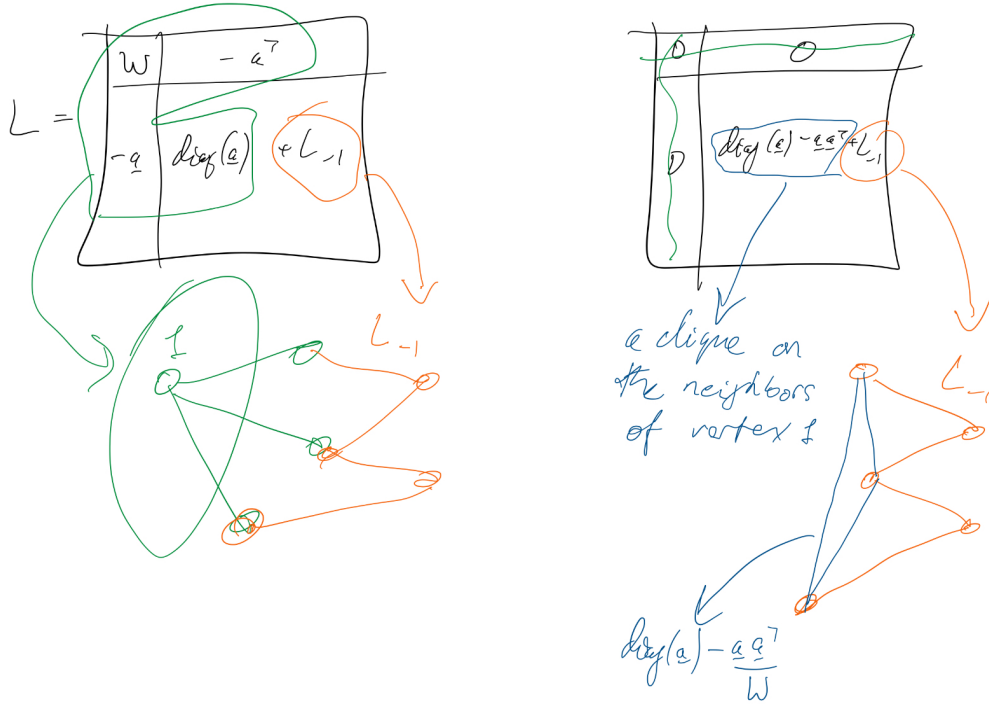


Figure 1: Gaussian Elimination: $\text{CLIQUE}(1, L) = \text{STAR}(1, L) - \frac{1}{L(1,1)} L(:, 1) L(:, 1)^T$.

In Lecture 7, we proved that $\text{CLIQUE}(v, \mathbf{S})$ is a graph Laplacian – it follows from the proof of Claim 1.1 in that lecture. Thus we have that following.

Claim 3.1. *If \mathbf{S} is the Laplacian of a connected graph, then $\text{CLIQUE}(v, \mathbf{S})$ is a graph Laplacian.*

Note that in Algorithm 1, we have $\mathbf{l}_i \mathbf{l}_i^T = \text{STAR}(v_i, \mathbf{S}_{i-1}) - \text{CLIQUE}(v_i, \mathbf{S}_{i-1})$. The update rule can be rewritten as

$$\mathbf{S}_i = \mathbf{S}_{i-1} - \text{STAR}(v_i, \mathbf{S}_{i-1}) + \text{CLIQUE}(v_i, \mathbf{S}_{i-1}),$$

This also provides way to understand why Gaussian Elimination is slow in some cases. At each step, one vertex is eliminated, but a clique is added to the subgraph on the remaining vertices, making the graph denser. And at the i th step, computing $\text{STAR}(v_i, \mathbf{S}_{i-1})$ takes around $\text{deg}(v_i)$ time, but computing $\text{CLIQUE}(v_i, \mathbf{S}_{i-1})$ requires around $\text{deg}(v_i)^2$ time. In order to speed up Gaussian Elimination, the algorithmic idea of [KS16] is to plug in a sparser approximate of the intended clique instead of the entire one.

The following procedure $\text{CLIQUESAMPLE}(v, \mathbf{S})$ produces a sparse approximation of $\text{CLIQUE}(v, \mathbf{S})$. Let V be the vertex set of the graph associated with \mathbf{S} and E the edge set. We define $\mathbf{b}_{i,j} \in \mathbb{R}^V$ to be the vector with

$$\mathbf{b}_{i,j}(i) = 1 \text{ and } \mathbf{b}_{i,j}(j) = -1 \text{ and } \mathbf{b}_{i,j}(k) = 0 \text{ for } k \neq i, j.$$

Given weights $\mathbf{w} \in \mathbb{R}^E$ and a vertex $v \in V$, we let

$$\mathbf{w}_v = \sum_{(u,v) \in E} \mathbf{w}(u, v).$$

Algorithm 2: CLIQUESAMPLE(v, \mathbf{S})

Input: Graph Laplacian $\mathbf{S} \in \mathbb{R}^{V \times V}$, of a graph with edge weights \mathbf{w} , and vertex $v \in V$

Output: $\mathbf{Y}_v \in \mathbb{R}^{V \times V}$ sparse approximation of CLIQUE(v, \mathbf{S})

$\mathbf{Y}_v \leftarrow \mathbf{0}_{n \times n}$;

foreach Multiedge $e = (v, i)$ from v to a neighbor i **do**

 Randomly pick a neighbor j of v with probability $\frac{\mathbf{w}(j, v)}{\mathbf{w}_v}$;

 If $i \neq j$, let $\mathbf{Y}_v \leftarrow \mathbf{Y}_v + \frac{\mathbf{w}(i, v)\mathbf{w}(j, v)}{\mathbf{w}(i, v) + \mathbf{w}(j, v)} \mathbf{b}_{i, j} \mathbf{b}_{i, j}^\top$;

return \mathbf{Y}_v ;

Remark 3.2. We can implement each sampling of a neighbor j in $O(1)$ time using a classical algorithm known as Walker's method (also known as the Alias method or Vose's method). This algorithm requires an additional $O(\deg_{\mathbf{S}}(v))$ time to initialize a data structure used for sampling. Overall, this means the total time for $O(\deg_{\mathbf{S}}(v))$ samples is still $O(\deg_{\mathbf{S}}(v))$.

Lemma 3.3. $\mathbb{E}[\mathbf{Y}_v] = \text{CLIQUE}(v, \mathbf{S})$.

Proof. Let $\mathbf{C} = \text{CLIQUE}(v, \mathbf{S})$. Observe that both $\mathbb{E}[\mathbf{Y}_v]$ and \mathbf{C} are Laplacians. Thus it suffices to verify $\mathbb{E} \mathbf{Y}_v(i, j) = \mathbf{C}(i, j)$ for $i \neq j$.

$$\mathbf{C}(i, j) = -\frac{\mathbf{w}(i, v)\mathbf{w}(j, v)}{\mathbf{w}_v},$$
$$\mathbb{E} \mathbf{Y}_v(i, j) = -\frac{\mathbf{w}(i, v)\mathbf{w}(j, v)}{\mathbf{w}(i, v) + \mathbf{w}(j, v)} \left(\frac{\mathbf{w}(j, v)}{\mathbf{w}_v} + \frac{\mathbf{w}(i, v)}{\mathbf{w}_v} \right) = -\frac{\mathbf{w}(i, v)\mathbf{w}(j, v)}{\mathbf{w}_v} = -\mathbf{C}(i, j).$$

□

Remark 3.4. Lemma 3.3 shows that CLIQUESAMPLE(v, \mathbf{L}) produces the original CLIQUE(v, \mathbf{L}) in expectation.

Now, we define *Approximate Gaussian Elimination*.

Algorithm 3: Approximate Gaussian Elimination / Cholesky Decomposition

Input: Graph Laplacian \mathbf{L}

Output: Lower triangular^a \mathcal{L} as given in Theorem 2.4

Let $\mathbf{S}_0 = \mathbf{L}$;

Generate a random permutation π on $[n]$;

for $i = 1$ to $i = n - 1$ **do**

$\mathbf{l}_i = \frac{1}{\sqrt{\mathbf{S}_{i-1}(\pi(i), \pi(i))}} \mathbf{S}_{i-1}(:, \pi(i))$;

$\mathbf{S}_i = \mathbf{S}_{i-1} - \text{STAR}(\pi(i), \mathbf{S}_{i-1}) + \text{CLIQUE}(\pi(i), \mathbf{S}_{i-1})$

$\mathbf{l}_n = \mathbf{0}_{n \times 1}$;

return $\mathcal{L} = [\mathbf{l}_1 \cdots \mathbf{l}_n]$ and π ;

^a \mathcal{L} is not actually lower triangular. However, if we let \mathbf{P}_π be the permutation matrix corresponding to π , then $\mathbf{P}_\pi \mathcal{L}$ is lower triangular. Knowing the ordering that achieves this is enough to let us implement forward and backward substitution for solving linear equations in \mathcal{L} and \mathcal{L}^\top .

Note that if we replace CLIQUESAMPLE($\pi(i), \mathbf{S}_{i-1}$) by CLIQUE($\pi(i), \mathbf{S}_{i-1}$) at each step, then we can recover Gaussian Elimination, but with a random elimination order.

4 Analyzing Approximate Gaussian Elimination

In this Section, we're going to analyze Approximate Gaussian Elimination, and see why it works.

Ultimately, the main challenge in proving Theorem 2.4 will be to prove for the output \mathcal{L} of Algorithm 3 that with high probability

$$0.5\mathbf{L} \preceq \mathcal{L}\mathcal{L}^\top \preceq 1.5\mathbf{L}. \quad (1)$$

We can reduce this to proving that with high probability

$$\left\| \mathbf{L}^{+1/2}(\mathcal{L}\mathcal{L}^\top - \mathbf{L})\mathbf{L}^{+1/2} \right\| \leq 0.5 \quad (2)$$

Ultimately, the proof is going to have a lot in common with our proof of Matrix Bernstein in Lecture 8. Overall, the lesson there was that when we have a sum of independent, zero-mean random matrices, we can show that the sum is likely to have small spectral norm if the spectral norm of each random matrix is small, and the matrix-valued variance is also small.

Thus, to replicate the proof, we need control over

1. The *sample norms*.
2. The *sample variance*.

But, there is seemingly another major obstacle: We are trying to analyze a process where the samples are far from independent. Each time we sample edges, we add new edges to the remaining graph, which we will later sample again. This creates a lot of dependencies between the samples, which we have to handle.

However, it turns out that independence is more than what is needed to prove concentration. Instead, it suffices to have a sequence of random variables such that each is mean-zero in expectation, conditional on the previous ones. This is called a martingale difference sequence. We'll now learn about those.

4.1 Normalization, a.k.a. Isotropic Position

Since our analysis requires frequently measuring matrices after right and left-multiplication by $\mathbf{L}^{+1/2}$, we reintroduce the “normalizing map” $\Phi : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^{n \times n}$ defined by

$$\Phi(\mathbf{A}) = \mathbf{L}^{+1/2} \mathbf{A} \mathbf{L}^{+1/2}.$$

We previously saw this in Lectures 7 and 8.

4.2 Martingales

A scalar martingale is a sequence of random variables Z_0, \dots, Z_k , such that

$$\mathbb{E}[Z_i \mid Z_0, \dots, Z_{i-1}] = Z_{i-1}. \quad (3)$$

That is, conditional on the outcome of all the previous random variables, the expectation of Z_i equals Z_{i-1} . If we unravel the sequence of conditional expectations, we get that *without conditioning* $\mathbb{E}[Z_k] = \mathbb{E}[Z_0]$.

Typically, we use martingales to show a statement along like “ Z_k is concentrated around $\mathbb{E}[Z_k]$ ”.

We can also think of a martingale in terms of the sequence of changes in the Z_i variables. Let $X_i = Z_i - Z_{i-1}$. The sequence of X_i s is called a martingale difference sequence. We can now state the martingale condition as

$$\mathbb{E}[X_i \mid Z_0, \dots, Z_{i-1}] = 0.$$

And because Z_0 and X_1, \dots, X_{i-1} completely determine Z_1, \dots, Z_{i-1} , we could also write the martingale condition equivalently as

$$\mathbb{E}[X_i \mid Z_0, X_1, \dots, X_{i-1}] = 0.$$

Crucially, we can write

$$Z_k = Z_0 + \sum_{i=1}^k Z_i - Z_{i-1} = Z_0 + \sum_{i=1}^k X_i$$

and when we are trying to prove concentration, the martingale difference property of the X_i 's is often “as good as” independence, meaning that $\sum_{i=1}^k X_i$ concentrates similarly to a sum of independent random variables.

Matrix-valued martingales. We can also define matrix-valued martingales. In this case, we replace the martingale condition of Equation (3), with the condition that the whole matrix stays the same in expectation. For example, we could have a sequence of random matrices $\mathbf{Z}_0, \dots, \mathbf{Z}_k \in \mathbb{R}^{n \times n}$, such that

$$\mathbb{E}[\mathbf{Z}_i \mid \mathbf{Z}_0, \dots, \mathbf{Z}_{i-1}] = \mathbf{Z}_{i-1}. \quad (4)$$

Lemma 4.1. *Let $\mathbf{L}_i = \mathbf{S}_i + \sum_{j=1}^i \mathbf{l}_j \mathbf{l}_j^\top$ for $i = 1, \dots, n$ and $\mathbf{L}_0 = \mathbf{S}_0 = \mathbf{L}$. Then*

$$\mathbb{E}[\mathbf{L}_i \mid \text{all random variables before } \text{CLIQUESAMPLE}(\pi(i), \mathbf{S}_{i-1})] = \mathbf{L}_{i-1}.$$

Proof. Let's only consider $i = 1$ here as other cases are similar.

$$\begin{aligned} \mathbf{L}_0 &= \mathbf{L} = \mathbf{l}_1 \mathbf{l}_1^\top + \text{CLIQUE}(v, \mathbf{L}) + \mathbf{L}_{-1} \\ \mathbf{L}_1 &= \mathbf{l}_1 \mathbf{l}_1^\top + \text{CLIQUESAMPLE}(v, \mathbf{L}) + \mathbf{L}_{-1} \\ \mathbb{E}[\mathbf{L}_1 \mid \pi(1)] &= \mathbf{l}_1 \mathbf{l}_1^\top + \mathbb{E}[\text{CLIQUESAMPLE}(v, \mathbf{L}) \mid \pi(1)] + \mathbf{L}_{-1} \\ &= \mathbf{l}_1 \mathbf{l}_1^\top + \text{CLIQUE}(v, \mathbf{L}) + \mathbf{L}_{-1} \\ &= \mathbf{L}_0 \end{aligned}$$

where we used Lemma 3.3 to get $\mathbb{E}[\text{CLIQUESAMPLE}(v, \mathbf{L}) \mid \pi(1)] = \text{CLIQUE}(v, \mathbf{L})$. □

Remark 4.2. $\sum_{j=1}^i \mathbf{l}_j \mathbf{l}_j^\top$ can be treated as what has already been eliminated by (Approximate) Gaussian Elimination, while \mathbf{S}_i is what still left or going to be eliminated. In Approximate Gaussian Elimination, $\mathbf{L}_n = \sum_{i=1}^n \mathbf{l}_i \mathbf{l}_i^\top$ and our goal is to show that $\mathbf{L}_n \approx_K \mathbf{L}$. Note that \mathbf{L}_i is always equal to the original Laplacian \mathbf{L} for all i in Gaussian Elimination. Lemma 4.1 demonstrates that $\mathbf{L}_0, \mathbf{L}_1, \dots, \mathbf{L}_n$ forms a matrix martingale.

Ultimately, our plan is to use this matrix martingale structure to show that “ \mathbf{L}_n is concentrated around \mathbf{L} ” in some appropriate sense. More precisely, the spectral approximation we would like to show can be established by showing that “ $\Phi(\mathbf{L}_n)$ is concentrated around $\Phi(\mathbf{L})$ ”

4.3 Martingale Difference Sequence as Edge-Samples

We start by taking a slightly different view of the observations we used to prove Lemma 4.1. Recall that $\mathbf{L}_i = \mathbf{S}_i + \sum_{j=1}^i \mathbf{l}_j \mathbf{l}_j^\top$, and $\mathbf{L}_{i-1} = \mathbf{S}_{i-1} + \sum_{j=1}^{i-1} \mathbf{l}_j \mathbf{l}_j^\top$ and

$$\mathbf{S}_i = \mathbf{S}_{i-1} - \text{STAR}(\pi(i), \mathbf{S}_{i-1}) + \text{CLIQUESAMPLE}(\pi(i), \mathbf{S}_{i-1}).$$

Putting these together, we get

$$\begin{aligned} \mathbf{L}_i - \mathbf{L}_{i-1} &= \mathbf{l}_i \mathbf{l}_i^\top + \text{CLIQUESAMPLE}(\pi(i), \mathbf{S}_{i-1}) - \text{STAR}(\pi(i), \mathbf{S}_{i-1}) \\ &= \text{CLIQUESAMPLE}(\pi(i), \mathbf{S}_{i-1}) - \text{CLIQUE}(\pi(i), \mathbf{S}_{i-1}) \\ &= \text{CLIQUESAMPLE}(\pi(i), \mathbf{S}_{i-1}) - \mathbb{E}[\text{CLIQUESAMPLE}(\pi(i), \mathbf{S}_{i-1}) \mid \text{preceding samples}] \\ &\hspace{15em} \text{by Lemma 3.3.} \end{aligned} \tag{5}$$

In particular, recall that by Lemma 3.3, conditional on the randomness before the call to $\text{CLIQUESAMPLE}(\pi(i), \mathbf{S}_{i-1})$, we have

$$\mathbb{E}[\text{CLIQUESAMPLE}(\pi(i), \mathbf{S}_{i-1}) \mid \text{preceding samples}] = \text{CLIQUE}(\pi(i), \mathbf{S}_{i-1})$$

Adopting the notation of Lemma 3.3 we write

$$\mathbf{Y}_{\pi(i)} = \text{CLIQUESAMPLE}(\pi(i), \mathbf{S}_{i-1})$$

and we further introduce notation each multi-edge sample for $e \in \text{STAR}(\pi(i), \mathbf{S}_{i-1})$, as $\mathbf{Y}_{\pi(i),e}$, denoting the random edge Laplacian sampled when the algorithm is processing multi-edge e . Thus, conditional on preceding samples, we have

$$\mathbf{Y}_{\pi(i)} = \sum_{e \in \text{STAR}(\pi(i), \mathbf{S}_{i-1})} \mathbf{Y}_{\pi(i),e} \tag{6}$$

Note that even the number of multi-edges in $\text{STAR}(\pi(i), \mathbf{S}_{i-1})$ depends on the preceding samples. We also want to associate zero-mean variables with each edge. Conditional on preceding samples, we also define

$$\mathbf{X}_{i,e} = \Phi(\mathbf{Y}_{\pi(i),e} - \mathbb{E}[\mathbf{Y}_{\pi(i),e}]) \text{ and } \mathbf{X}_i = \sum_{e \in \text{STAR}(\pi(i), \mathbf{S}_{i-1})} \mathbf{X}_{i,e}$$

and combining this with Equations (5) and (6)

$$\mathbf{X}_i = \Phi(\mathbf{Y}_{\pi(i)} - \mathbb{E}[\mathbf{Y}_{\pi(i)}]) = \Phi(\mathbf{L}_i - \mathbf{L}_{i-1})$$

Altogether, we can write

$$\Phi(\mathbf{L}_n - \mathbf{L}) = \sum_{i=1}^n \Phi(\mathbf{L}_i - \mathbf{L}_{i-1}) = \sum_{i=1}^n \mathbf{X}_i = \sum_{i=1}^n \sum_{e \in \text{STAR}(\pi(i), \mathbf{S}_{i-1})} \mathbf{X}_{i,e}$$

Note that the $\mathbf{X}_{i,e}$ variables form a martingale difference sequence, because the linearity of Φ ensures they are zero-mean conditional on preceding randomness.

4.4 Stopped Martingales

Unfortunately, directly analyzing the concentration properties of the \mathbf{L}_i martingale that we just introduced turns out to be difficult. The reason is that we're trying to prove some very delicate multiplicative error guarantees. And, if we analyze \mathbf{L}_i , we find that the multiplicative error is not easy to control, *after it's already gotten big*. But that's not really what we care about anyway: We want to say it never gets big in the first place, with high probability. So we need to introduce another martingale, that lets us ignore the bad case when the error has already gotten too big. At the same time, we also need to make sure that statements about our new martingale can help us prove guarantees about \mathbf{L}_i . Fortunately, we can achieve both at once. The technique we use is related to the much broader topic of martingale *stopping times*, which we only scratch the surface of here. We're also going to be quite informal about it, in the interest of brevity. Lecture notes by Tropp [Tro19] give a more formal introduction for those who are interested.

We define the stopped martingale sequence $\tilde{\mathbf{L}}_i$ by

$$\tilde{\mathbf{L}}_i = \begin{cases} \mathbf{L}_i & \text{if for all } j < i \text{ we have } \mathbf{L}_j \preceq 1.5\mathbf{L} \\ \mathbf{L}_{j^*} & \text{for } j^* \text{ being the least } j \text{ such that } \mathbf{L}_j \not\preceq 1.5\mathbf{L} \end{cases} \quad (7)$$

Figure 2 shows the $\tilde{\mathbf{L}}_i$ martingale getting stuck at the first time $\mathbf{L}_{j^*} \not\preceq 1.5\mathbf{L}$.

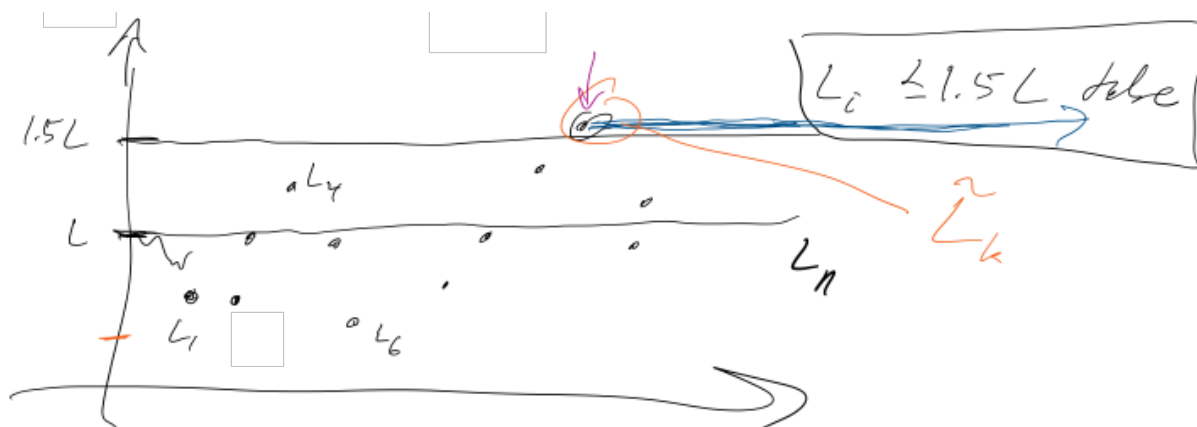


Figure 2: Gaussian Elimination : $\text{CLIQUE}(1, \mathbf{L}) = \text{STAR}(1, \mathbf{L}) - \frac{1}{L(1,1)} \mathbf{L}(:, 1) \mathbf{L}(:, 1)^\top$.

We state the following without proof:

Claim 4.3.

1. The sequence $\{\tilde{\mathbf{L}}_i\}$ for $i = 0, \dots, n$ is a martingale.
2. $\left\| \mathbf{L}^{+1/2} (\tilde{\mathbf{L}}_i - \mathbf{L}) \mathbf{L}^{+1/2} \right\| \leq 0.5$ implies $\left\| \mathbf{L}^{+1/2} (\mathbf{L}_i - \mathbf{L}) \mathbf{L}^{+1/2} \right\| \leq 0.5$

The martingale property also implies that the unconditional expectation satisfies $\mathbb{E}[\tilde{\mathbf{L}}_n] = \mathbf{L}$. The proof of the claim is easy to sketch: For Part 1, each difference is zero-mean if the condition has not been violated, and is identically zero (and hence zero-mean) if it has been violated. For Part 2,

if the martingale $\{\tilde{\mathbf{L}}_i\}$ has stopped, then $\left\| \mathbf{L}^{+2}(\tilde{\mathbf{L}}_i - \mathbf{L})\mathbf{L}^{+2} \right\| \leq 0.5$ is false, and the implication is vacuously true. If the, on the other hand, if the martingale has not stopped, the quantities are equal, because $\tilde{\mathbf{L}}_i = \mathbf{L}_i$, and again it's easy to see the implication holds.

Thus, ultimately, our strategy is goin to be to show that $\left\| \mathbf{L}^{+2}(\tilde{\mathbf{L}}_i - \mathbf{L})\mathbf{L}^{+2} \right\| \leq 0.5$ with high probability. Expressed using the normalizing map $\Phi(\cdot)$, our goal is to show that with high probability

$$\left\| \Phi(\tilde{\mathbf{L}}_n - \mathbf{L}) \right\| \leq 0.5.$$

Stopped martingale difference sequence. In order to prove the spectral norm bound, we want to express the $\{\tilde{\mathbf{L}}_i\}$ martingale in terms of a sequence of martingale differences. To this end, we define $\tilde{\mathbf{X}}_i = \Phi(\tilde{\mathbf{L}}_i - \tilde{\mathbf{L}}_{i-1})$. This ensures that

$$\tilde{\mathbf{X}}_i = \begin{cases} \mathbf{X}_i & \text{if for all } j < i \text{ we have } \mathbf{L}_j \preceq 1.5\mathbf{L} \\ \mathbf{0} & \text{otherwise} \end{cases} \quad (8)$$

Whenever the modified martingale $\tilde{\mathbf{X}}_i$ has not yet stopped, we also introduce individual modified edge samples $\tilde{\mathbf{X}}_{i,e} = \mathbf{X}_{i,e}$. If the martingale *has* stopped, i.e. $\tilde{\mathbf{X}}_i = \mathbf{0}$, then we can take these edge samples $\tilde{\mathbf{X}}_{i,e}$ to be zero. We can now write

$$\Phi(\tilde{\mathbf{L}}_n - \mathbf{L}) = \sum_{i=1}^n \Phi(\tilde{\mathbf{L}}_i - \tilde{\mathbf{L}}_{i-1}) = \sum_{i=1}^n \tilde{\mathbf{X}}_i = \sum_{i=1}^n \sum_{e \in \text{STAR}(\pi(i), \mathbf{S}_{i-1})} \tilde{\mathbf{X}}_{i,e}.$$

Thus, we can see that Equation (2) is implied by

$$\left\| \sum_{i=1}^n \tilde{\mathbf{X}}_i \right\| \leq 0.5. \quad (9)$$

4.5 Sample Norm Control

In this Subsection, we're going to see that the norms of each multi-edge sample is controlled throughout the algorithm.

Lemma 4.4. *Given two Laplacians \mathbf{L} and \mathbf{S} on the same vertex set.¹ If each multiedge e of $\text{STAR}(v, \mathbf{S})$ has bounded norm in the following sense,*

$$\left\| \mathbf{L}^{+2} \mathbf{w}_{\mathbf{S}}(e) \mathbf{b}_e \mathbf{b}_e^\top \mathbf{L}^{+2} \right\| \leq R,$$

then each possible sampled multiedge e' of $\text{CLIQUE_SAMPLE}(v, \mathbf{S})$ also satisfies

$$\left\| \mathbf{L}^{+2} \mathbf{w}_{\text{new}}(e') \mathbf{b}_{e'} \mathbf{b}_{e'}^\top \mathbf{L}^{+2} \right\| \leq R.$$

¹ \mathbf{L} can be regarded as the original Laplacian we care about, while \mathbf{S} can be regarded as some intermediate Laplacian appearing during Approximate Gaussian Elimination.

Proof. Let $\mathbf{w} = \mathbf{w}_S$ for simplicity. Consider a sampled edge between i and j with weight $\mathbf{w}_{\text{new}}(i, j) = \mathbf{w}(i, v)\mathbf{w}(j, v)/(\mathbf{w}(i, v) + \mathbf{w}(j, v))$.

$$\begin{aligned}
\left\| \mathbf{L}^{+1/2} \mathbf{w}_{\text{new}}(i, j) \mathbf{b}_{ij} \mathbf{b}_{ij}^\top \mathbf{L}^{+1/2} \right\| &= \mathbf{w}_{\text{new}}(i, j) \left\| \mathbf{L}^{+1/2} \mathbf{b}_{ij} \mathbf{b}_{ij}^\top \mathbf{L}^{+1/2} \right\| \\
&= \mathbf{w}_{\text{new}}(i, j) \left\| \mathbf{L}^{+1/2} \mathbf{b}_{ij} \right\|^2 \\
&\leq \mathbf{w}_{\text{new}}(i, j) \left(\left\| \mathbf{L}^{+1/2} \mathbf{b}_{iv} \right\|^2 + \left\| \mathbf{L}^{+1/2} \mathbf{b}_{jv} \right\|^2 \right) \\
&= \frac{\mathbf{w}(j, v)}{\mathbf{w}(i, v) + \mathbf{w}(j, v)} \left\| \mathbf{L}^{+1/2} \mathbf{w}(i, v) \mathbf{b}_{iv} \mathbf{b}_{iv}^\top \mathbf{L}^{+1/2} \right\| + \\
&\quad \frac{\mathbf{w}(i, v)}{\mathbf{w}(i, v) + \mathbf{w}(j, v)} \left\| \mathbf{L}^{+1/2} \mathbf{w}(j, v) \mathbf{b}_{jv} \mathbf{b}_{jv}^\top \mathbf{L}^{+1/2} \right\| \\
&\leq \frac{\mathbf{w}(j, v)}{\mathbf{w}(i, v) + \mathbf{w}(j, v)} R + \frac{\mathbf{w}(i, v)}{\mathbf{w}(i, v) + \mathbf{w}(j, v)} R \\
&= R
\end{aligned}$$

The first inequality uses the triangle inequality of effective resistance in \mathbf{L} , in that effective resistance is a distance as proved in lecture 6. The second inequality just uses the conditions of this lemma. \square

Remark 4.5. Lemma 4.4 only requires that each single multiedge has small norm instead of that the sum of all edges between a pair of vertices have small norm. And this lemma tells us, after sampling, each multiedge in the new graph still satisfies the bounded norm condition.

From the Lemma, we can conclude that each edge sample $\mathbf{Y}_{\pi(i),e}$ satisfies $\|\Phi(\mathbf{Y}_{\pi(i),e})\| \leq R$ provided the assumptions of the Lemma hold. Let's record this observation as a Lemma.

Lemma 4.6. *If for all $e \in \text{STAR}(v, \mathbf{S}_i)$,*

$$\left\| \Phi(\mathbf{w}_{\mathbf{S}_i}(e) \mathbf{b}_e \mathbf{b}_e^\top) \right\| \leq R.$$

then all $e \in \text{STAR}(\pi(i), \mathbf{S}_i)$,

$$\left\| \Phi(\mathbf{Y}_{\pi(i),e}) \right\| \leq R.$$

Preprocessing by multi-edge splitting. In the original graph of Laplacian \mathbf{L} of graph $G = (V, E, \mathbf{w})$, we have for each edge \hat{e} that

$$\mathbf{w}(\hat{e}) \mathbf{b}_{\hat{e}} \mathbf{b}_{\hat{e}}^\top \preceq \sum_e \mathbf{w}(e) \mathbf{b}_e \mathbf{b}_e^\top = \mathbf{L}$$

This also implies that

$$\left\| \mathbf{L}^{+1/2} \mathbf{w}(\hat{e}) \mathbf{b}_{\hat{e}} \mathbf{b}_{\hat{e}}^\top \mathbf{L}^{+1/2} \right\| \leq 1.$$

Now, that means that if we split every original edge e of the graph into K multi-edges e_1, \dots, e_K , with a fraction $1/K$ of the weight, we get a new graph $G' = (V, E', \mathbf{w}')$ such that

Claim 4.7.

1. G' and G have the same graph Laplacian.

2. $|E'| = K |E|$
3. For every multi-edge in G'

$$\left\| \mathbf{L}^{+/2} \mathbf{w}'(e) \mathbf{b}_e \mathbf{b}_e^\top \mathbf{L}^{+/2} \right\| \leq 1/K.$$

Before we run Approximate Gaussian Elimination, we are going to do this multi-edge splitting to ensure we have control over multi-edge sample norms. Combined with Lemma 4.4 immediately establishes the next lemma, because we start off with all multi-edges having bounded norm and only produce multi-edges with bounded norm.

Lemma 4.8. *When Algorithm 3 is run on the (multi-edge) Laplacian of G' , arising from splitting edges of G into K multi-edges, the every edge sample $\mathbf{Y}_{\pi(i),e}$ satisfies*

$$\|\Phi(\mathbf{Y}_{\pi(i),e})\| \leq 1/K.$$

As we will see later $K = 200 \log^2 n$ suffices.

4.6 Random Matrix Concentration from Trace Exponentials

Let us recall how matrix-valued variances come into the picture when proving concentration following the strategy from Matrix Bernstein in Lecture 8.

For some matrix-valued random variable $\mathbf{X} \in S^n$, we'd like to show $\Pr[\|\mathbf{X}\| \leq 0.5]$. Using Markov's inequality, and some observations about matrix exponentials and traces, we saw that for all $\theta > 0$,

$$\Pr[\|\mathbf{X}\| \geq 0.5] \leq \exp(-0.5\theta) (\mathbb{E}[\text{Tr}(\exp(\theta\mathbf{X}))] + \mathbb{E}[\text{Tr}(\exp(-\theta\mathbf{X}))]). \quad (10)$$

We then want to bound $\mathbb{E}[\text{Tr}(\exp(\theta\mathbf{X}))]$ using Lieb's theorem. We can handle $\mathbb{E}[\text{Tr}(\exp(-\theta\mathbf{X}))]$ similarly.

Theorem 4.9 (Lieb). *Let $f : S_{++}^n \rightarrow \mathbb{R}$ be a matrix function given by*

$$f(\mathbf{A}) = \text{Tr}(\exp(\mathbf{H} + \log(\mathbf{A})))$$

for some $\mathbf{H} \in S^n$. Then $-f$ is convex (i.e. f is concave).

As observed by Tropp, this is useful for proving matrix concentration statements. Combined with Jensen's inequality, it gives that for a random matrix $\mathbf{X} \in S^n$ and a fixed $\mathbf{H} \in S^n$

$$\mathbb{E}[\text{Tr}(\exp(\mathbf{H} + \mathbf{X}))] \leq \text{Tr}(\exp(\mathbf{H} + \log(\mathbb{E}[\exp(\mathbf{X})]))).$$

The next crucial step was to show that it suffices to obtain an upper bound on the matrix $\mathbb{E}[\exp(\mathbf{X})]$ w.r.t the Loewner order. Using the following three lemmas, this conclusion is an immediate corollary.

Lemma 4.10. *If $\mathbf{A} \preceq \mathbf{B}$, then $\text{Tr}(\exp(\mathbf{A})) \leq \text{Tr}(\exp(\mathbf{B}))$.*

Lemma 4.11. *If $0 \prec \mathbf{A} \preceq \mathbf{B}$, then $\log(\mathbf{A}) \preceq \log(\mathbf{B})$.*

Lemma 4.12. $\log(\mathbf{I} + \mathbf{A}) \preceq \mathbf{A}$ for $\mathbf{A} \succ -\mathbf{I}$.

Corollary 4.13. *For a random matrix $\mathbf{X} \in S^n$ and a fixed $\mathbf{H} \in S^n$, if $\mathbb{E}[\exp(\mathbf{X})] \preceq \mathbf{I} + \mathbf{U}$ where $\mathbf{U} \succ -\mathbf{I}$, then*

$$\mathbb{E}[\text{Tr}(\exp(\mathbf{H} + \mathbf{X}))] \leq \text{Tr}(\exp(\mathbf{H} + \mathbf{U})).$$

4.7 Mean-Exponential Bounds from Variance Bounds

To use Corollary 4.13, we need to construct useful upper bounds on $\mathbb{E}[\exp(\mathbf{X})]$. This can be done, starting from the following lemma.

Lemma 4.14. $\exp(\mathbf{A}) \preceq \mathbf{I} + \mathbf{A} + \mathbf{A}^2$ for $\|\mathbf{A}\| \leq 1$.

If \mathbf{X} is zero-mean and $\|\mathbf{X}\| \leq 1$, this means that $\mathbb{E}[\exp(\mathbf{X})] \preceq \mathbf{I} + \mathbb{E}[\mathbf{X}^2]$, which is how we end up wanting to bound the matrix-valued variance $\mathbb{E}[\mathbf{X}^2]$. In the rest of this Subsection, we're going to see the matrix-valued variance of the stopped martingale is bounded throughout the algorithm.

Firstly, we note that for a single edge sample $\tilde{\mathbf{X}}_{i,e}$, by Lemma 4.8, we have that

$$\|\tilde{\mathbf{X}}_{i,e}\| \leq \|\Phi(\mathbf{Y}_{\pi(i),e} - \mathbb{E}[\mathbf{Y}_{\pi(i),e}])\| \leq 1/K,$$

using that $\|\mathbf{A} - \mathbf{B}\| \leq \max(\|\mathbf{A}\|, \|\mathbf{B}\|)$, for $\mathbf{A}, \mathbf{B} \succeq \mathbf{0}$, and $\|\mathbb{E}[\mathbf{A}]\| \leq \mathbb{E}[\|\mathbf{A}\|]$ by Jensen's inequality.

Thus, if $0 < \theta \leq K$, we have that

$$\begin{aligned} \mathbb{E} \left[\exp(\theta \tilde{\mathbf{X}}_{i,e}) \mid \text{preceding samples} \right] &\preceq \mathbf{I} + \mathbb{E} \left[(\theta \tilde{\mathbf{X}}_{i,e})^2 \mid \text{preceding samples} \right] \\ &\preceq \mathbf{I} + \frac{1}{K} \theta^2 \cdot \mathbb{E} \left[\Phi(\mathbf{Y}_{\pi(i),e}) \mid \text{preceding samples} \right] \end{aligned} \quad (11)$$

4.8 The Overall Mean-Trace-Exponential Bound

We will use $\mathbb{E}_{(<i)}$ to denote expectation over variables preceding the i th elimination step. We are going to refrain from explicitly writing out conditioning in our expectations, but any *inner* expectation that appears inside another *outer* expectation should be taken as conditional on the outer expectation. We are going to use d_i to denote the multi-edge degree of vertex $\pi(i)$ in \mathcal{S}_{i-1} . This is exactly the number of edge samples in the i th elimination. Note that there is no elimination at step n (the algorithm is already finished). As a notational convenience, let's write $\hat{n} = n - 1$. With all that in mind, we bound the mean-trace-exponential for some parameter $0 < \theta \leq 0.5/\sqrt{K}$

$$\mathbb{E} \text{Tr} \left(\exp\left(\theta \sum_{i=1}^{\hat{n}} \tilde{\mathbf{X}}_i\right) \right) \quad (12)$$

$$\begin{aligned} &= \mathbb{E}_{(<\hat{n})} \mathbb{E}_{\pi(\hat{n})} \mathbb{E}_{\tilde{\mathbf{X}}_{\hat{n},1}} \cdots \mathbb{E}_{\tilde{\mathbf{X}}_{\hat{n},d_{\hat{n}}-1}} \mathbb{E}_{\tilde{\mathbf{X}}_{\hat{n},d_{\hat{n}}}} \text{Tr} \exp \left(\underbrace{\sum_{i=1}^{\hat{n}-1} \theta \tilde{\mathbf{X}}_i + \sum_{e=1}^{d_{\hat{n}}-1} \theta \tilde{\mathbf{X}}_{\hat{n},e} + \theta \tilde{\mathbf{X}}_{\hat{n},d_{\hat{n}}}}_{\mathbf{H}} \right) \\ &\quad \tilde{\mathbf{X}}_{\hat{n},1}, \dots, \tilde{\mathbf{X}}_{\hat{n},d_{\hat{n}}} \text{ are independent conditional on } (<\hat{n}), \pi(\hat{n}) \\ &\leq \mathbb{E}_{(<\hat{n})} \mathbb{E}_{\pi(\hat{n})} \mathbb{E}_{\tilde{\mathbf{X}}_{\hat{n},1}} \cdots \mathbb{E}_{\tilde{\mathbf{X}}_{\hat{n},d_{\hat{n}}-1}} \text{Tr} \exp \left(\sum_{i=1}^{\hat{n}-1} \theta \tilde{\mathbf{X}}_i + \sum_{e=1}^{d_{\hat{n}}-1} \theta \tilde{\mathbf{X}}_{\hat{n},e} + \frac{1}{K} \theta^2 \cdot \mathbb{E}_{\tilde{\mathbf{X}}_{\hat{n},d_{\hat{n}}}} \Phi(\mathbf{Y}_{\pi(\hat{n}),d_{\hat{n}}}) \right) \\ &\quad \text{By Equation (11) and Corollary 4.13 .} \end{aligned}$$

\vdots Repeat for each multi-edge sample $\tilde{\mathbf{X}}_{\hat{n},1}, \dots, \tilde{\mathbf{X}}_{\hat{n},d_{\hat{n}}-1}$

$$\begin{aligned}
&\leq \mathbb{E}_{\langle \hat{n} \rangle} \mathbb{E}_{\pi(\hat{n})} \operatorname{Tr} \exp \left(\sum_{i=1}^{\hat{n}-1} \theta \tilde{\mathbf{X}}_i + \sum_{e=1}^{d_{\hat{n}}} \frac{1}{K} \theta^2 \cdot \mathbb{E}_{\tilde{\mathbf{X}}_{\hat{n},e}} \Phi(\mathbf{Y}_{\pi(\hat{n}),e}) \right) \\
&= \mathbb{E}_{\langle \hat{n} \rangle} \mathbb{E}_{\pi(\hat{n})} \operatorname{Tr} \exp \left(\sum_{i=1}^{\hat{n}-1} \theta \tilde{\mathbf{X}}_i + \frac{1}{K} \theta^2 \operatorname{CLIQUE}(\pi(\hat{n}), \mathbf{S}_{\hat{n}-1}) \right)
\end{aligned}$$

To further bound this quantity, we now need to deal with the random choice of $\pi(\hat{n})$. We'll be able to use this to bound the trace-exponential in a very strong way. From a random matrix perspective, it's the following few steps that give the analysis its surprising strength.

We can treat $\frac{1}{K} \theta^2 \operatorname{CLIQUE}(\pi(\hat{n}), \mathbf{S}_{\hat{n}-1})$ as a random matrix. It is not zero-mean, but we can still bound the trace-exponential using Corollary 4.13.

We can also bound the expected matrix exponential in that case, using a simple corollary of Lemma 4.14.

Corollary 4.15. $\exp(\mathbf{A}) \preceq \mathbf{I} + (1 + R)\mathbf{A}$ for $\mathbf{0} \preceq \mathbf{A}$ with $\|\mathbf{A}\| \leq R$.

Proof. The conclusion follows after observing that for $\mathbf{0} \preceq \mathbf{A}$ with $\|\mathbf{A}\| \leq R$, we have $\mathbf{A}^2 \preceq R\mathbf{A}$. We can see this by considering the spectral decomposition of \mathbf{A} and dealing with each eigenvalue separately. \square

Next, we need a simple structural observation about the cliques created by elimination:

Claim 4.16.

$$\operatorname{CLIQUE}(\pi(i), \mathbf{S}_i) \preceq \operatorname{STAR}(\pi(i), \mathbf{S}_i) \preceq \mathbf{S}_i$$

Proof. The first inequality is immediate from $\operatorname{CLIQUE}(\pi(i), \mathbf{S}_i) \preceq \operatorname{CLIQUE}(\pi(i), \mathbf{S}_i) + \mathbf{l}_i \mathbf{l}_i^\top = \operatorname{STAR}(\pi(i), \mathbf{S}_i)$. The latter inequality $\operatorname{STAR}(\pi(i), \mathbf{S}_i) \preceq \mathbf{S}_i$ follows from the star being a subgraph of the whole Laplacian \mathbf{S}_i . \square

Next we make use of the fact that $\tilde{\mathbf{X}}_i$ is from the difference sequence of the *stopped* martingale. This means we can assume

$$\mathbf{S}_i \preceq 1.5\mathbf{L},$$

since otherwise $\tilde{\mathbf{X}}_i = \mathbf{0}$ and we get an even better bound on the trace-exponential. To make this formal, in Equation (12), we ought to do a case analysis that also includes the case $\tilde{\mathbf{X}}_i = \mathbf{0}$ when the martingale has stopped, but we omit this.

Thus we can conclude by Claim 4.16 that

$$\|\Phi(\operatorname{CLIQUE}(\pi(i), \mathbf{S}_i))\| \leq 1.5.$$

By our assumption $0 < \theta \leq 0.5/\sqrt{K}$, we have $\|\frac{1}{K} \theta^2 \Phi(\operatorname{CLIQUE}(\pi(i), \mathbf{S}_{i-1}))\| \leq 1$, so that by Corollary 4.15,

$$\begin{aligned}
\mathbb{E}_{\pi(i)} \exp \left(\frac{1}{K} \theta^2 \Phi(\operatorname{CLIQUE}(\pi(i), \mathbf{S}_{i-1})) \right) &\preceq \mathbf{I} + \frac{2}{K} \theta^2 \mathbb{E}_{\pi(i)} \Phi(\operatorname{CLIQUE}(\pi(i), \mathbf{S}_{i-1})) \quad (13) \\
&\preceq \mathbf{I} + \frac{2}{K} \theta^2 \mathbb{E}_{\pi(i)} \Phi(\operatorname{STAR}(\pi(i), \mathbf{S}_{i-1})) \quad \text{by Claim 4.16.}
\end{aligned}$$

Next we observe that, because every multi-edge appears in exactly two stars, and $\pi(i)$ is chosen uniformly at random among the $n + 1 - i$ vertices that \mathbf{S}_{i-1} is supported on, we have

$$\mathbb{E}_{\pi(i)} \text{STAR}(\pi(i), \mathbf{S}_{i-1}) = 2 \frac{1}{n+1-i} \mathbf{S}_{i-1}.$$

And, since we assume $\mathbf{S}_i \preceq 1.5\mathbf{L}$, we further get

$$\mathbb{E}_{\pi(i)} \exp \left(\frac{1}{K} \theta^2 \Phi(\text{CLIQUE}(\pi(i), \mathbf{S}_{i-1})) \right) \preceq \mathbf{I} + \frac{6\theta^2}{K(n+1-i)} \mathbf{I}.$$

We can combine this with Equation (12) and Corollary 4.13 to get

$$\begin{aligned} & \mathbb{E} \text{Tr} \left(\exp \left(\theta \sum_{i=1}^{\hat{n}} \tilde{\mathbf{X}}_i \right) \right) \\ & \leq \mathbb{E}_{\langle \hat{n} \rangle} \mathbb{E}_{\pi(\hat{n})} \text{Tr} \exp \left(\sum_{i=1}^{\hat{n}-1} \theta \tilde{\mathbf{X}}_i + \frac{1}{K} \theta^2 \text{CLIQUE}(\pi(\hat{n}), \mathbf{S}_{\hat{n}-1}) \right) \\ & \leq \mathbb{E}_{\langle \hat{n} \rangle} \text{Tr} \exp \left(\sum_{i=1}^{\hat{n}-1} \theta \tilde{\mathbf{X}}_i + \frac{6\theta^2}{K(n+1-i)} \mathbf{I} \right) \end{aligned}$$

And by repeating this analysis for each term $\tilde{\mathbf{X}}_i$, we get

$$\begin{aligned} \mathbb{E} \text{Tr} \left(\exp \left(\theta \sum_{i=1}^{\hat{n}} \tilde{\mathbf{X}}_i \right) \right) & \leq \text{Tr} \exp \left(\sum_{i=1}^{\hat{n}} \frac{6\theta^2}{K(n+1-i)} \mathbf{I} \right) \\ & \leq \text{Tr} \exp \left(\frac{7\theta^2 \log(n)}{K} \mathbf{I} \right) \\ & = n \exp \left(\frac{7\theta^2 \log(n)}{K} \right) \end{aligned}$$

Then, by choosing $K = 200 \log^2 n$ and $\theta = 0.5\sqrt{K}$, we get

$$\exp(-0.5\theta) \mathbb{E} \text{Tr} \left(\exp \left(\theta \sum_{i=1}^{\hat{n}} \tilde{\mathbf{X}}_i \right) \right) \leq \exp(-0.5\theta)n \exp \left(\frac{7\theta^2 \log(n)}{K} \right) \leq 1/n^5.$$

$\mathbb{E} \text{Tr} \left(\exp(-\theta \sum_{i=1}^{\hat{n}} \tilde{\mathbf{X}}_i) \right)$ can be bounded by an identical argument, so that Equation (10) gives

$$\Pr \left[\left\| \sum_{i=1}^{\hat{n}} \tilde{\mathbf{X}}_i \right\| \geq 0.5 \right] \leq 2/n^5.$$

Thus we have established $\left\| \sum_{i=1}^{\hat{n}} \tilde{\mathbf{X}}_i \right\| \leq 0.5$ with high probability (Equation (9)), and this in turn implies Equation (2), and finally Equation (1):

$$0.5\mathbf{L} \preceq \mathcal{L}\mathcal{L}^\top \preceq 1.5\mathbf{L}.$$

Now, all that's left to note is that the running time is linear in the multi-edge degree of the vertex being eliminated in each iteration (and this also bounds the number of non-zero entries being

created in \mathcal{L}). The total number of multi-edges left in the remaining graph stays constant at $Km = O(m \log^2 n)$. Thus the expected degree in the i th elimination is $Km/(n + i - 1)$, because the remaining number of vertices is $n + i - 1$. Hence the total running time and total number of non-zero entries created can both be bounded as

$$Km \sum_i 1/(n + i - 1) = O(m \log^3 n).$$

We can further prove that the bound $O(m \log^3 n)$ on running time and number of non-zeros in \mathbf{L} holds with high probability (e.g. $1 - 1/n^5$). To show this, we essentially need a scalar Chernoff bound, in except the degrees are in fact not independent, and so we need a scalar martingale concentration result, e.g. Azuma's Inequality. This way, we complete the proof of Theorem 2.4.

References

- [KS16] R. Kyng and S. Sachdeva. Approximate gaussian elimination for laplacians - fast, sparse, and simple. In *2016 IEEE 57th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 573–582, 2016.
- [ST04] Daniel A. Spielman and Shang-Hua Teng. Nearly-linear time algorithms for graph partitioning, graph sparsification, and solving linear systems. In *Proceedings of the Thirty-Sixth Annual ACM Symposium on Theory of Computing, STOC '04*, page 81–90, New York, NY, USA, 2004. Association for Computing Machinery.
- [Tro19] Joel A Tropp. Matrix concentration & computational linear algebra. 2019.