

Computing the log-determinant of symmetric, diagonally dominant matrices in near-linear time

WORKING DRAFT

Timothy Hunter

Ahmed El Alaoui

Alexandre M. Bayen

Department of Electrical Engineering and Computer Sciences

University of California

Berkeley, CA 94720-1776, USA

TJHUNTER@EECS.BERKELEY.EDU

ELALAOUI@EECS.BERKELEY.EDU

BAYEN@BERKELEY.EDU

Editor: Editor

Outline

1. Sampling log-determinants (note: this all works for PD matrices)
Should be moved to the end?
 - (a) Bernstein inequality
 - (b) Algo preconditioned monte carlo
 - (c) Sampling theorem
 - (d) Corollary with SDD
2. General method for SDD and Laplacian systems
 - (a) determinant of Laplacian
 - (b) idea of preconditioner for determinant
3. Using a tree preconditioner
 - (a) Results from Abraham
 - (b) Application to log-det
4. Using an ST preconditioner
 - (a) Explanation from Koutis-Miller
 - (b) Explanation about recursive formulation of determinant
 - (c) Using an approximate preconditioner
 - i. Approximate series
 - ii. Convergence proofs
5. Bounds using the stretch

1. Introduction

We consider the problem of computing the determinant of symmetric, diagonally dominant (SDD) matrices, i.e. real symmetric matrices A for which:

$$A_{ii} \geq \sum_{j \neq i} |A_{ij}|$$

The set of all such matrices of size $n \times n$ is denoted SDD_n , and the set of all symmetric real matrices is called \mathcal{S}_n . Call m the number of non-zero entries in A . We are interested in computing the determinant of sparse matrices, i.e. matrices for which $m \ll n^2$.

The best algorithm known for computing the determinant of general matrices, the Cholesky factorization, runs in a cubic complexity $O(n^3)$. Computing the factorization can be sped up for a few specific patterns such as trees, but no algorithm has been shown to work in a generic way for SDD_n , let alone general symmetric matrices. We present an algorithm that returns an approximation of the logarithm of the determinant in time quasi-linear with the number of non-zero entries of A . More specifically, we show that our algorithm, **UltraLogDet**, computes an ϵ -approximation of the logarithm of the determinant in expected time¹:

$$\tilde{O}\left(m(\log n)^7 \epsilon^{-2} \log\left(\frac{n\kappa_A}{\epsilon}\right)\right)$$

where κ_A is the condition number of A . This algorithm builds upon the work of Spielmann and Teng on *ultra-sparsifiers* Spielman and Teng (2009), and it critically exploits the recent improvements from Koutis, Miller and Peng Koutis et al. (2010). This is to our knowledge the first algorithm that presents a nearly linear complexity which depends neither on the condition number of A (except through a log-term) nor on a specific pattern for the non-zero coefficients of A .

The high complexity of the algorithm transpires through the large exponent of $\log n$. However, our algorithm will directly benefit from any improvement on ultra-sparsifiers. Given the considerable practical importance of such preconditioners, we expect some fast improvements in this area. Also, the bulk of the work is performed in a Monte Carlo procedure that is straightforward to parallelize. Furthermore, we also present simpler, non-optimal algorithms that compute upper and lower bounds of the logarithm of the determinant, and that may be of more immediate practical interest.

Background There are two approaches in numerical linear algebra to compute approximately a determinant (or the log of the determinant): by performing a (partial) Cholesky factorization of A , or by considering the trace of some power series.

As mentioned above, the Cholesky factorization performs a decomposition of the form: $A = PLDL^T P^T$ with P a permutation matrix, L a low-triangular matrix with 1 on the diagonal and D a diagonal matrix of non-negative coefficients. Then the log-determinant of A is simply²:

$$\log |A| = \sum_i \log D_{ii}$$

The complexity of dense Cholesky factorization for dense matrices is $O(n^3)$. Unfortunately, Cholesky factorization usually does not gain much from the knowledge of the sparsity pattern due to the *fill-in problem* (see Meurant (1999), section 3.2). There is one case, though, for which Cholesky factorization is efficient: if the sparsity pattern of A is a tree, then performing Cholesky factorization takes $O(n)$ time, and the matrix L is a banded matrix Liu (1990). If the sparsity pattern of A is not a tree, however, this advantageous decomposition does not hold anymore.

1. We use the notation \tilde{O} to hide a factor at most $(\log \log n)^8$

2. We will use the $|\cdot|$ operator to denote the determinant, it will be clear from the context that it is different from the absolute value.

When the matrix A is close to the identity, more precisely when the spectral radius of $M = A - I$ is less than 1, one can use the remarkable Martin expansion of the log-determinant Martin (1992):

$$\log |A| = \text{Tr}(\log A) \quad (1)$$

where $\log A$ is the matrix logarithm defined by the series expansion:

$$\log A = \sum_{i=0}^{\infty} \frac{(-1)^i}{i+1} M^i \quad (2)$$

The determinant can then be computed by a sum of traces of the power of M , and the rate of convergence of this series is driven by the spectral radius M . This line of reasoning has lead to looking for decompositions of A of the form $A = U + V$ with the determinant of U being easier to compute and $U^{-1}V + I$ having a small spectral radius. Then $\log |A| = \log |U| + \log |U^{-1}V + I|$. The most common decomposition U, V is in terms of block diagonal and off-diagonal terms, which can then use Hadamard inequalities on the determinant to bound the error Ipsen and Lee (2006). Diagonal blocks also have the advantage of having determinants easy to compute. However, this approach requires some strong assumptions on the condition number of A , which may not hold in practice.

The trace approach is driven by *spectral properties* (the condition number) while the Cholesky approach is driven by *graphical properties* (the non-zero pattern). We propose to combine these two approaches by decomposing the problem with one component that is close to a tree (and is more amenable to Cholesky methods), and one component that has a bounded condition number. Our solution is to use a *spectral sparsifier* introduced by Spielman in Spielman and Teng (2008).

The problem of estimating determinants has important applications in spatial data analysis, statistical physics and statistics. In spatial statistics, it is often convenient to interpolate measurements in a 2-, 3- or 4-dimensional volume using a sparse Gaussian process, a technique known in the geospatial community as *kriging* Zhang and Wang (2010); Li and Sudjianto (2005). Computing the optimal parameters of this Gaussian process involves repeated evaluations of the partition function, which is a log-determinant. In this context, a diagonally dominant matrix for the Gram matrix of the process corresponds to weak interactions between points of measure (which is verified in some contexts Pace and Barry (1997)). Determinants also play a crucial role in quantum physics and in theoretical physics. The wave function of a system of multiple fermion particles is an antisymmetric function which can be described as a determinant (Slatter determinant, Atkins and Friedman (2011); Löwdin (1955)). In the theory of quantum chromodynamics (QCD), the interaction between fundamental particles can be discretized as a lattice, and the energy level of particles is the determinant of some functional operators over this lattice Duncan et al. (1998). It is itself a very complex problem because of the size of the matrices involved for any non-trivial problem, for which the number of variables is typically in the millions Bernardson et al. (1994). In this respect, the restriction to diagonally dominant matrices can be interpreted as an interaction between relatively massive particles de Forcrand and Gupta (1989), or as a bound on the propagation of interactions between sites in the lattice Bernardson et al. (1994).

For these reasons, computing estimates of the log-determinant has been an active problem in physics and statistics. In particular, the Martin expansion presented in Equation 1 is extensively used in quantum physics Ipsen and Lee (2006), and it can be combined with sampling method to estimate the trace of a matrix series (Zhang et al. (2008), McCourt (2008), Zhang and Leithead (2007)). Another different line of research has worked on bounds on the values of the determinant itself. This is deeply connected to simplifying statistical models using variational methods. Such a relaxation using a message-passing technique is presented in Wainwright and Jordan (2006). Our method is close in spirit to Reuksen's work Reuksen (2002) by the use of a preconditioner. However, Reuksen considers preconditioners based on a clever approximation of the Cholesky decomposition, and its interaction with the eigenvalues of the complete matrix is not well understood. Using simpler

methods based on sampling, we are able to carefully control the spectrum of the reminder, which in turn leads to strong convergence guarantees.

Main results. We first present some general results about the preconditioning of determinants. Consider $A \succ 0$, and some other matrix $B \succ 0$ that is close to A in the spectral sense. All the results of this article stem from observing that:

$$\begin{aligned} \log |A| &= \log |B| + \log |B^{-1}A| \\ &= \log |B| + \text{Tr}(\log(B^{-1}A)) \end{aligned}$$

The first section is concerned with estimating the reminder term using the Martin expansion. The exact inverse B^{-1} is usually not available, but we are given instead a linear operator C that is an ϵ -approximation of B^{-1} . We show in this section that if the precision of this approximation is high enough, we can estimate the reminder with high probability and with a reasonable number of calls to the operator C (these terms will be made precise in the rather technical Theorem 7). Using this general framework, the subsequent section shows that spectral sparsifiers make excellent preconditioners that are close enough to A so that computing the Martin expansion is not too expensive. In particular, we build upon the recursive structure of Spielman-Teng ultra-sparsifiers to obtain our main result:

Theorem 1 *On input $A \in \text{SDD}_n$ with m non-zeros, $\epsilon > 0$, $\eta > 0$, the algorithm XXX returns a scalar z so that:*

$$\mathbb{P} [|z - n^{-1} \log |A|| > \epsilon] \leq \eta$$

and this algorithm completes in expected time $\tilde{O}\left(m(\log n)^7 \epsilon^{-2} \log\left(\frac{n\kappa_A}{\epsilon}\right)\right)$.

A note on scaling. Unlike other characteristics of linear operators, the determinant and the log-determinant are very sensitive to dimensionality. We will follow the approach of Reuksen Reusken (2002) and consider the *regularized log-determinant* $f(A) = n^{-1} \log |A|$ instead of the log-determinant. The regularized determinant has appealing properties with respect to dimensionality. In particular, its sensitivity to perturbations does not increase with the dimensionality, but only depends on spectral properties of the operator A . For example, calling λ_{\min} and λ_{\max} the minimum and maximum eigenvalues of A , respectively:

$$\begin{aligned} \log \lambda_{\min} &\leq f(A) \leq \log \lambda_{\max} \\ |f(A + \epsilon I) - f(A)| &\leq \epsilon \|A^{-1}\|_2 + O(\epsilon^2) \end{aligned}$$

The last inequality in particular shows that any perturbation to $\log |A|$ will be in the order $O(n)$, and so why all the interesting log-determinants in practice will be dominated by some $O(n)$.

The rest of the article is structured as follows. In the next section, we present some results about estimating the log-determinant from a truncated expansion. These results will justify the use of *preconditioners* to compute the determinant of a matrix. The techniques developed by Spielman et al. work on the Laplacians of weighted graphs. Section 3 introduces some new concepts to expand the notion of determinants to Laplacian matrices, and presents a few straightforward results in the relations between graph Laplacians and SDD matrices. Section 4 will use these new concepts to introduce a first family of preconditioners based on low-stretch spanning trees. Finally, Section 5 contains the proof of our main result, an algorithm to compute determinants in near-linear time.

2. Preconditioned log-determinants

We begin by considering a simple sampling algorithm to compute log-determinants, presented first in Barry and Pace (1999). We will first present some error bounds on this algorithm that expand on bounds previously presented in Bai et al. (1996) and Barry and Pace (1999).

Consider a real symmetric matrix $S \in \mathcal{S}_n^+$ such that its spectral radius is less than 1: $0 \preceq S \preceq (1 - \delta)I$ for some $\delta \in (0, 1)$. We want to compute $\log |I - S|$ up to precision ϵ and with high probability. From the Martin expansion:

$$\log |I - S| = -\text{Tr} \left(\sum_{k=1}^{\infty} \frac{1}{k} S^k \right) \quad (3)$$

This series of traces can be estimated by Monte Carlo sampling, up to precision ϵ with high probability. In order to bound the errors, we will bound the large deviation errors using Hoeffding inequality. We use a modified version of Proposition 4.2 from Bai et al. (1996):

Lemma 2 Consider $H \in \mathcal{S}_n$ lower- and upper-bounded by λ_{\min} and λ_{\max} : $\lambda_{\min}I \preceq H \preceq \lambda_{\max}I$. Consider p vectors sampled from the standard Normal distribution: $\mathbf{u}_i \sim \mathcal{N}(\mathbf{0}, I_n)$ for $i = 1 \dots p$. Then for all $\nu > 0$:

$$\mathbb{P} \left[\left| \frac{1}{p} \sum_{i=1}^p \mathbf{u}_i^T H \mathbf{u}_i - \text{Tr}(H) \right| \geq \frac{\nu}{p} \right] \leq 2 \exp \left(-\frac{2\nu^2}{p(\lambda_{\max} - \lambda_{\min})^2} \right)$$

Proof This is a simple adaptation from Bai et al. (1996). ■

We use the following Bernstein inequality:

Lemma 3 $X_1 \dots X_U$ are independant random variables with $\mathbb{E}[X_u] = 0$, $|X_u| < c$. Call $\sigma^2 = \frac{1}{U} \sum_u \text{Var}(X_u)$, then for all $\epsilon > 0$:

$$\mathbb{P} \left[\frac{1}{U} \sum_u X_u \geq \epsilon \right] \leq \exp \left(-\frac{U\epsilon^2}{2\sigma^2 + 2c\epsilon/3} \right)$$

The proof of this result can be found in a lecture note from Peter Bartlett's class (stat241B/spring 2003). We can adapt some results from Barry and Pace (1999) to prove this bound on the deviation from the trace.

Lemma 4 Consider $H \in \mathcal{S}_n$, and call λ_i the eigenvalues of H . Consider p vectors sampled from the standard Normal distribution: $\mathbf{u}_i \sim \mathcal{N}(\mathbf{0}, I_n)$ for $i = 1 \dots p$. Then for all $\epsilon > 0$:

$$\mathbb{P} \left[\left| \frac{1}{p} \sum_{i=1}^p \frac{\mathbf{u}_i^T H \mathbf{u}_i}{\mathbf{u}_i^T \mathbf{u}_i} - \frac{1}{n} \text{Tr}(H) \right| \geq \epsilon \right] \leq \exp \left(-\frac{p\epsilon^2}{2n^{-1}(\lambda_{\max} - \lambda_{\min})^2 + (\lambda_{\max} - \lambda_{\min})\epsilon} \right)$$

with $\lambda_{\max} = \max \lambda_i$ and $\lambda_{\min} = \min \lambda_i$.

Proof The distribution of \mathbf{u}_i is invariant through a rotation, so we can consider H diagonal. We assume without loss of generality that $H = \text{diag}(\lambda_1 \dots \lambda_n)$. Again without loss of generality, we assume that $\lambda'_{\max} = \lambda_{\max} - \lambda_{\min}$ and $\lambda'_{\min} = 0$ (by considering $H' = H - \lambda_{\min}I$). Call $V_i = \frac{\mathbf{u}_i^T H \mathbf{u}_i}{\mathbf{u}_i^T \mathbf{u}_i} - n^{-1} \text{Tr}(H)$. Using results from Barry and Pace (1999), we have: $|V_i| \leq \lambda_{\max} - \lambda_{\min}$, $\mathbb{E}[V_i] = 0$ and:

$$\text{Var}(V_i) = \frac{2}{n^2} \sum_{i=1}^n (\lambda_i - n^{-1} \text{Tr}(H))^2$$

Each of the variables V_i is independent, so invoking Lemma 3 gives:

$$\mathbb{P} \left[\frac{1}{p} \sum_{i=1}^p V_i \geq \epsilon \right] \leq \exp \left(-\frac{p\epsilon^2}{2\sigma^2 + 2(\lambda_{\max} - \lambda_{\min})\epsilon/3} \right)$$

with $\sigma^2 = \frac{2}{n^2} \sum_{i=1}^n (\lambda_i - n^{-1} \text{Tr}(H))^2$ ■

The previous lemma shows that if the eigenspectrum of a matrix is bounded, we can obtain a Hoeffding bound on the error done by sampling the trace. Furthermore, the convergence of the series 3 is also determined by the extremal eigenvalues of S . If we truncate the series (3), we can bound the truncation error using the extremal eigenvalues. We formalize this intuition in the following theorem, which is adapted from the main Theorem in Barry and Pace (1999). While that main Theorem in Barry and Pace (1999) only considered a confidence interval based on the covariance properties of Gaussian distribution, we generalize this result to more general Hoeffding bounds.

Theorem 5 Consider $S \in \mathcal{S}_n^+$ with $0 \preceq S \preceq (1 - \delta)I$ for some $\delta \in (0, 1)$. Call y the quantity to estimate:

$$y = \log |I - S|$$

and consider $\mathbf{u}_i \sim \mathcal{N}(\mathbf{0}, I_n)$ for $i = 1 \dots p$. Call $\hat{y}_{p,l}$ an estimator of the truncated series of l elements computed by sampling the trace using p samples:

$$\hat{y}_{p,l} = \frac{n}{p} \sum_{j=1}^p \sum_{k=1}^l \frac{1}{k} \frac{\mathbf{u}_j^T S^k \mathbf{u}_j}{\mathbf{u}_j^T \mathbf{u}_j}$$

By choosing $p \geq 4 \log(\eta^{-1}) \epsilon^{-2} (2n^{-1} \log^2(\delta^{-1}) + \epsilon \log(\delta^{-1}))$ and $l \geq 2\delta^{-1} \log(\frac{n}{\delta\epsilon})$, the error probability is small:

$$\mathbb{P}[|y - \hat{y}_{p,l}| \geq n\epsilon] \leq \eta$$

Proof The proof of this theorem follows the proof of the Main Theorem in Barry and Pace (1999) with some slight modifications. Using triangular inequality:

$$|y - \hat{y}_{p,l}| \leq |\mathbb{E}[\hat{y}_{p,l}] - \hat{y}_{p,l}| + |y - \mathbb{E}[\hat{y}_{p,l}]|$$

Since S is upper-bounded, we have for all $k \in \mathbb{N}$:

$$|\text{Tr}(S^k)| \leq n(1 - \delta)^k$$

Using again triangle inequality, we can bound the error with respect to the expected value:

$$\begin{aligned} |y - \mathbb{E}[\hat{y}_{p,l}]| &= \left| \sum_{i=l+1}^{\infty} \frac{(-1)^i}{i} \text{Tr}(S^k) \right| \\ &\leq \sum_{i=l+1}^{\infty} \frac{1}{i} |\text{Tr}(S^k)| \\ &\leq \frac{1}{l+1} \sum_{i=l+1}^{\infty} |\text{Tr}(S^k)| \\ &\leq \frac{n}{l+1} \sum_{i=l+1}^{\infty} (1 - \delta)^k \\ &\leq \frac{n}{l+1} \frac{(1 - \delta)^{l+1}}{\delta} \\ &\leq n \frac{(1 - \delta)^{l+1}}{\delta} \end{aligned}$$

Hence, for a choice of $l \geq \delta^{-1} \log(\frac{2n}{\epsilon\delta})$, the latter part is less than $\epsilon/2$. We now bound the first part using Lemma 4. Call H the truncated series:

$$H = - \sum_{i=1}^m \frac{1}{i} S^i$$

This truncated series is upper-bounded by 0. The lowest eigenvalue of the truncated series can be lower-bounded in terms of δ :

$$H = -\sum_{i=1}^m \frac{1}{i} S^i \succeq -\sum_{i=1}^m \frac{1}{i} (1-\delta)^i I \succeq -\sum_{i=1}^{+\infty} \frac{1}{i} (1-\delta)^i I \succeq -\sum_{i=1}^{+\infty} \frac{1}{i} (1-\delta)^i I = (\log \delta) I$$

We can now invoke Lemma 4 to conclude:

$$\mathbb{P} \left[\left| \frac{1}{p} \sum_{i=1}^p (\mathbf{u}_i^T \mathbf{u}_i)^{-1} \mathbf{u}_i^T H \mathbf{u}_i - n^{-1} \text{Tr}(H) \right| \geq \frac{\epsilon}{2} \right] \leq \exp \left(-\frac{p\epsilon^2}{8n^{-1}(\log \delta)^2 + 4 \log(1/\delta)\epsilon} \right)$$

and $\exp \left(-\frac{p\epsilon^2}{8n^{-1}(\log \delta)^2 + 4 \log(1/\delta)\epsilon} \right) \leq \eta$. This latter inequality is verified for any

$$p \geq 4 \log(\eta^{-1}) \epsilon^{-2} (2n^{-1} \log^2(\delta^{-1}) + \epsilon \log(\delta^{-1}))$$

■

We immediately get the following result that justifies the notion of preconditioners for determinants. The corresponding algorithm, which we call **PreconditionedLogDetMonteCarlo**, is presented in Algorithm ??.

Corollary 6 *If A and B are positive definite matrices so that B is a κ -approximation of A :*

$$A \preceq B \preceq \kappa A \tag{4}$$

*then the algorithm **PreconditionedLogDetMonteCarlo** computes an ϵ -approximation of $\frac{1}{n} \log |B^{-1}A|$ with high probability, by performing $O\left(\kappa \epsilon^{-2} (\log \kappa)^2 \log\left(\frac{n^2 \kappa}{\epsilon}\right)\right)$ vector inversions from B and vector multiplies from A .*

Usually, computing the exact inverse by an SDD matrix is too expensive. We can instead extend the previous result to consider a black box that approximatively inverts compute $B^{-1}x$. If the error introduced by the approximate inversion is limited, the result from the previous corollary still holds. This is what the following theorem establishes:

Theorem 7 *Consider $A, B \in S_n^{+*}$ with B a κ -approximation of A . Furthermore, assume there exists a function C so that for all $y \in \mathbb{R}^n$, C returns an ϵ -approximation of $B^{-1}y$:*

$$\|C(y) - B^{-1}y\|_B \leq \|B^{-1}y\|_B \epsilon$$

*Given $\eta > 0$ and $\mu > 0$, if $\epsilon \leq \frac{\mu^3 \eta}{16\kappa^3 \kappa(B)(1+n\epsilon) \log(\eta^{-1}) \log^2 \kappa}$, then the algorithm **PreconditionedLogDetMonteCarlo** returns a scalar z so that:*

$$\mathbb{P} [|z - n^{-1} \log |B^{-1}A|| \geq \mu] \leq \eta$$

by performing $O\left(\kappa \mu^{-2} \log\left(\frac{n\kappa}{\mu}\right) \log^2(\kappa) \log(\eta^{-1})\right)$ vector calls to the operator C and vector multiplies from A .

The proof of this result is detailed in Appendix A.

This last theorem shows that we can compute a good approximation of the log-determinant if the preconditioner B : (a) is close to A in the spectral sense, and (b) can be approximately inverted and the error introduced by the approximate inversion can be controlled. This happens to be the case for symmetric, diagonally dominant matrices.

3. Ultra-sparsifiers as determinant preconditioners

3.1 Making the problem laplacian

From now on, we consider all matrices to be in SDD_n . The techniques we will develop work on Laplacian matrices instead of SDD matrices. An SDD matrix is positive semi-definite while a Laplacian matrix is always singular, since its nullspace is spanned by $\mathbf{1}$. We generalize the definition of the determinant to handle this technicality.

Definition 8 Pseudo-log-determinant (PLD): Let $A \in \mathcal{S}^{n+}$ be a non-null positive semi-definite matrix. The pseudo-log-determinant is defined by the sum of the logarithms of all the positive eigenvalues:

$$ld(A) = \sum_{\lambda_i > 0} \log(\lambda_i)$$

where λ_i are the eigenvalues of A .

The interest of the PLD lies in the connection between SDD matrices and their associated Laplacian. Recall that an SDD matrix can be transformed into a Laplacian (which has all off-diagonal terms non-positive, and for which the sum of each row and each column is zero). As underlined in Gremban (1996), solving an equation involving an SDD can be turned into solving an equation with an SDD in which all the off-diagonal terms are non-positive of twice the same size. This transform is also useful for computing determinants. The following lemma shows that computing the log-determinant of an SDD with positive off-diagonal entries can be turned into computing the log-determinants of two SDD matrices with non-positive off-diagonal entries.

Lemma 9 Consider the decomposition of a SDD matrix into a positive and a negative sum:

$$A = P - N$$

with $P_{ij} \leq 0$ for $i \neq j$ and $N_{ij} \leq 0$ for all i, j so that $|A_{ij}| = |P_{ij}| + |N_{ij}|$. Then the matrices $\begin{pmatrix} P & N \\ N & P \end{pmatrix} \in SDD_{2n}$ and $P + N \in SDD_n$ with all the off-diagonal coefficients being negative, and we have the relation:

$$\log \left| \begin{pmatrix} P & N \\ N & P \end{pmatrix} \right| = \log |P + N| + \log |P - N|$$

Proof Call $C = \begin{pmatrix} P & N \\ N & P \end{pmatrix}$

Consider μ an eigenvalue of A , and x an associated eigenvector: $Ax = \mu x$. Then:

$$\begin{pmatrix} P & N \\ N & P \end{pmatrix} \begin{pmatrix} x \\ -x \end{pmatrix} = \begin{pmatrix} Ax \\ -Ax \end{pmatrix} = \mu \begin{pmatrix} x \\ -x \end{pmatrix}$$

which shows that the vector $\begin{pmatrix} x \\ -x \end{pmatrix}$ is an eigenvector of C with associated eigenvalue μ . Similarly, if η is an eigenvector of B with associated eigenvalue y , then η is also an eigenvector of C with associated eigenvalue $\begin{pmatrix} y \\ y \end{pmatrix}$. Since C is exactly of size $2n$, the set of eigenvalues of C is exactly the union of sets of eigenvalues of A and B . Thus $\log |C| = \sum_i \log \mu_i + \sum_i \log \eta_i = \log |A| + \log |B|$ ■

Furthermore, computing the determinant of SDD matrices with non-negative off-diagonal entries is equivalent to computing the PLD of Laplacian. Recall that a Laplacian matrix can be constructed from any SDD with non-positive off-diagonal entries matrix by adding a “ground” node (see Gremban’s thesis Gremban (1996), Chapter 4).

Definition 10 Augmented Laplacian of an SDD matrix (Definition XXX in Gremban (1996)): Let $A \in SDD_n$. Let L, D be the decomposition of A into L a Laplacian matrix and D a (non-negative) diagonal matrix such that $A = L + D$. The augmented matrix of A is:

$$L_A = \begin{pmatrix} A & -\mathbf{d} \\ -\mathbf{d}^T & h \end{pmatrix}$$

with $\mathbf{d} \in \mathbb{R}^n$: $\mathbf{d}_i = D_{ii}$ and $h = \sum_i D_{ii}$.

From now on, given any $A \in SDD_n$, we will implicitly define a Laplacian L_A associated to it, using the augmented Laplacian. This process is also reversible: to any Laplacian L we can associate a unique positive definite matrix F_L (up to a permutation), and this transform preserves eigenvalues and matrix inequalities. We call this process “floating” of the Laplacian, by opposition to “grounding” the SDD matrix as a Laplacian.

Definition 11 Floating a Laplacian. Consider L a Laplacian matrix. Call F_L the matrix formed by removing the last row and the last column from L .

The following lemma shows that the Laplacian matrix overdetermines a system, and that no information is lost by floating it.

Lemma 12 Consider Z a (weighted) Laplacian matrix, then:

- $L_{F_Z} = Z$
- The eigenvalues of F_Z are the positive eigenvalues of Z , and the corresponding eigenvectors for F_Z are the same eigenvectors, truncated by the last coefficient.
- $\text{ld}(Z) = \log |F_Z|$
- Given Z_1, Z_2 Laplacian matrices, we have $Z_1 \succeq Z_2 \Rightarrow F_{Z_1} \succeq F_{Z_2}$

The proof of this lemma is in Appendix XXX.

Using the floating and the grounding procedures, we can transform the PLD of a Laplacian matrix into the log-determinant of a SDD matrix. The following proposition establishes this equivalence.

Proposition 13 Let $A \in SDD_n$, and L_A its augmented Laplacian. Let $P = (I_n \mathbf{0})$ be a projection matrix over the first n columns of L_A . Then:

- If λ is an eigenvalue of A with associated eigenvector x , it is also an eigenvalue of PL_AP^T with the same eigenvector
- $\text{ld}(L_A) = \log |A|$
- Let $A, B \in SDD_n$. $\text{ld}(L_AL_B) = \text{ld}(L_A) + \text{ld}(L_B)$ and $\text{ld}(L_B^+L_A) = \log |B^{-1}A|$ with L_A^+ the pseudo inverse of L_A .

Proof XXX The proofs of this proposition are technical and are contained in Appendix B. ■

A Laplacian matrix can be considered either for its graphical properties, or for its algebraic properties. Recent results has shown a deep connection between these two aspects, and let us develop a general framework for computing determinants: given $G \in SDD_n$, we consider the graph associated to its Laplacian L_G , through the grounding procedure. Using graphical properties of L_G , we can construct a weighted subgraph $L_H \subseteq L_G$ for which the PLD is easier to compute and that is

a good approximation of L_G in the spectral sense. We can float the subgraph L_H and apply results of Section 2 to approximate the remainder with high probability. More precisely:

$$\begin{aligned}\log |G| &= \text{ld}(L_H) - \log |F_{L_H}^{-1}| + \log |G| \\ &= \text{ld}(L_H) + \log |F_{L_H}^{-1}G|\end{aligned}$$

The first term $\text{ld}(L_H)$ is usually easier to compute by considering the graphical properties of L_H , while the remainder $\log |F_{L_H}^{-1}G|$ is approximated by sampling. Preconditioner graphs L_H are typically efficient to factorize using Cholesky factorization, and close enough to G that the sampling procedure from the previous section can be applied to compute $\log |F_{L_H}^{-1}G|$. We will see how to adapt Spielman and Teng's remarkable work on *ultra-sparsifiers* to produce good preconditioners H for the determinant.

3.2 A first preconditioner

While the results in this section are not the main claims of this paper, we hope they will provide some intuition, and an easier path towards an implementation.

We present a first preconditioner that is not optimal, but that will motivate our results for stronger preconditioners: a tree that spans the graph G . Every graph has a low-stretch spanning tree, as discovered by Alon et al. (1995). The bound of Alon et al. was then improved by Abraham et al. (2008). We restate their main result.

Lemma 14 *Consider L_G a weighted graph. There exists a tree L_T that is a subgraph of L_G so that:*

$$L_G \preceq L_T \preceq \kappa L_G$$

with $\kappa = Cm \log n \log \log n (\log \log \log n)^3$ for some constant $C > 0$.

Proof This follows directly from Abraham et al. (2008). L_T is a subgraph of L_G (with the same weights on the edges), so $L_T \preceq L_G$. Furthermore, we have $L_T \preceq \text{st}_G(T) L_G$ following a result of Boman and Hendrickson (2004) cited by Spielman et al. in Spielman (2010). ■

Trees enjoy a lot of convenient properties for Gaussian elimination. The Cholesky factorization of a tree can be computed in linear time, and furthermore this factorization has a linear number of non-zero elements Spielman and Teng (2009). This factorization can be expressed as:

$$L_T = PLDL^T P^T$$

where P is a permutation matrix, L is a lower-triangular matrix with the diagonal being all ones, and D a diagonal matrix in which all the elements but the last one are positive, the last element being 0. These well-known facts about trees are presented in Spielman and Teng (2009). Once the Cholesky factorization of the tree is performed, the log-determinant of the original graph is an immediate by-product:

$$\log |L_T| = \sum_{i=1}^{n-1} \log D_{ii}$$

Furthermore, computing $L_T^+ x$ also takes $O(n)$ computations by forward-backward substitution. Applying Corollary 6 gives immediately:

Theorem 15 *Let G be a graph with n vertices and m edges. Its PLD can be computed up to a precision of ϵ and with high probability in time $\tilde{O}\left(mn (\log m)^3 (\log n)^2 \frac{\log \epsilon^{-1}}{\epsilon^2}\right)$.*

Proof We assume G is connected, hence $m \geq n$ and $\log \kappa = \tilde{O}(\log m)$. ■

This bound may be of independent interest since it requires relatively little machinery to compute, and it is an improvement already for graphs with small vertex degree ($m = \mathcal{O}(n^{1+o(1)})$) over the Cholesky factorization of G . Also, note that the PLD of the tree constructed above provides an upper bound to the log-determinant of G since $L_G \preceq L_T$. We will see in Subsection XXX that we can also compute a non-trivial lower bound as well.

3.3 Incremental sparsifiers

We can do better and achieve near-linear time by using ultra-sparsifiers. The main insight of our result is that the class preconditioners presented by Spielman and Teng are based on incomplete Cholesky factorization, and hence have a determinant that is relatively easy to compute, and that furthermore they are excellent spectral preconditioners, so the procedure `PreconditionedLogDetMonteCarlo` is efficient to apply. We reintroduce some concepts presented in Koutis et al. (2010) to present a self-contained result.

The central idea is to sample $\mathcal{O}(n)$ edges from the graph A , to form a subgraph B that is close to a tree (hence it is easy to compute some partial Cholesky factorization), yet it is close to the original A in the spectral sense ($A \preceq B \preceq \kappa A$), thanks to the additional edges. The partial Cholesky factorization is computed using the `GreedyElimination` algorithm presented in Koutis et al. (2010). In order for this section to be self-contained, we repeat the main results of Section 4 in Spielman and Teng (2009).

Consider a Laplacian matrix G . There exists an algorithm that computes the partial Cholesky factorization:

$$B = PLCL^T P^T$$

where:

- P is a permutation matrix
- L is a non-singular, low triangular matrix of the form

$$L = \begin{pmatrix} L_{1,1} & 0 \\ L_{2,1} & I_{n_1} \end{pmatrix}$$

with the diagonal of $L_{1,1}$ being all ones.

- C has the form

$$C = \begin{pmatrix} D_{n-n_1} & 0 \\ 0 & A_1 \end{pmatrix}$$

and every row and column of A_1 has at least 3 non-zero coefficients. Furthermore, A_1 is itself Laplacian and:

$$\text{ld}(B) = \sum_1^{n-n_1} \log D_{ii} + \text{ld}(A_1)$$

The exact algorithm is detailed that achieves this factorization is presented in Koutis et al. (2010). Using this factorization, the PLD of the original Laplacian A is:

$$\begin{aligned} \text{ld}(A) &= \text{ld}(B) + \text{ld}(B^+ A) \\ &= \sum_1^{n-n_1} \log D_{ii} + \text{ld}(A_1) + \text{ld}(B^+ A) \end{aligned}$$

Thus, we are left with solving a smaller problem A_1 , and approximate the value of $\text{ld}(B^+A)$ using the algorithm **SampleLogDet**. ST preconditioners are appealing for this task: they guarantee that A_1 is substantially smaller than A , so the recursion completes in $O(\log \log n)$ steps. Furthermore, computing the vector product B^+Ax is itself efficient (it can be done in near-linear time), so the sampling algorithm can be run in reasonable time. We formalize the notion of chain of preconditioners by reintroducing some material from Koutis et al. (2010).

Definition 16 *Definition 7.1 from Koutis et al. (2010). $\kappa(n)$ -good chain. $\mathcal{C} = \{A_1 = A, B_1, A_2, \dots, A_d\}$ a chain of graphs with n_i and m_i the number of vertices and edges of A_i . \mathcal{C} is $\kappa(n)$ -good for A if:*

1. $A_i \preceq B_i \preceq \kappa(n_i) A_i$
2. $A_{i+1} = \text{GreedyElimination}(B_i)$
3. $m_i/m_{i+1} \geq c_r \sqrt{\kappa(n_i)}$ for some constant c_r .

Good chains exist, as found by Koutis, Miller and Peng:

Lemma 17 *(Lemma 7.3 from Koutis et al. (2010)) Given a graph A , **BuildChain**(A, p) from Koutis et al. (2010) produces a $\tilde{O}(\log^4 n)$ -good chain for A with probability $1 - p$. The algorithm runs in time*

$$\tilde{O}((m \log n + n \log^2 n) \log(1/p))$$

These chains furthermore can be used as good preconditioners for conjugate gradient and lead to near-linear algorithms for approximate inversion (Lemma 7.2 from Koutis et al. (2010))

Lemma 18 *(Lemma 7.2 from Koutis et al. (2010)) Given a $\kappa(n)$ -good chain for A , a vector x such that $\|x - L_A^+ b\|_A < \epsilon \|L_A^+ b\|_A$ can be computed in $O(m_d^3 m_1 \sqrt{\kappa(n_1)} \log(1/\epsilon))$*

It should now become clear how we can combine a $\kappa(n)$ -good chain with the Algorithm **PreconditionedLogDetMonteCarlo**. We start by building a chain. The partial Cholesky factorizations provide an upper bound to $\text{ld}(A)$. We then refine this upper bound by running **PreconditionedLogDetMonteCarlo** at each state of the chain to approximate $\text{ld}(B_i^+ A_i)$ with high probability. The complete algorithm is presented in Algorithm ???. Using

3.4 Stretch bounds on preconditioners

How good is the estimate provided by the preconditioner? Intuitively, this depends on how well the preconditioner L_H approximates the graph L_G . This notion of quality of approximation can be formalized by the notion of *stretch*. This section presents a deterministic bound on the PLD of L_G based on the PLD of L_H and the stretch of G relative to H . This may be useful in practice as it gives a (tight) interval for the PLD before performing any Monte-Carlo estimation of the residual.

The stretch of a graph is usually defined with respect to a (spanning) tree. In our analysis, it is convenient and straightforward to generalize this definition to arbitrary graphs. To our knowledge, this straightforward extension is not considered in the literature, so we feel compelled to properly introduce it.

Definition 19 *Generalized stretch. Consider \mathcal{V} a set of vertices, $G = (\mathcal{V}, \mathcal{E}_G)$, $H = (\mathcal{V}, \mathcal{E}_H)$ connected graphs over the same set of vertices, and L_G, L_H their respective Laplacians. The stretch of G with respect to H is the sum of the effective resistances of each edge of graph G with respect to graph H ,*

$$\text{st}_H(G) = \sum_{(u,v) \in \mathcal{E}_G} L_G(u,v) (\mathcal{X}_u - \mathcal{X}_v)^T L_H^+ (\mathcal{X}_u - \mathcal{X}_v)$$

with $\mathcal{X}_u \in \mathbb{R}^n$ the unit vector that is 1 at position u , and zero otherwise.

If the graph H is a tree, this is a standard definition of stretch, because the effective resistance $(\mathcal{X}_u - \mathcal{X}_v)^T L_H^+ (\mathcal{X}_u - \mathcal{X}_v)$ between vertices u and v is the sum of all resistances over the unique path between u and v (see Lemma 2.4 in Spielman and Woo (2009)). Furthermore, the arguments to prove Theorem 2.1 in Spielman and Woo (2009) carry over to our definition of stretch. For the sake of completeness, we include this result:

Lemma 20 *(Straightforward generalization of Theorem 2.1 in Spielman and Woo (2009)) Let $G = (\mathcal{V}, \mathcal{E}_G)$, $H = (\mathcal{V}, \mathcal{E}_H)$ be connected graphs over the same set of vertices, and L_G , L_H their respective Laplacians. Then:*

$$st_H(G) = \text{Tr}(L_H^+ L_G)$$

with L_H^+ the pseudo-inverse of L_H .

Proof

We denote $E(u, v)$ the unit matrix that is 1 in position u, v . This is the same arguments as the original proof:

$$\begin{aligned} \text{Tr}(L_H^+ L_G) &= \sum_{(u,v) \in \mathcal{E}_G} L_G(u, v) \text{Tr}(E(u, v) L_H^+) \\ &= \sum_{(u,v) \in \mathcal{E}_G} L_G(u, v) \text{Tr}\left((\mathcal{X}_u - \mathcal{X}_v)(\mathcal{X}_u - \mathcal{X}_v)^T L_H^+\right) \\ &= \sum_{(u,v) \in \mathcal{E}_G} L_G(u, v) (\mathcal{X}_u - \mathcal{X}_v)^T L_H^+ (\mathcal{X}_u - \mathcal{X}_v) \\ &= st_H(G) \end{aligned}$$

■

A number of properties of the stretch extend to general graphs using this definition. In particular, the stretch inequality (Lemma 8.2 in Spielman and Teng (2009)) can be generalized to arbitrary graphs (instead of spanning trees).

Lemma 21 *Let $G = (\mathcal{V}, \mathcal{E}_G)$, $H = (\mathcal{V}, \mathcal{E}_H)$ be connected graphs over the same set of vertices, and L_G , L_H their respective Laplacians. Then:*

$$L_G \preceq st_H(G) L_H$$

Proof The proof is very similar to that of Lemma 8.2 in Spielman and Woo (2009), except that the invocation of Lemma 8.1 is replaced by invoking Lemma 28 in Appendix B. The Laplacian G can be written as a linear combination of edge Laplacian matrices:

$$L_G = \sum_{e \in \mathcal{E}_G} \omega_e L(e) = \sum_{(u,v) \in \mathcal{E}_G} \omega_{(u,v)} (\mathcal{X}_u - \mathcal{X}_v)(\mathcal{X}_u - \mathcal{X}_v)^T$$

and using Lemma 28 gives:

$$(\mathcal{X}_u - \mathcal{X}_v)(\mathcal{X}_u - \mathcal{X}_v)^T \preceq (\mathcal{X}_u - \mathcal{X}_v)^T L_H^+ (\mathcal{X}_u - \mathcal{X}_v) L_H$$

By summing all the edge inequalities, we get:

$$\begin{aligned} L_G &\preceq \sum_{(u,v) \in \mathcal{E}_G} \omega_{(u,v)} (\mathcal{X}_u - \mathcal{X}_v)^T L_H^+ (\mathcal{X}_u - \mathcal{X}_v) L_H \\ &\preceq st_H(G) L_H \end{aligned}$$

■

This bound is remarkable as it relates any pair of (connected) graphs, as opposed to spanning trees or subgraphs. An approximation of the generalized stretch can be quickly computed using a construct detailed in Spielman and Srivastava (2011). We now introduce the main result of this section: a bound on the PLD of L_G using the PLD of L_H and the stretch.

Theorem 22 *Let $G = (\mathcal{V}, \mathcal{E}_G)$, $H = (\mathcal{V}, \mathcal{E}_H)$ be connected graphs over the same set of vertices, and L_G, L_H their respective Laplacians. Assuming $L_H \preceq L_G$, then:*

$$\text{ld}(L_H) + (n-1) \log \left(\frac{\text{st}_H(G)}{n-1} \right) \geq \text{ld}(L_G) \geq \text{ld}(L_H) + \log(\text{st}_H(G) - n + 2) \quad (5)$$

This inequality is an equality for $G = H$.

Proof This is an application of Jensen's inequality on $\text{ld}(L_H^+ L_G)$. We have $\text{ld} L_G = \text{ld} L_H + \text{ld}(L_H^+ L_G)$ and $\text{ld}(L_H^+ L_G) = \text{ld}(\sqrt{L_H}^+ L_G \sqrt{L_H}^+)$ with \sqrt{T} the matrix square root of T . From Lemma 29, we have the following inequality:

$$\begin{aligned} \text{ld}(\sqrt{L_H}^+ L_G \sqrt{L_H}^+) &\leq (n-1) \log \left(\frac{\text{Tr}(\sqrt{L_H}^+ L_G \sqrt{L_H}^+)}{n-1} \right) \\ &= (n-1) \log \left(\frac{\text{Tr}(L_H^+ L_G)}{n-1} \right) \\ &= (n-1) \log \left(\frac{\text{st}_H(G)}{n-1} \right) \end{aligned}$$

The latter equality is an application of Lemma 20.

The lower bound is slightly more involved. Call λ_i the positive eigenvalues of $\sqrt{L_H}^+ L_G \sqrt{L_H}^+$ and $\sigma = \text{st}_H(G)$. We have $1 \leq \lambda_i$ from the assumption $L_H \preceq L_G$. By definition: $\text{ld}(L_H^+ L_G) = \sum_i \log \lambda_i$. Furthermore, we know from Lemma 20 that $\sum_i \lambda_i = \sigma$. The upper and lower bounds on λ_i give:

$$\text{ld}(L_H^+ L_G) \geq \min_{\substack{1 \leq \lambda_i \\ \sum_i \lambda_i = \sigma}} \sum_i \log \lambda_i$$

Since there are precisely $n-1$ positive eigenvalues λ_i , one can show that the minimization problem above has a unique minimum which is $\log(\sigma - n + 2)$.

To see that, consider the equivalent problem of minimizing $\sum_i \log(1 + u_i)$ under the constraints $\sum_i u_i = \sigma - (n-1)$ and $u_i \geq 0$. Note that:

$$\sum_i \log(1 + u_i) = \log \left(\prod_i [1 + u_i] \right) = \log \left(1 + \sum_i u_i + \text{Poly}(u) \right)$$

with $\text{Poly}(u) \geq 0$ for all $u_i \geq 0$, so we get: $\sum_i \log(1 + u_i) \geq \log(1 + \sum_i u_i)$ and this inequality is tight for $u_1 = \sigma - (n-1)$ and $u_{i \geq 2} = 0$. Thus the vector $\lambda^* = (\sigma - n + 2, 1 \cdots 1)^T$ is (a) a solution to the minimization problem above, and (b) the objective value of any feasible vector λ is higher or equal. Thus, this is the solution (unique up to a permutation). Hence we have $\text{ld}(L_H^+ L_G) \geq \sum_i \log \lambda_i^* = \log(\sigma - n + 2)$. ■

Note that this inequality is a stronger result than the bound $L_H \preceq L_G \Rightarrow \text{st}_H(G) \geq n - 1$ presented in Spielman and Woo (2009) (a direct consequence of Lemma 20). Interestingly, it does not make assumption on the topology of the graphs (such as L_H being a subset of L_G). Research on conditioners has focused so far on low-stretch approximations that are subgraphs of the original graph. It remains to be seen if some better preconditioners can be found with stretches in $\mathcal{O}(n)$ by considering more general graphs. In this case, the machinery developped in Section 3 would not be necessary.

From a practical perspective, the stretch can be calculated also linearly with respect to the number of non-zero entries.

Lemma 23 *Let $G = (\mathcal{V}, \mathcal{E}_G)$, $H = (\mathcal{V}, \mathcal{E}_H)$ be connected graphs over the same set of vertices, and L_G , L_H their respective Laplacians. Call $r = \max_e L_H(e) / \min_e L_H(e)$. Given $\epsilon > 0$, there exists an algorithm that returns a scalar y so that:*

$$(1 - \epsilon) \text{st}_H(G) \leq y \leq (1 + \epsilon) \text{st}_H(G)$$

with high probability and in expected time $\tilde{O}(m\epsilon^{-2} \log(rn))$.

Proof This is a straightforward consequence of Theorem 2 in Spielman and Srivastava (2011). Once the effective resistance of an edge can be approximated in time $O(\log n / \epsilon^2)$, we can sum it and weight it by the conductance in G for each edge. ■

Comments

A first approximation upper-bound of the log-determinant follows immediately from the computation of the κ -good chain. We presented in Section 3.2 a first analysis to bound the value of the residue PLD. This analysis was done using trees as preconditioners, it can be carried on on more general preconditioners by introducing a generalization of the stretch over a subgraph (see the Appendix, section ??). Since the bulk of the computations are performed in estimating the residue PLD, it would be interesting to see if this could be bypassed using better bounds based on the stretch.

When looking at each step of the analysis, one can see that the ϵ^{-2} factor comes from the approximation of the trace of a matrix by sampling. This seems to be a fundamental limitation of this method and it is shared with other algorithms that rely on random projections. This result is absolute. It would be interesting to see if the analysis could be tightened to present a relative bound that does not depend on the condition number of A . Also, even if this algorithm presents a linear bound, it requires a fairly advanced machinery (ST solvers) that may limit its practicality. Some heuristic implementation, for example based on algebraic multi-grid methods, could be a first step in this direction.

The authors are much indebted to Satish Rao and Jim Demmel for suggesting the original idea and their helpful comments on the draft of this article.

Appendix A: Proofs of Section 2

3.5 Proof of Corollary6

Proof We introduce some notations that will prove useful for the rest of the article:

$$\begin{aligned} H &= I - B^{-1}A \\ S &= I - B^{-1/2}AB^{-1/2} \end{aligned}$$

with $B^{-1/2}$ the inverse of the square root of the positive-definite matrix B^3 . The inequality 4 is equivalent to $\kappa^{-1}B \preceq A \preceq B$, or also:

$$(1 - \kappa^{-1})I \succeq I - B^{-1/2}AB^{-1/2} \succeq 0$$

$$(1 - \kappa^{-1})I \succeq S \succeq 0$$

The matrix S is a contraction, and its spectral radius is determined by κ . Furthermore, computing the determinant of $B^{-1}A$ is equivalent to computing the determinant of $I - S$:

$$\begin{aligned} \log |I - S| &= \log |B^{-1/2}AB^{-1/2}| \\ &= \log |A| - \log |B| \\ &= \log |B^{-1}A| \\ &= \log |I - H| \end{aligned}$$

and invoking Theorem 5 gives us bounds on the number of calls to matrix-vector multiplies with respect to S . It would seem at this point that computing the inverse square root of B is required, undermining our effort. However, we can reorganize the terms in the series expansion to yield only full inverses of B . Indeed, given $l \in \mathbb{N}^*$, consider the truncated series:

$$\begin{aligned} y_l &= -\text{Tr} \left(\sum_{i=1}^l \frac{1}{i} S^i \right) \\ &= -\sum_{i=1}^l \frac{1}{i} \text{Tr} (S^i) \\ &= -\sum_{i=1}^l \frac{1}{i} \text{Tr} \left(\sum_j \binom{j}{i-j} (B^{-1/2}AB^{-1/2})^j \right) \\ &= -\sum_{i=1}^l \frac{1}{i} \sum_j \binom{j}{i-j} \text{Tr} \left((B^{-1/2}AB^{-1/2})^j \right) \\ &= -\sum_{i=1}^l \frac{1}{i} \sum_j \binom{j}{i-j} \text{Tr} \left((B^{-1}A)^j \right) \\ &= -\sum_{i=1}^l \frac{1}{i} \text{Tr} \left(\sum_j \binom{j}{i-j} (B^{-1}A)^j \right) \\ &= -\sum_{i=1}^l \frac{1}{i} \text{Tr} (H^i) \end{aligned}$$

Hence, the practical computation of the latter sum can be done on $A^{-1}B$. To conclude, if we compute $p = 4\epsilon^{-2} (\log \kappa)^2 \log n$ truncated chains of length $l = \kappa (2 \log n + \log \kappa + \log (\epsilon^{-1}))$, we get our result. This requires lp multiplies by A and inversions by B . \blacksquare

3. Given a real PSD matrix X , which can be diagonalized: $X = Q\Delta Q^T$ with Δ diagonal, and $\Delta_{ii} \geq 0$. Call $Y = Q\sqrt{\Delta}Q^T$ the square root of X , then $Y^2 = X$.

3.6 Proof of Theorem 7

We prove here the main result of Section 2. In the following, A and B are positive-definite matrices in \mathcal{S}_n , and B is a κ -approximation of A ($A \preceq B \preceq \kappa A$). The following notations will prove useful:

$$S = I - B^{-1/2} A B^{-1/2} \quad (6)$$

$$R = I - B^{-1} A \quad (7)$$

$$\varphi = \kappa^{-1}$$

Lemma 24 *S and R are contractions for the Euclidian and B -norms:*

$$\begin{aligned} \|S\| &\leq 1 - \varphi \\ \|R\| &\leq 1 - \varphi \\ \|R\|_B &\leq (1 - \varphi)^2 \end{aligned}$$

Proof Recall the definition of the matrix norm: $\|S\| = \max_{x^T x \leq 1} \sqrt{x^T S x}$. Since we know that $S \preceq (1 - \varphi) I$, we get the first inequality.

The second inequality is a consequence of Proposition 3.3 from Spielman and Teng (2009): A and B have the same nullspace and we have the linear matrix inequality $A \preceq B \preceq \kappa A$, which implies that the eigenvalues of $B^{-1} A$ lie between $\kappa^{-1} = \varphi$ and 1. This implies that the eigenvalues of $I - B^{-1} A$ are between $1 - \varphi$ and 0.

Recall the definition of the matrix norm induced by the B -norm over \mathbb{R}^n :

$$\begin{aligned} \|R\|_B &= \max_{x \neq 0} \frac{\|Rx\|_B}{\|x\|_B} \\ &= \max_{\|x\|_B^2 \leq 1} \sqrt{x^T R^T B R x} \\ &= \max_{x^T B x \leq 1} \sqrt{x^T R^T B R x} \\ &= \max_{y^T y \leq 1} \sqrt{y^T B^{-1/2} R^T B R B^{-1/2} y} \end{aligned}$$

and the latter expression simplifies:

$$\begin{aligned} B^{-1/2} R^T B R B^{-1/2} &= B^{-1/2} (I - A B^{-1}) B (I - B^{-1} A) B^{-1/2} \\ &= (I - B^{-1/2} A B^{-1/2}) (I - B^{-1/2} A B^{-1/2}) \\ &= S^2 \end{aligned}$$

so we get:

$$\|R\|_B = \|S^2\| \leq \|S\|^2 \leq (1 - \varphi)^2$$

■

The approximation of the log-determinant is performed by computing sequences of power series $(R^k x)_k$. These chains are computed approximately by repeated applications of the R operator on the previous element of the chain, starting from a Gaussian white noise element $x_0 \sim \mathcal{N}(0, \mathbf{I})$. We formalize the notion of an approximate chain.

Definition 25 Approximate power sequence. *Given a linear operator H , a start point $x^{(0)} \in \mathbb{R}^n$, and a positive-definite matrix D , we define an ϵ -approximate power sequence as a sequence that does not deviate too much from the power sequence:*

$$\left\| x^{(k+1)} - Hx^{(k)} \right\|_D \leq \epsilon \left\| Hx^{(k)} \right\|_D$$

We now prove the following result that is quite intuitive: if the operator H is a contraction and if the relative error ϵ is not too great, the sum of all the errors on the chain is bounded.

Lemma 26 *Given the previous hypothesis, and assuming furthermore that $\|H\|_D \leq 1 - \eta$ and that $2\epsilon \leq \eta \leq 1/2$, the total error is bounded by $\mathcal{O}(\eta^{-3}\epsilon)$:*

$$\sum_{k=0}^{\infty} \left\| x^{(k)} - H^k x^{(0)} \right\|_D \leq 4\epsilon\eta^{-3} \left\| x^{(0)} \right\|_D$$

Proof Call $\omega_k = \left\| x^{(k)} - H^k x^{(0)} \right\|_D$ and $\theta_k = \left\| Hx^{(k)} \right\|_D$. We are going to bound the rate of convergence of these two series. We have first using triangular inequality on the D norm and then the definition of the induced matrix norm.

$$\begin{aligned} \theta_k &\leq \left\| Hx^{(k)} - H^k x^{(0)} \right\|_D + \left\| H^k x^{(0)} \right\|_D \\ &= \omega_k + \left\| H^k x^{(0)} \right\|_D \\ &\leq \omega_k + \|H\|_D^k \left\| x^{(0)} \right\|_D \end{aligned}$$

We now bound the error on the ω_k sequence:

$$\begin{aligned} \omega_{k+1} &= \left\| x^{(k+1)} - Hx^{(k)} + Hx^{(k)} - H^{k+1} x^{(0)} \right\|_D \\ &\leq \left\| Hx^{(k)} - H^{k+1} x^{(0)} \right\|_D + \left\| x^{(k+1)} - Hx^{(k)} \right\|_D \\ &\leq \|H\|_D \left\| x^{(k)} - H^k x^{(0)} \right\|_D + \epsilon \left\| Hx^{(k)} \right\|_D \\ &= \|H\|_D \omega_k + \epsilon \theta_k \\ &\leq \|H\|_D \omega_k + \epsilon \left(\omega_k + \|H\|_D^k \left\| x^{(0)} \right\|_D \right) \\ &\leq \left[(1 - \eta)^2 + \epsilon \right] \omega_k + \epsilon (1 - \eta)^{2k} \left\| x^{(0)} \right\|_D \end{aligned}$$

Note that the condition $(1 - \eta)^2 + \epsilon \leq 1 - \eta$ is equivalent to $\epsilon \leq \eta - \eta^2$. We can assume without loss of generality that $0 \leq \eta \leq 1/2$. Under this condition, $\eta/2 \leq \eta - \eta^2$. So the condition $2\epsilon \leq \eta$ implies $(1 - \eta)^2 + \epsilon \leq 1 - \eta$ which leads to:

$$\omega_{k+1} \leq (1 - \eta) \omega_k + \epsilon (1 - \eta)^{2k} \left\| x^{(0)} \right\|_D$$

By induction, one obtains:

$$\omega_k \leq \omega_1 (1 - \eta)^k + \frac{\epsilon \left\| x^{(0)} \right\|_D}{\eta} (1 - \eta)^{2k-1}$$

and since $\omega_1 = \left\| x^{(1)} - Hx^{(0)} \right\|_D \leq \epsilon \left\| Hx^{(0)} \right\|_D \leq \epsilon \|H\|_D \left\| x^{(0)} \right\|_D \leq \epsilon (1 - \eta)^2 \left\| x^{(0)} \right\|_D$, we have a final bound that depends on ϵ :

$$\omega_k \leq \epsilon (1 - \eta)^2 \left\| x^{(0)} \right\|_D (1 - \eta)^k + \frac{\epsilon \left\| x^{(0)} \right\|_D}{\eta} (1 - \eta)^{2k-1}$$

This is the sum of two geometric series, which give the bound:

$$\sum_k \omega_k \leq \epsilon \eta^{-1} \left[1 + \frac{1}{\eta^2 (1 - \eta)^2} \right] \|x^{(0)}\|_D$$

Since $\eta \leq 1/2$, it implies $(1 - \eta)^{-2} \leq 4$ and $1 \leq \eta^{-2}$, so we can further simplify:

$$\sum_k \omega_k \leq 4\epsilon \eta^{-3} \|x^{(0)}\|_D$$

■

We can use the bound on the norm of A to compute bound the error with a preconditioner:

Lemma 27 *Consider A, B with the same hypothesis as above, $x_0 \in \mathbb{R}^n$, $\epsilon \in (0, \varphi/2)$ and $(x_u)_u$ an ϵ -approximate power sequence for the operator R with start vector x_0 . Then:*

$$\left| \sum_{i=1}^l \frac{1}{i} x_0^T R^i x_0 - \sum_{i=1}^l \frac{1}{i} x_0^T x_i \right| \leq 4\epsilon \kappa^3 \kappa(B) \|x_0\|^2$$

where $\kappa(B)$ is the condition number of B .

Proof

Call \hat{z} the truncated sequence:

$$\hat{z} = \sum_{i=1}^l \frac{1}{i} (x_0)^T x_i$$

This sequence is an approximation of the exact sequence z :

$$z = \sum_{i=1}^l \frac{1}{i} x_0^T R^i x_0$$

We now bound the error between the two sequences:

$$|\hat{z} - z| \leq \sum_{i=1}^l \frac{1}{i} |x_0^T (R^i x_0 - x_i)| \leq \sum_{i=1}^l |x_0^T (R^i x_0 - x_i)| \leq \sum_{i=1}^l \left| (B^{-1} x_0)^T B (R^i x_0 - x_i) \right|$$

Using the Cauchy-Schwartz inequality, we obtain:

$$\left| (B^{-1} x_0)^T B (R^i x_0 - x_i) \right| = \left| \langle B^{-1} x_0, R^i x_0 - x_i \rangle_B \right| \leq \|B^{-1} x_0\|_B \|R^i x_0 - x_i\|_B$$

so that we can bound the deviation:

$$|\hat{z} - z| \leq \|B^{-1} x_0\|_B \sum_{i=1}^l \|R^i x_0 - x_i\|_B \leq 4\epsilon \kappa^3 \|B^{-1} x_0\|_B \|x_0\|_B \leq 4\epsilon \kappa^3 \kappa(B) \|x_0\|^2$$

where $\kappa(B)$ is the condition number of B .

■

We now have all the elements required for the proof of Theorem 7.

Proof

Consider $\mathbf{u}_j \sim \mathcal{N}(0, I_n)$ for $j = 1 \cdots p$, and $x_{i,j} = \begin{cases} \mathbf{u}_j & i = 0 \\ x_{i-1,j} - C(Ax_{i-1,j}) & i > 0 \end{cases}$

Call

$$z_{p,l} = \frac{1}{p} \sum_{j=1}^p \sum_{i=1}^l \frac{1}{i} \frac{(x_{0,j})^T x_{i,j}}{x_{0,j}^T x_{0,j}}$$

$$\hat{y}_{p,l} = \frac{1}{p} \sum_{j=1}^p \sum_{k=1}^l \frac{1}{k} \frac{\mathbf{u}_j^T S^k \mathbf{u}_j}{\mathbf{u}_j^T \mathbf{u}_j}$$

By construction, $(x_{i,j})_i$ is an ϵ -approximate chain for the operator R . Applying Lemma 27, we get:

$$|z_{p,l} - \hat{y}_{p,l}| \leq 4\epsilon\kappa^3\kappa(B) \left[\frac{1}{p} \sum_{j=1}^p \|\mathbf{u}_j\|^2 \right]$$

Call $\beta = 4p^{-1}\kappa^3\kappa(B)$. While this bound depends on the (random) value of the squared norms of the start vectors, it is enough to bound this term with high probability using a Chebyshev bound. Call $Z = \sum_{j=1}^p \|\mathbf{u}_j\|^2$ and $Y = \epsilon\beta Z$. The random variable Z is a sum of square norms of np independent scaled, centered Gaussian distributions, hence it follows a Chi-square distribution with np degrees of freedom. We have $\mathbb{E}[Z] = np$ and $\text{Cov}[Z] = 2np$, so $\mathbb{P}[|Z - np| \geq \mu] \leq 2\frac{np}{\mu^2}$.

Assume $\epsilon \leq \min\left(\frac{1}{4\kappa^3\kappa(B)}\frac{\mu}{n}, \frac{1}{4\kappa^3\kappa(B)}\frac{\mu}{n\sqrt{p}}\eta, \frac{1}{2\kappa}\right)$. Then $\mathbb{E}[Y] = 4p^{-1}\epsilon\kappa^3\kappa(B)np \leq \mu$, which means that

$$\begin{aligned} \mathbb{P}[Y \geq \mu] &\leq \mathbb{P}[|Y - \mathbb{E}[Y]| \geq \mu] \\ &\leq \frac{\text{Var}(Y)}{\mu^2} \\ &= \epsilon^2 \frac{\beta^2 np}{\mu^2} \end{aligned}$$

and our choice of ϵ ensures that $\epsilon^2 \frac{\beta^2 np}{\mu^2} \leq \eta$. Finally, we can simplify the bound on ϵ since $\frac{\eta}{\sqrt{p}} \leq 1$.

Thus, we have $\mathbb{P}[|z_{p,l} - \hat{y}_{p,l}| \geq \mu] \leq \mathbb{P}[Y \geq \mu] \leq \eta$ ■

Proofs of section 3.1

We put here the proofs that pertain to Section 3.1.

This lemma generalizes Lemma 8.1 in Spielman and Teng (2009).

Lemma 28 *Consider $A \in \mathcal{S}_n$ positive semi-definite, and $x \in \mathbb{R}^n$. Then $xx^T \preceq (x^T A^+ x) A$*

Proof Without loss of generality, consider $x^T x = 1$. Consider the eigenvalue decomposition of A : $A = \sum_i \lambda_i u_i u_i^T$. Since $(u_i)_i$ is an orthonormal basis of \mathbb{R}^n , we only need to establish that

$(u_i^T x)^2 \leq (x^T A^+ x) u_i^T A u_i$ for all i . The latter term can be simplified:

$$\begin{aligned} (x^T A^+ x) u_i^T A u_i &= \left(x^T \left[\sum_j \lambda_j^{-1} u_j u_j^T \right] x \right) \lambda_i \\ &= \lambda_i \sum_j \lambda_j^{-1} (u_j^T x)^2 \\ &\geq (u_i^T x)^2 \end{aligned}$$

which is the inequality we wanted. ■

Lemma 29 *Jensen inequality for the matrix logarithm.* Let $A \in \mathcal{S}_n$ be a positive semi-definite matrix with p positive eigenvalues. Then

$$\text{ld}(A) \leq p \log \left(\frac{\text{Tr}(A)}{p} \right)$$

Proof This is a direct application of Jensen's inequality. Call $(\lambda_i)_i$ the positive eigenvalues of A . Then $\text{ld}(A) = \sum_i \log \lambda_i$. By concavity of the logarithm:

$$\sum_i \log \lambda_i \leq p \log \left(\frac{\sum_i \lambda_i}{p} \right) = p \log \left(\frac{\text{Tr}(A)}{p} \right)$$
■

References

- Ittai Abraham, Yair Bartal, and Ofer Neiman. Nearly tight low stretch spanning trees. *Foundations of Computer ...*, pages 781–790, 2008. URL http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=4691010.
- Noga Alon, RM Karp, D Peleg, and Douglas West. A graph-theoretic game and its application to the k-server problem. *SIAM Journal on Computing*, 24(1):78–100, 1995. doi: 10.1137/S0097539792224474. URL <http://epubs.siam.org/doi/abs/10.1137/S0097539792224474>. <http://epubs.siam.org/doi/pdf/10.1137/S0097539792224474>.
- Peter W Atkins and Ronald S Friedman. *Molecular quantum mechanics*. Oxford university press, 2011.
- Zhaojun Bai, Mark Fahey, and Gene H. Golub. Some large-scale matrix computation problems. *Journal of Computational and Applied ...*, 74(1):71–89, 1996. doi: 10.1.1.56.8150. URL <http://www.sciencedirect.com/science/article/pii/0377042796000180>.
- Ronald Paul Barry and R. Kelley Pace. Monte Carlo estimates of the log determinant of large sparse matrices. *Linear Algebra and its Applications*, 1999. URL <http://www.sciencedirect.com/science/article/pii/S002437959710009X>.

- Shannon Bernardson, Paul McCarty, and Chris Thron. Monte carlo methods for estimating linear combinations of inverse matrix entries in lattice {QCD}. *Computer Physics Communications*, 78 (3):256 – 264, 1994. ISSN 0010-4655. doi: [http://dx.doi.org/10.1016/0010-4655\(94\)90004-3](http://dx.doi.org/10.1016/0010-4655(94)90004-3). URL <http://www.sciencedirect.com/science/article/pii/0010465594900043>.
- Erik G. Boman, Doron Chen, Bruce Hendrickson, and Sivan Toledo. Maximum-weight-basis preconditioners. *Numerical Linear Algebra with Applications*, 11(89):695–721, October 2004. ISSN 1070-5325. doi: 10.1002/nla.343. URL <http://doi.wiley.com/10.1002/nla.343>.
- Philippe de Forcrand and Rajan Gupta. Multigrid techniques for quark propagator. *Nuclear Physics B - Proceedings Supplements*, 9(0):516 – 520, 1989. ISSN 0920-5632. doi: [http://dx.doi.org/10.1016/0920-5632\(89\)90153-9](http://dx.doi.org/10.1016/0920-5632(89)90153-9). URL <http://www.sciencedirect.com/science/article/pii/0920563289901539>.
- A Duncan, E Eichten, and H Thacker. Efficient algorithm for qcd with light dynamical quarks. *Physical Review D*, 59(1):014505, 1998.
- Keith D. Gremban. *Combinatorial Preconditioners for Sparse, Symmetric, Diagonally Dominant Linear Systems*. PhD thesis, Carnegie Mellon University, 1996. URL www.cs.cmu.edu/~glmiller/Publications/GrembanPHD.ps.gz.
- Ilse C F Ipsen and Dean J Lee. Determinant approximations. *Numerical Linear Algebra with Applications (under ...)*, (X), 2006. URL <http://www.ncsu.edu/crsc/reports/ftp/pdf/crsc-tr03-30.pdf>.
- Ioannis Koutis, Gary L Miller, and Richard Peng. Approaching optimality for solving SDD linear systems. pages 1–16, 2010.
- Runze Li and Agus Sudjianto. Analysis of computer experiments using penalized likelihood in gaussian kriging models. *Technometrics*, 47(2), 2005.
- J Liu. The Role of Elimination Trees in Sparse Factorization. *SIAM Journal on Matrix Analysis and Applications*, 11(1):134–172, 1990. doi: 10.1137/0611010. URL <http://epubs.siam.org/doi/abs/10.1137/0611010>.
- Per-Olov Löwdin. Quantum theory of many-particle systems. iii. extension of the hartree-fock scheme to include degenerate systems and correlation effects. *Physical review*, 97(6):1509, 1955.
- R J Martin. Approximations to the determinant term in Gaussian maximum likelihood estimation of some spatial models. *Communications in Statistics-Theory and Methods*, 22(1):189–205, 1992.
- M McCourt. A Stochastic Simulation for Approximating the log-Determinant of a Symmetric Positive Definite Matrix. *compare*, 2:1–10, 2008. URL <http://www.thefutureofmath.com/mathed/logdet.pdf>.
- Gérard A Meurant. *Computer Solution of Large Linear Systems*. North-Holland: Amsterdam, 1999.
- R. Kelley Pace and Ronald Barry. Sparse spatial autoregressions. *Statistics and Probability Letters*, 33(3):291 – 297, 1997. ISSN 0167-7152. doi: [http://dx.doi.org/10.1016/S0167-7152\(96\)00140-X](http://dx.doi.org/10.1016/S0167-7152(96)00140-X). URL <http://www.sciencedirect.com/science/article/pii/S016771529600140X>.
- Arnold Reusken. Approximation of the Determinant of Large Sparse Symmetric Positive Definite Matrices. *SIAM Journal on Matrix Analysis and Applications*, 23(3):799, 2002. ISSN 08954798. doi: 10.1137/S089547980036869X. URL <http://link.aip.org/link/SJMAEL/v23/i3/p799/s1&Agg=doi>.

- Daniel A Spielman. Algorithms , Graph Theory , and Linear Equations in Laplacian Matrices. Technical report, Proceedings of the International Congress of Mathematicians, Hyderabad, India, 2010.
- Daniel A Spielman and Nikhil Srivastava. Graph Sparsification by Effective Resistances. *SIAM Journal on Computing*, 40(6):1913–1926, 2011. doi: 10.1137/080734029. URL <http://epubs.siam.org/doi/abs/10.1137/080734029>.
- Daniel A Spielman and Shang-Hua Teng. A Local Clustering Algorithm for Massive Graphs and its Application to Nearly-Linear Time Graph Partitioning. 2008.
- Daniel A Spielman and Shang-Hua Teng. Nearly-Linear Time Algorithms for Preconditioning and Solving Symmetric , Diagonally Dominant Linear Systems. pages 1–48, 2009.
- Daniel A Spielman and Jaeh Woo. A Note on Preconditioning by Low-Stretch Spanning Trees. pages 1–4, 2009.
- M.J. Wainwright and M.I. Jordan. Log-determinant relaxation for approximate inference in discrete Markov random fields. *IEEE Transactions on Signal Processing*, 54(6):2099–2109, June 2006. ISSN 1053-587X. doi: 10.1109/TSP.2006.874409. URL <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=1634807>.
- Hao Zhang and Yong Wang. Kriging and cross-validation for massive spatial data. *Environmetrics*, 21(3-4):290–304, 2010.
- Y. Zhang and W. E. Leithead. Approximate implementation of the logarithm of the matrix determinant in Gaussian process regression. *Journal of Statistical Computation and Simulation*, 77(4):329–348, April 2007. ISSN 0094-9655. doi: 10.1080/10629360600569279. URL <http://www.tandfonline.com/doi/abs/10.1080/10629360600569279>.
- Yunong Zhang, W.E. Leithead, D.J. Leith, and L. Walshe. Log-det approximation based on uniformly distributed seeds and its application to Gaussian process regression. *Journal of Computational and Applied Mathematics*, 220(1-2): 198–214, October 2008. ISSN 03770427. doi: 10.1016/j.cam.2007.08.012. URL <http://linkinghub.elsevier.com/retrieve/pii/S0377042707004360>.