

# Fast spatial Gaussian process maximum likelihood estimation via skeletonization factorizations

Victor Minden\*    Anil Damle<sup>†</sup>    Kenneth L. Ho<sup>‡</sup>    Lexing Ying<sup>§</sup>

## Abstract

Maximum likelihood estimation for parameter-fitting given observations from a Gaussian process in space is a computationally-demanding task that restricts the use of such methods to moderately-sized datasets. We present a framework for unstructured observations in two spatial dimensions that allows for evaluation of the log-likelihood and its gradient (*i.e.*, the score equations) in  $\tilde{O}(n^{3/2})$  time under certain assumptions, where  $n$  is the number of observations. Our method relies on the skeletonization procedure described by Martinsson & Rokhlin [30] in the form of the recursive skeletonization factorization of Ho & Ying [24]. Combining this with an adaptation of the matrix peeling algorithm of Lin et al. [28] for constructing  $\mathcal{H}$ -matrix representations of black-box operators, we obtain a framework that can be used in the context of any first-order optimization routine to quickly and accurately compute maximum-likelihood estimates.

## 1 Introduction

Gaussian processes are commonly used in the applied sciences as a statistical model for spatially-indexed observations. In such applications, each observation  $z_i \in \mathbb{R}$  of some quantity is associated with a corresponding location  $x_i \in \Omega \subset \mathbb{R}^d$  with  $d = 2$  or  $3$ . Given a prescribed

---

\*Institute for Computational and Mathematical Engineering, Stanford University, Stanford, CA 94305 ([vminden@stanford.edu](mailto:vminden@stanford.edu)). **Funding:** Stanford Graduate Fellowship in Science & Engineering and U.S. Department of Energy Computational Science Graduate Fellowship (grant number DE-FG02-97ER25308).

<sup>†</sup>Department of Computer Science, Cornell University, Ithaca, NY 14850 ([damle@cornell.edu](mailto:damle@cornell.edu)). **Funding:** Simons Graduate Research Assistantship and National Science Foundation Graduate Research Fellowship (grant number DGE-1147470).

<sup>‡</sup>Department of Mathematics, Stanford University, Stanford, CA 94305. Current address: TSMC Technology Inc., 2851 Junction Ave., San Jose, CA 95134. ([klho@alumni.caltech.edu](mailto:klho@alumni.caltech.edu)). **Funding:** National Science Foundation Mathematical Sciences Postdoctoral Research Fellowship (grant number DMS-1203554).

<sup>§</sup>Department of Mathematics and Institute for Computational and Mathematical Engineering, Stanford University, Stanford, CA 94305 ([lexing@math.stanford.edu](mailto:lexing@math.stanford.edu)). **Funding:** National Science Foundation (grant number DMS-1328230 and DMS-1521830) and U.S. Department of Energy Advanced Scientific Computing Research program (grant number DE-FC02-13ER26134 and de-sc0009409).

covariance kernel  $K(\cdot, \cdot; \theta)$  that maps  $\mathbb{R}^d \times \mathbb{R}^d$  to  $\mathbb{R}$  and is specified up to some parameter vector  $\theta \in \mathbb{R}^p$ , any vector of observations  $z = [z_1, \dots, z_n] \in \mathbb{R}^n$  (with associated locations  $\{x_i\}_{i=1}^n$ ) is assumed to be randomly distributed as a multivariate Gaussian

$$z \sim N(0, \Sigma(\theta)) \text{ with covariances } [\Sigma(\theta)]_{ij} = K(x_i, x_j; \theta). \quad (1)$$

We assume the mean of the process to be 0 (*i.e.*, known) for simplicity though, as we discuss later, this is not strictly necessary.

Neglecting  $\theta$ , common choices for the covariance kernel include the family of rational quadratic kernels

$$K(x, y) = \left(1 + \frac{\|x - y\|^2}{2\alpha}\right)^{-\alpha}, \quad (2)$$

which has corresponding processes that are infinitely differentiable in the mean-squared sense for all  $\alpha > 0$ . Notably, this family includes the Gaussian or squared exponential kernel as a limiting case as  $\alpha \rightarrow \infty$ . Another popular family is the Matérn family of kernels

$$K(x, y) = \frac{1}{\Gamma(\nu)2^{\nu-1}} \left(\sqrt{2\nu} \cdot \|x - y\|\right)^\nu K_\nu \left(\sqrt{2\nu} \cdot \|x - y\|\right), \quad (3)$$

where  $K_\nu(\cdot)$  is the modified second-kind Bessel function of order  $\nu$ ,  $\Gamma(\cdot)$  is the gamma function, and the corresponding process is  $\lfloor \nu \rfloor$  times mean-squared differentiable. To explicitly parameterize  $K(\cdot, \cdot; \theta)$ , we might introduce a correlation length parameter  $\rho$  leading to

$$K(x, y; \theta) = \left(1 + \frac{\|x - y\|^2}{2\alpha\rho^2}\right)^{-\alpha}$$

in the case of the rational quadratic kernel. The fundamental parameters of the kernel family *e.g.*,  $\alpha$  or  $\nu$ , may be considered as part of  $\theta$  or fixed *a priori*.

Typically, the parameter vector  $\theta$  is unknown and must be estimated from the data. For example, given a parameterized family of kernels and a set of observations, we might want to infer  $\theta$  for later use in estimating the value of the field at other spatial locations as in kriging (see Stein [37]). In this paper we consider the general Gaussian process maximum likelihood estimation (MLE) problem for  $\theta$ : given an observation vector  $z \in \mathbb{R}^n$ , find  $\hat{\theta}_{\text{MLE}}$  maximizing the Gaussian process log-likelihood

$$\ell(\theta) \equiv -\frac{1}{2}z^T \Sigma^{-1}z - \frac{1}{2} \log |\Sigma| - \frac{n}{2} \log 2\pi, \quad (4)$$

where we have dropped the explicit dependence of  $\Sigma$  on  $\theta$  for notational convenience.

If  $\theta$  is unconstrained, then  $\hat{\theta}_{\text{MLE}}$  is given by maximizing (4) over all of  $\mathbb{R}^p$ . In this case, it is possible under certain assumptions to obtain the maximum likelihood estimate by solving the score equations  $g(\theta) = 0$ , where the gradient of the log-likelihood  $g(\theta) \in \mathbb{R}^p$  is given component-wise by

$$g_i \equiv \frac{\partial \ell(\theta)}{\partial \theta_i} = \frac{1}{2}z^T \Sigma^{-1} \Sigma_i \Sigma^{-1}z - \frac{1}{2} \text{Tr}(\Sigma^{-1} \Sigma_i), \quad i = 1, \dots, p, \quad (5)$$

where  $\Sigma_i \equiv \frac{\partial}{\partial \theta_i} \Sigma$ . This may be accomplished without evaluating  $\ell(\theta)$  as is done by, *e.g.*, Anitescu et al. [2] and Stein et al. [39]. In contrast, we consider in this paper the use of first-order methods for nonlinear optimization that, at each iteration, use both gradient and log-likelihood evaluations to find a local optimum. This allows for the treatment of constraints if desired, though we note the methods here are equally applicable to the unconstrained case.

For any given  $\theta$ , both  $\ell(\theta)$  and the gradient  $g(\theta)$  contain a number of terms whose evaluation is traditionally computationally expensive. For example, the Cholesky decomposition of  $\Sigma$  may be used to calculate  $z^T \Sigma^{-1} z$  and  $z^T \Sigma^{-1} \Sigma_i \Sigma^{-1} z$  as well as the log-determinant  $\log |\Sigma|$  and trace  $\text{Tr}(\Sigma^{-1} \Sigma_i)$ , but the asymptotic computation and storage complexities are  $O(n^3)$  and  $O(n^2)$  respectively. This is prohibitively expensive for datasets with a large number of observations, necessitating alternative approaches.

## 1.1 Our method

The contribution of this paper is a framework for efficiently finding  $\hat{\theta}_{\text{MLE}}$  by taking advantage of fast hierarchical matrix algorithms developed in the numerical linear algebra community. Such algorithms exploit the fact that linear operators defined in terms of pairwise kernel evaluations between points embedded in  $\mathbb{R}^d$  frequently exhibit *hierarchical rank structure*, as we discuss in Section 2 (briefly, many different-sized off-diagonal blocks of the matrix are close to low-rank and thus compressible). Our framework has two key parts:

- (I) Construct a fast approximate hierarchical factorization of the covariance matrix  $\Sigma$  and its derivatives  $\Sigma_i$  for  $i = 1, \dots, p$  in a matrix-free fashion using the kernel function  $K(\cdot, \cdot; \theta)$  and the points  $\{x_i\}_{i=1}^n$ .
- (II) Additionally factor the derivatives of the covariance matrix  $\Sigma_i$  for  $i = 1, \dots, p$  through the same approach as in (I). Use the hierarchical factorizations as a fast black-box operator for approximately applying  $\Sigma^{-1} \Sigma_i$  for  $i = 1, \dots, p$ , and use this operator and the points  $\{x_i\}_{i=1}^n$  to compute the traces  $\text{Tr}(\Sigma^{-1} \Sigma_i)$  through a randomized “matrix peeling” scheme for efficiently extracting the trace of a hierarchically rank-structured linear operator.

In both parts, the approximation accuracy is well-controlled by specified tolerance parameters intrinsic to the algorithms.

For any  $\theta$ , the approximate factorization of part (I) can be used to efficiently evaluate the terms composing  $\ell(\theta)$  (including the log-determinant), which overlaps with recent work by Ambikasaran et al. [1] that addresses the use of hierarchical factorizations for kernelized covariance matrices for computing these terms (see also Khoromskij et al. [27] and Börm & Garcke [4] for earlier work on  $\mathcal{H}$ -matrix techniques for fast computation of matrix-vector products with  $\Sigma$  in a Gaussian process context). This piece of the framework alone gives sufficient machinery to perform black-box optimization using numerical derivatives. However, using finite differences of an approximate log-likelihood can magnify approximation errors, which can lead to larger inaccuracies in the approximate gradient depending on, *e.g.*, the conditioning of  $\Sigma$  (see Section 5.4 for a relatively benign example). Therefore, central to our

framework is the computation in (II) of the gradient components  $g_i$  for  $i = 1, \dots, p$ , which requires the trace terms  $\text{Tr}(\Sigma^{-1}\Sigma_i)$ .

In the simplest form detailed in this paper, our framework employs the *recursive skeletonization factorization* [24] as the approximate hierarchical factorization variant of choice for  $\Sigma$  and  $\Sigma_i$ ,  $i = 1, \dots, p$ . We then compute the trace terms in the gradient using an adaptation of the matrix peeling algorithm of Lin et al. [28]. Combining these two tools, we obtain an efficient method for evaluating (4) and (5)—and, ultimately, finding  $\hat{\theta}_{\text{MLE}}$ —with high and controllable accuracy using a black-box first-order optimization package (*e.g.*, `fminunc` or `fmincon` in MATLAB<sup>®</sup>).

While the framework of this paper technically applies to observations in  $d$  dimensions for general  $d$ , the computational complexity increases in high dimensions due to how the runtime of hierarchically rank-structured factorizations depends on the numerical rank of off-diagonal blocks. For example, applying these methods in the case  $d = 1$  is essentially optimal in the sense that off-diagonal matrix blocks have numerical rank that is not strongly dependent on the number of observations. We direct the reader to Ambikasaran et al. [1] for extensive numerical examples of factoring kernel matrices in this case. In contrast, for  $d = 3$  the observed rank growth is in general much larger and leads to greater asymptotic complexities. We focus on the case  $d = 2$  in the remainder of this paper, but note the broader applicability.

## 1.2 Alternative approaches

Due to the prevalence of Gaussian process models, a number of methods exist in the literature for fast computations involving kernelized covariance matrices. To decrease apply, solve, and storage costs,  $\Sigma$  can be replaced with a sparser “tapered” approximant as described by Furrer et al. [12], wherein the desired covariance kernel is multiplied pointwise with a compactly-supported tapering function to attain sparsity. Of course, the computational benefit of tapering depends on the sparsity of the resulting approximant, which is limited by the desired accuracy if the correlation length of the kernel is not small.

If  $\Sigma$  decomposes naturally into the sum of a diagonal and a numerically low-rank matrix then such a decomposition can be quite efficient for computation (see Cressie & Johannesson [10]), but this representation is too simple for the applications and kernel functions we consider. Extending this to a general sparse-plus-low-rank model by replacing the diagonal piece with a tapered covariance kernel can perform better than either a tapered or low-rank  $\Sigma$  alone, as shown by Sang & Huang [35] (see also Vanhatal et al. [42]).

For cases where the underlying process is stationary and the observations lie on a regular grid it is possible to directly approximate the log-likelihood using spectral approximations due to Whittle [44] or to quickly apply the covariance matrix in Fourier space to solve linear systems with an iterative method. Further, in such cases these systems can be preconditioned using the method of Stein et al. [38] yielding efficient methods for many important problem classes. For irregularly spaced data such as we consider in this paper, however, these approaches do not apply directly. One approach for log-likelihood evaluation (and thus derivative-free optimization) with generally distributed data is that of Aune et al. [3],

which offers an involved framework for approximating the log-determinant using Krylov methods, assuming the covariance matrix or its inverse can be efficiently applied. More recently, Castrillón-Candás et al. [5] demonstrate a combination of multi-level preconditioning and tapering for fast derivative-free restricted maximum likelihood estimation, though gradient computation is not discussed.

An alternative approach to approximating the Gaussian process log-likelihood directly is to explicitly construct and solve a different set of estimating equations that is less computationally cumbersome. For example, the Hutchinson-like sample average approximation (SAA) estimator [2, 39] falls into this category, as do the composite likelihood methods described by, *e.g.*, Vecchia [43] and Stein et al. [40] and their extension, block composite likelihood methods (see, *e.g.*, Eidsvik et al. [11]). Another notable effort based on a modified model is the work of Lindgren et al. [29], which gives a way of approximating Gaussian processes by Gaussian Markov random fields for specific kernels in the Matérn family. In practice, these methods perform quite well, but in our approach, we consider maximizing the true likelihood as opposed to alternative models or estimators.

### 1.3 Outline

The remainder of the paper is organized as follows. In Section 2 we review hierarchical matrix structure and outline the recursive skeletonization factorization [24, section 3], which is our hierarchical matrix format of choice for fast MLE. In Section 3, we discuss a modification of the matrix peeling algorithm by Lin et al. [28], which we use to quickly evaluate the gradient of the log-likelihood. In Section 4, we succinctly summarize our framework. In Section 5 we present numerical results on a number of test problems and demonstrate the scaling of our approach. Finally, in Section 6, we make some concluding remarks.

## 2 Factorization of the covariance matrix

Consider the kernelized covariance matrix  $\Sigma$  as in (1), and assume for simplicity of exposition that the points  $\{x_i\}_{i=1}^n$  are uniformly distributed inside a rectangular domain  $\Omega$ . Partitioning  $\Omega$  into four equal rectangular subdomains  $\Omega_{1;i}$  for  $i = 1, \dots, 4$ , it has been observed that the corresponding block partitioning of  $\Sigma$  exposes low-rank structure of off-diagonal blocks when  $K(x, y; \theta)$  is sufficiently nice as a function of  $\|x - y\|$ .

Concretely, we define the set  $[n] \equiv \{1, 2, \dots, n\}$  and let  $\mathcal{I}_{1;i} \subset [n]$  denote the index set indexing degrees of freedoms (DOFs) located inside  $\Omega_{1;i}$  for  $i = 1, \dots, 4$  such that  $\{x_j\}_{j \in \mathcal{I}_{1;i}} \subset \Omega_{1;i}$  and  $\bigcup_{i=1}^4 \mathcal{I}_{1;i} = [n]$ . In a self-similar fashion, we further partition  $\Omega_{1;1}$  into the four subdomains  $\Omega_{2;i}$  (with corresponding DOFs  $\mathcal{I}_{2;i}$ ) for  $i = 1, \dots, 4$  and obtain the decomposition shown in Fig. 1 (left). Assuming that the covariance kernel is smooth away from  $x = y$  and does not exhibit high-frequency oscillations, the off-diagonal blocks  $\Sigma(\mathcal{I}_{1;i}, \mathcal{I}_{1;j})$  for  $i \neq j$  and  $\Sigma(\mathcal{I}_{2;i}, \mathcal{I}_{2;j})$  for  $i \neq j$  in the corresponding partitioning shown in the same figure (right) tend to be *numerically low-rank* and thus compressible.

**Definition 1** (Numerically low-rank). We call a matrix  $\mathbf{A} \in \mathbb{R}^{m_1 \times m_2}$  numerically low-rank with respect to a specified tolerance  $\epsilon$  if for some  $r < \min(m_1, m_2)$  there exist matrices  $\mathbf{V}_1 \in \mathbb{R}^{m_1 \times r}$  and  $\mathbf{V}_2 \in \mathbb{R}^{m_2 \times r}$  such that  $\|\mathbf{A} - \mathbf{V}_1 \mathbf{V}_2^T\|_2 \leq \epsilon \|\mathbf{A}\|_2$ .

The rank structure of Fig. 1 includes numerically low-rank blocks at multiple spatial scales independent of the length-scale of the underlying Gaussian process, i.e., we may continue to recursively subdivide the domain and expose more compressible blocks of  $\Sigma$ . Explicitly representing each of these blocks in low-rank form leads to the so-called hierarchical off-diagonal low-rank (HODLR) matrix format that has been used by Ambikasaran et al. [1] to compress various families of covariance kernels with application to Gaussian processes.

**Remark 1.** A sufficient condition to ensure this rank structure is that for any pair of distinct subdomains on the same level  $\Omega_{\ell;i}$  and  $\Omega_{\ell;j}$  there exists an approximation  $K(x_i, x_j; \theta) \approx \sum_{k=1}^r f_k(x_i)g_k(x_j)$  for any  $x_i \in \mathcal{I}_{\ell;i}$  and  $x_j \in \mathcal{I}_{\ell;j}$  where the number of terms  $r$  is relatively small (via, e.g., Chebyshev polynomials). An important example where this is not typically the case is periodic kernels with short period relative to the size of the domain. Further, the techniques we discuss here are less relevant to compactly-supported kernels, which are typically already efficient to compute with using standard sparse linear algebra.

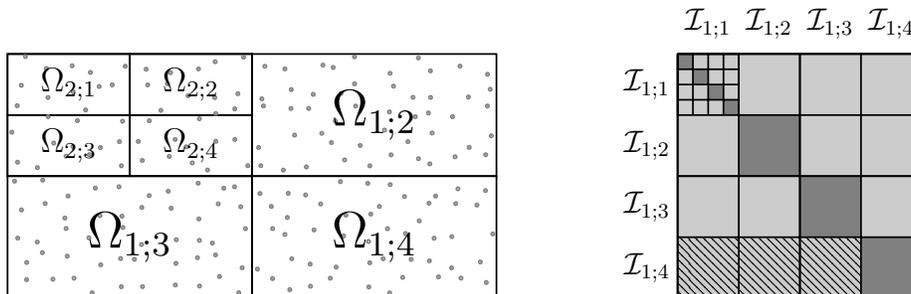


Figure 1: To expose low-rank structure in the matrix  $\Sigma$ , consider the partitions  $\Omega = \cup_{i=1}^4 \Omega_{1;i}$  and  $\Omega_{1;1} = \cup_{i=1}^4 \Omega_{2;i}$  (left), with corresponding index sets  $\mathcal{I}_{\ell;i}$  for  $\ell = 1, \dots, 2$ ,  $i = 1, \dots, 4$ . Due to smoothness properties of  $K(\cdot, \cdot; \theta)$ , the blocks of  $\Sigma$  are rank-structured (right), where dark gray blocks on the diagonal are full-rank and light gray off-diagonal blocks are numerically low-rank. Representing each off-diagonal block in low-rank form independently leads to the HODLR format. In contrast, the HBS format uses one low-rank representation for all off-diagonal blocks for each block row, e.g., the patterned blocks in the bottom row are aggregated into a single low-rank matrix.

The HODLR format is only one of many hierarchical matrix formats, appearing as a special case of the  $\mathcal{H}$ - and  $\mathcal{H}^2$ -matrices of Hackbusch and collaborators [17, 20, 19]. This format is particularly simple as at each level it compresses *all off-diagonal blocks* of  $\Sigma$ , including those corresponding to domains that are adjacent (e.g.,  $\Omega_{1;2}$  and  $\Omega_{1;3}$ ). Matrices compressible in this way are referred to as *weakly-admissible* [21], in contrast to *strongly-admissible* matrices which compress only a subset of off-diagonal blocks at each level. Closely related literature includes work on hierarchically semiseparable (HSS) matrices [46, 7, 6] and

hierarchically block separable (HBS) matrices [30, 13] which offer simplified representations for weakly-admissible matrices with improved runtime.

For matrices where the entries are explicitly generated by an underlying kernel function such as  $\Sigma$  and its derivatives, specific factorization algorithms have been developed to exploit this additional structure for increased efficiency [15, 13, 23, 24, 9]. These “skeletonization-based” algorithms, based on the framework introduced by Martinsson & Rokhlin [30] stemming from observations by Starr & Rokhlin [36] and Greengard & Rokhlin [16], construct low-rank representations of certain off-diagonal blocks using the skeletonization process described by Cheng et al. [8].

Our framework is agnostic to the choice of hierarchical factorization used for  $\Sigma$  and its derivatives, provided that the factorization admits fast linear algebra computations (including computation of the log-determinant) with the underlying operator. The recursive skeletonization factorization that we use in this paper was first introduced by Ho & Ying [24, section 3] as a multiplicative factorization based on skeletonization [30]. In the remainder of this section we provide a brief review of the algorithm.

**Remark 2.** *In what follows, we will continue to assume that the DOFs  $\{x_i\}_{i=1}^n$  are uniformly distributed inside a rectangular subdomain  $\Omega$  for simplicity of exposition. Further, we will describe the algorithm as though the quadtree representing the hierarchical partitioning of space is perfect, i.e., every subdomain is subdivided into four child subdomains at every level. In practice an adaptive decomposition of space is used to avoid subdividing the domain in regions of low observation density.*

## 2.1 Block compression through skeletonization

We begin by recursively subdividing  $\Omega$  into four subdomains until each leaf-level subdomain contains a constant number of observations independent of  $n$ . This leads to a quadtree data structure with levels labeled  $\ell = 0$  through  $\ell = L$ , where  $\ell = 0$  refers to the entire domain  $\Omega$  and  $\ell = L$  refers to the collection of subdomains  $\Omega_{L;i}$  for  $i = 1, \dots, 4^L$ . Considering factorization of the covariance matrix  $\Sigma$ , the basic intuition of the method is to first compress all blocks of  $\Sigma$  corresponding to covariances between observations in distinct subdomains at the leaf level, and then to recurse in a bottom-up traversal.

Consider first a single leaf-level subdomain containing observations indexed by  $\mathcal{I} \subset [n]$  and define the complement DOF set  $\mathcal{I}^c \equiv [n] \setminus \mathcal{I}$ . Given a specified tolerance  $\epsilon > 0$ , the algorithm proceeds by compressing the off-diagonal blocks  $\Sigma(\mathcal{I}, \mathcal{I}^c)$  and  $\Sigma(\mathcal{I}^c, \mathcal{I})$  as in the HBS format (see Fig. 1) through the use of an *interpolative decomposition* (ID) [8].

**Definition 2** (Interpolative decomposition). *Given a matrix  $\mathbf{A} \in \mathbb{R}^{m \times |\mathcal{I}|}$  with columns indexed by  $\mathcal{I}$  and a tolerance  $\epsilon > 0$ , an  $\epsilon$ -accurate interpolative decomposition of  $\mathbf{A}$  is a partitioning of  $\mathcal{I}$  into DOF sets associated with so-called skeleton columns  $\mathcal{S} \subset \mathcal{I}$  and redundant columns  $\mathcal{R} = \mathcal{I} \setminus \mathcal{S}$  and a corresponding interpolation matrix  $\mathbf{T}_{\mathcal{I}}$  such that  $\|\mathbf{A}(:, \mathcal{R}) - \mathbf{A}(:, \mathcal{S})\mathbf{T}_{\mathcal{I}}\|_2 \leq \epsilon \|\mathbf{A}\|_2$ , where  $\mathbf{A}(:, \mathcal{R}) \in \mathbb{R}^{m \times |\mathcal{R}|}$  is given by subselecting the columns of  $\mathbf{A}$  indexed by  $\mathcal{R}$ , and  $\mathbf{A}(:, \mathcal{S})$  is defined analogously.*

It is desirable in Definition 2 to take  $|\mathcal{S}|$  as small as possible for a given  $\epsilon$ .

Given an ID of  $\Sigma(\mathcal{I}^c, \mathcal{I})$  such that  $\Sigma(\mathcal{I}^c, \mathcal{R}) \approx \Sigma(\mathcal{I}^c, \mathcal{S})\mathsf{T}_{\mathcal{I}}$ ,  $\Sigma$  can be written in block form (up to a permutation) as

$$\left[ \begin{array}{cc|c} \Sigma(\mathcal{I}^c, \mathcal{I}^c) & \Sigma(\mathcal{I}^c, \mathcal{S}) & \Sigma(\mathcal{I}^c, \mathcal{R}) \\ \Sigma(\mathcal{S}, \mathcal{I}^c) & \Sigma(\mathcal{S}, \mathcal{S}) & \Sigma(\mathcal{S}, \mathcal{R}) \\ \hline \Sigma(\mathcal{R}, \mathcal{I}^c) & \Sigma(\mathcal{R}, \mathcal{S}) & \Sigma(\mathcal{R}, \mathcal{R}) \end{array} \right] \approx \left[ \begin{array}{cc|c} \Sigma(\mathcal{I}^c, \mathcal{I}^c) & \Sigma(\mathcal{I}^c, \mathcal{S}) & \Sigma(\mathcal{I}^c, \mathcal{S})\mathsf{T}_{\mathcal{I}} \\ \Sigma(\mathcal{S}, \mathcal{I}^c) & \Sigma(\mathcal{S}, \mathcal{S}) & \Sigma(\mathcal{S}, \mathcal{R}) \\ \hline \mathsf{T}_{\mathcal{I}}^T \Sigma(\mathcal{S}, \mathcal{I}^c) & \Sigma(\mathcal{R}, \mathcal{S}) & \Sigma(\mathcal{R}, \mathcal{R}) \end{array} \right].$$

Using a sequence of block row and column operations, we first eliminate the blocks  $\Sigma(\mathcal{I}^c, \mathcal{S})\mathsf{T}_{\mathcal{I}}$  and  $\mathsf{T}_{\mathcal{I}}^T \Sigma(\mathcal{S}, \mathcal{I}^c)$  and then decouple the bottom-right block to obtain

$$\mathsf{L}_{\mathcal{I}}^{-1} \Sigma \mathsf{L}_{\mathcal{I}}^{-T} \approx \left[ \begin{array}{cc|c} \Sigma(\mathcal{I}^c, \mathcal{I}^c) & \Sigma(\mathcal{I}^c, \mathcal{S}) & \\ \Sigma(\mathcal{S}, \mathcal{I}^c) & \Sigma(\mathcal{S}, \mathcal{S}) & \mathsf{X}_{\mathcal{S}\mathcal{R}} \\ \hline & \mathsf{X}_{\mathcal{R}\mathcal{S}} & \mathsf{X}_{\mathcal{R}\mathcal{R}} \end{array} \right] = \mathsf{U}_{\mathcal{I}} \left[ \begin{array}{cc|c} \Sigma(\mathcal{I}^c, \mathcal{I}^c) & \Sigma(\mathcal{I}^c, \mathcal{S}) & \\ \Sigma(\mathcal{S}, \mathcal{I}^c) & \mathsf{X}_{\mathcal{S}\mathcal{S}} & \\ \hline & & \mathsf{X}_{\mathcal{R}\mathcal{R}} \end{array} \right] \mathsf{U}_{\mathcal{I}}^T,$$

where  $\mathsf{L}_{\mathcal{I}}$  and  $\mathsf{U}_{\mathcal{I}}$  are block unit-triangular matrices that are fast to apply or invert and the  $\mathsf{X}$  subblocks are linear combinations of the  $\Sigma$  subblocks.

## 2.2 The recursive skeletonization factorization

Defining the collection of DOF sets corresponding to subdomains at level  $\ell = L$  as  $\mathcal{L}_L \equiv \{\mathcal{I}_{L;1}, \mathcal{I}_{L;2}, \dots, \mathcal{I}_{L;4^L}\}$ , we use the skeletonization process of Section 2.1 to compress the corresponding blocks of  $\Sigma$  for each  $\mathcal{I} \in \mathcal{L}_L$ , yielding  $\Sigma \approx (\prod_{\mathcal{I} \in \mathcal{L}_L} \mathsf{L}_{\mathcal{I}} \mathsf{U}_{\mathcal{I}}) \tilde{\Sigma}_L (\prod_{\mathcal{I} \in \mathcal{L}_L} \mathsf{U}_{\mathcal{I}}^T \mathsf{L}_{\mathcal{I}}^T)$ , where the order taken in the product over  $\mathcal{L}_L$  does not matter due to the structure of the  $\mathsf{L}$  and  $\mathsf{U}$  matrices. Using  $\mathcal{R}_{L;i}$  and  $\mathcal{S}_{L;i}$  to denote the redundant DOFs and skeleton DOFs associated with  $\mathcal{I}_{L;i}$  for each  $i = 1, \dots, 4^L$  and defining  $\mathcal{R}_L \equiv \cup_{i=1}^{4^L} \mathcal{R}_{L;i}$ , the blocks of the matrix  $\tilde{\Sigma}_L$  have the following structure for each  $i = 1, \dots, 4^L$ :

- The modified block  $\tilde{\Sigma}_L(\mathcal{R}_{L;i}, \mathcal{R}_{L;i})$  has been decoupled from the rest of  $\tilde{\Sigma}_L$ .
- The block  $\tilde{\Sigma}_L(\mathcal{S}_{L;i}, \mathcal{S}_{L;i})$  has been modified.
- The blocks  $\tilde{\Sigma}_L(\mathcal{S}_{L;i}, \mathcal{I}_{L;i}^c \setminus \mathcal{R}_L) = \Sigma(\mathcal{S}_{L;i}, \mathcal{I}_{L;i}^c \setminus \mathcal{R}_L)$  and  $\tilde{\Sigma}_L(\mathcal{I}_{L;i}^c \setminus \mathcal{R}_L, \mathcal{S}_{L;i}) = \Sigma(\mathcal{I}_{L;i}^c \setminus \mathcal{R}_L, \mathcal{S}_{L;i})$  remain unmodified from what they were in  $\Sigma$ .

In other words, we have identified and decoupled all redundant DOFs at the leaf level while leaving unchanged the blocks of  $\Sigma$  corresponding to kernel evaluations between skeleton DOFs in distinct leaf-level subdomains.

The recursive skeletonization factorization of  $\Sigma$  is given by repeating this process at each higher level of the quadtree. For example, at level  $\ell = L - 1$  of the quadtree, each subdomain contains skeleton DOFs corresponding to its four distinct child subdomains in the tree. However, the redundant DOFs of its child subdomains no longer need to be considered as they have already been decoupled. We thus define  $\tilde{\mathcal{I}}_{L-1;i} \equiv \mathcal{I}_{L-1;i} \setminus \mathcal{R}_L$  for

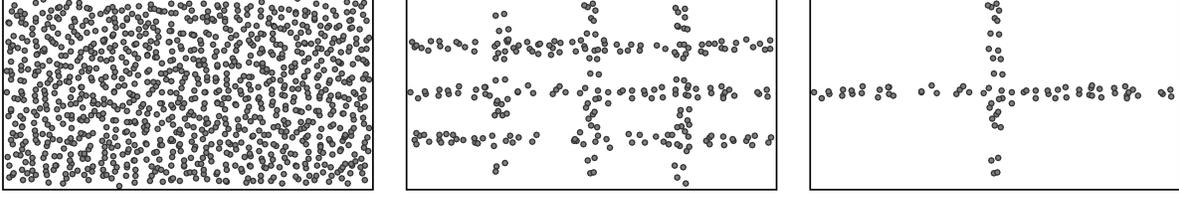


Figure 2: Shown here are the DOFs to be compressed at each level of the recursive skeletonization factorization. At the leaf level (left), all DOFs are involved in skeletonization. At each subsequent level (center, right), only skeleton DOFs from the previous level are involved in further skeletonization. We see that the skeleton DOFs tend to line the boundaries of their corresponding subdomains.

each  $i = 1, \dots, 4^{(L-1)}$  and write the collection of DOFs remaining at this level as  $\mathcal{L}_{L-1} \equiv \{\tilde{\mathcal{I}}_{L-1;1}, \tilde{\mathcal{I}}_{L-1;2}, \dots, \tilde{\mathcal{I}}_{L-1;4^{(L-1)}}\}$ . Due to the hierarchical block low-rank structure of  $\tilde{\Sigma}$ , off-diagonal blocks at this level are again compressible through skeletonization, yielding

$$\tilde{\Sigma}_L \approx \mathbf{P}_L \left( \prod_{\mathcal{I} \in \mathcal{L}_{L-1}} \mathbf{L}_{\mathcal{I}} \mathbf{U}_{\mathcal{I}} \right) \tilde{\Sigma}_{L-1} \left( \prod_{\mathcal{I} \in \mathcal{L}_{L-1}} \mathbf{U}_{\mathcal{I}}^T \mathbf{L}_{\mathcal{I}}^T \right) \mathbf{P}_L^T,$$

where  $\mathbf{P}_L$  is a global permutation matrix regrouping the DOF sets in  $\mathcal{L}_{L-1}$  to be contiguous. Proceeding in this fashion level-by-level and defining the ordered product  $\prod_{\ell=1}^L \mathbf{A}_{\ell} \equiv \mathbf{A}_L \mathbf{A}_{L-1} \dots \mathbf{A}_1$ , we obtain the full recursive skeletonization factorization  $\mathbf{F}$  of  $\Sigma$

$$\begin{aligned} \Sigma &\approx \left[ \prod_{\ell=1}^L \left( \prod_{\mathcal{I} \in \mathcal{L}_{\ell}} \mathbf{L}_{\mathcal{I}} \mathbf{U}_{\mathcal{I}} \right) \mathbf{P}_{\ell} \right] \tilde{\Sigma}_1 \left[ \prod_{\ell=1}^L \left( \prod_{\mathcal{I} \in \mathcal{L}_{\ell}} \mathbf{L}_{\mathcal{I}} \mathbf{U}_{\mathcal{I}} \right) \mathbf{P}_{\ell} \right]^T \\ &= \left[ \prod_{\ell=1}^L \left( \prod_{\mathcal{I} \in \mathcal{L}_{\ell}} \mathbf{L}_{\mathcal{I}} \mathbf{U}_{\mathcal{I}} \right) \mathbf{P}_{\ell} \right] \mathbf{C}\mathbf{C}^T \left[ \prod_{\ell=1}^L \left( \prod_{\mathcal{I} \in \mathcal{L}_{\ell}} \mathbf{L}_{\mathcal{I}} \mathbf{U}_{\mathcal{I}} \right) \mathbf{P}_{\ell} \right]^T \equiv \mathbf{F}, \end{aligned} \quad (6)$$

where  $\tilde{\Sigma}_1$  is block-diagonal with diagonal blocks corresponding to the sets of redundant DOFs at each level and  $\tilde{\Sigma}_1 = \mathbf{C}\mathbf{C}^T$  is the Cholesky decomposition of  $\tilde{\Sigma}_1$ .

**Remark 3.** *Because the factorization  $\mathbf{F} \approx \Sigma$  is approximate, using an extremely inaccurate tolerance  $\epsilon$  in the IDs of Definition 2 admits the possibility that  $\tilde{\Sigma}_1$  may be slightly indefinite due to approximation error. In practice, this is not an issue for any tolerance precise enough to be used for computing  $\ell(\theta)$  for optimization purposes. When factoring the derivative matrices  $\tilde{\Sigma}_i$  for  $i = 1, \dots, p$  (which may themselves be indefinite), the Cholesky decomposition may be replaced with, e.g., an  $\text{LDL}^T$  factorization.*

### 2.3 Computational complexity

The computational cost of the recursive skeletonization factorization is in theory dominated by the cost of computing IDs  $\Sigma(\mathcal{I}^c, \mathcal{R}) \approx \Sigma(\mathcal{I}^c, \mathcal{S}) \mathbf{T}_{\mathcal{I}}$  of  $\Sigma(\mathcal{I}^c, \mathcal{I})$  in Section 2.1 for each  $\mathcal{I}$ .

This is because the typical algorithm to compute an ID is based on a rank-revealing QR factorization, such that the  $m \times |\mathcal{I}|$  ID of Definition 2 has complexity  $O(m|\mathcal{I}|^2)$  [8]. This dependence on  $m$  is prohibitively expensive during initial levels of the factorization, since  $m = |\mathcal{I}^c| = O(n)$  and there are  $O(n)$  such IDs to compute.

The original application of skeletonization was to boundary integral equations arising from elliptic partial differential equations, in which case the so-called “proxy trick” described by Martinsson & Rokhlin [30] can be applied to accelerate the computation of an ID through the use of integral identities. These integral identities do not strictly apply in the case where  $K(\cdot, \cdot; \theta)$  is a general covariance function, but we find that a variant of this proxy trick works well to obtain similar acceleration.

### 2.3.1 Modified proxy trick

Suppose that the index set  $\mathcal{I}$  corresponds to points inside the subdomain  $B \subset \Omega$  in Fig. 3, where we use  $B$  as a stand-in for an arbitrary subdomain  $\Omega_{\ell,i}$  in our quadtree. The purpose of computing an ID of  $\Sigma(\mathcal{I}^c, \mathcal{I})$  is to find a small set of skeleton DOFs  $\mathcal{S} \subset \mathcal{I}$  such that the range of  $\Sigma(\mathcal{I}^c, \mathcal{S})$  approximately captures the range of  $\Sigma(\mathcal{I}^c, \mathcal{I})$ . The key to computational acceleration using the proxy trick is to accomplish this without using all rows of  $\Sigma(\mathcal{I}^c, \mathcal{I})$  in the computation.

As detailed by Ho & Ying [24, section 3.3], we can partition the DOFs  $\mathcal{I}^c$  into those that are near to  $B$  and those that are far from  $B$ , denoted  $\mathcal{N}$  and  $\mathcal{F}$  respectively. For example, we may take  $\mathcal{N}$  to be all points  $x_i \in \mathcal{I}^c$  such that the distance between  $x_i$  and the center of  $B$  is less than some radius. The proxy trick proceeds by finding a surrogate representation  $\mathbf{M}(\Gamma, \mathcal{I})$  for  $\Sigma(\mathcal{F}, \mathcal{I})$  in the ID computation, such that  $\mathbf{M}(\Gamma, \mathcal{I})$  has many fewer rows than  $\Sigma(\mathcal{F}, \mathcal{I})$  and

$$\begin{bmatrix} \Sigma(\mathcal{N}, \mathcal{R}) \\ \mathbf{M}(\Gamma, \mathcal{R}) \end{bmatrix} \approx \begin{bmatrix} \Sigma(\mathcal{N}, \mathcal{S}) \\ \mathbf{M}(\Gamma, \mathcal{S}) \end{bmatrix} \mathbb{T}_{\mathcal{I}} \implies \begin{bmatrix} \Sigma(\mathcal{N}, \mathcal{R}) \\ \Sigma(\mathcal{F}, \mathcal{R}) \end{bmatrix} \approx \begin{bmatrix} \Sigma(\mathcal{N}, \mathcal{S}) \\ \Sigma(\mathcal{F}, \mathcal{S}) \end{bmatrix} \mathbb{T}_{\mathcal{I}} \quad (7)$$

such that we may compute the left ID in (7) and get the right ID for “free”.

In the modified proxy trick, we let  $\Gamma$  be a set of  $n_{\text{prox}}$  points discretizing the gray annulus in the right of Fig. 3. Crucially, this differs from the original proxy trick due to the fact that we discretize a two-dimensional region (the annulus), whereas if our kernel satisfied some form of a Green’s identity we could instead discretize a quasi-one-dimensional curve (a circle) around  $B$  as in the original proxy trick. Defining the matrix  $\mathbf{M}(\Gamma, \mathcal{I})$  to have entries  $K(y, x_i; \theta)$  with rows indexed by  $y \in \Gamma$  and columns indexed by  $x_i \in \mathcal{I}$ , we observe that (7) holds without significant loss in accuracy even with  $n_{\text{prox}}$  relatively small. This brings the complexity of computing the right ID down to  $O(|\mathcal{I}|^3 + n_{\text{prox}}|\mathcal{I}|^2)$ , which is beneficial when  $n_{\text{prox}}$  is small compared to  $|\mathcal{I}^c|$ . In practice, we take  $n_{\text{prox}}$  to be constant with respect to the total number of points  $n$ .

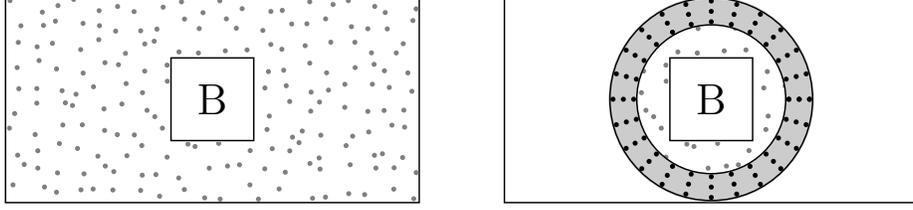


Figure 3: Because of the underlying kernel, computing an ID of the submatrix  $\Sigma(\mathcal{I}^c, \mathcal{I})$  can be accelerated, where  $\mathcal{I} \subset [n]$  indexes observations inside the subdomain B (left). Rather than considering all of  $\mathcal{I}^c$ , the (modified) *proxy trick* involves neglecting rows of  $\Sigma(\mathcal{I}^c, \mathcal{I})$  corresponding to points  $\mathcal{F} \subset \mathcal{I}^c$  far from B. Instead, the ID is computed by considering only points  $\mathcal{N} \subset \mathcal{I}^c$  near B combined with a small number of so-called “proxy points” discretizing an annulus around B (right).

### 2.3.2 Complexity sketch using the modified proxy trick

Using the modified proxy trick, the cost of the recursive skeletonization factorization is essentially determined by the number of DOFs interior to each skeletonized subdomain, *i.e.*,  $|\mathcal{S}_{\ell;i}|$  for each  $\ell = 1, \dots, L$  and  $i = 1, \dots, 4^\ell$ . As seen in Fig. 2, the skeleton DOFs tend to line the boundaries of their corresponding subdomains. This is statistically intuitive: due to the fact that our kernels of interest  $K(\cdot, \cdot; \theta)$  decay smoothly, the subset of DOFs that best represent the covariance structure of a subdomain with the rest of the domain is the subset closest to the rest of the domain.

Assuming a uniform distribution of points  $\{x_i\}_{i=1}^n$  and perfect quadtree, the average number of skeleton DOFs per subdomain at level  $\ell$  is on the order of the sidelength of a subdomain in the quadtree at level  $\ell$ . In other words, the number of skeleton DOFs per box grows by roughly a factor of two each time we step up a level in the tree and thus  $s_\ell \equiv \frac{1}{4^\ell} \sum_{i=1}^{4^\ell} |\mathcal{S}_{\ell;i}| = O(2^{-\ell})$ . The assumptions that lead to this rank growth bound are described in more detail by Ho & Greengard [23, section 4]; we do not go into them here.

**Theorem 1** ([30, 24, 23]). *Assuming that the size of the skeleton sets behaves like  $s_\ell = O(2^{-\ell})$  for  $\ell = 1, \dots, L$  and  $L \sim \log n$ , the computational complexity of the recursive skeletonization factorization  $\mathbf{F} \approx \Sigma$  (with constants depending on the tolerance  $\epsilon$  in Definition 2) is  $T_{\text{factor}} = O(n^{3/2})$  and  $T_{\text{apply}} = T_{\text{solve}} = O(n \log n)$ , where  $T_{\text{factor}}$  is the complexity of the factorization and  $T_{\text{apply}}$  and  $T_{\text{solve}}$  are the complexities of applying  $\mathbf{F}$  or  $\mathbf{F}^{-1}$  to a vector. The storage complexity is  $O(n \log n)$ .*

From (6) we see that the application of the factorization  $\mathbf{F}$  to a vector  $x \in \mathbb{R}^n$  simply requires application of the block unit-triangular matrices  $\mathbf{L}_\mathcal{I}$  and  $\mathbf{U}_\mathcal{I}$  corresponding to each subdomain at each level as well as the block-diagonal Cholesky factor  $\mathbf{C}$  of  $\tilde{\Sigma}_1$ . Further, the inverse of  $\mathbf{F}$  can be applied by noting that

$$\mathbf{F}^{-1} = \left[ \prod_{\ell=L}^1 \left( \prod_{\mathcal{I} \in \mathcal{L}_\ell} \mathbf{P}_\ell^T \mathbf{U}_\mathcal{I}^{-1} \mathbf{L}_\mathcal{I}^{-1} \right) \right]^T \mathbf{C}^{-T} \mathbf{C}^{-1} \left[ \prod_{\ell=L}^1 \left( \prod_{\mathcal{I} \in \mathcal{L}_\ell} \mathbf{P}_\ell^T \mathbf{U}_\mathcal{I}^{-1} \mathbf{L}_\mathcal{I}^{-1} \right) \right].$$

Table 1: Complexity of the 2D recursive skeletonization factorization

Operation	Complexity
Construct $\mathbf{F} \approx \Sigma$	$O(n^{3/2})$
Apply $\mathbf{F}$ or $\mathbf{F}^{-1}$ to a vector	$O(n \log n)$
Apply $\mathbf{F}^{1/2}$ or $\mathbf{F}^{-1/2}$ to a vector	$O(n \log n)$
Compute $\log  \mathbf{F} $ from $\mathbf{F}$	$O(n \log n)$

Additionally, a generalized square root  $\mathbf{F}^{1/2}$  such that  $\mathbf{F} = \mathbf{F}^{1/2}(\mathbf{F}^{1/2})^T$  can be applied (as can its transpose or inverse) by taking

$$\mathbf{F}^{1/2} = \left[ \prod_{\ell=1}^L \left( \prod_{\mathcal{I} \in \mathcal{L}_\ell} \mathbf{L}_{\mathcal{I}} \mathbf{U}_{\mathcal{I}} \right) \mathbf{P}_\ell \right] \mathbf{C}.$$

Finally, the log-determinant of  $\Sigma$  can be approximated by  $\log |\Sigma| \approx \log |\mathbf{F}| = 2 \log |\mathbf{C}|$ . Table 1 summarizes the computational complexities for these operations, which essentially follow from Theorem 1.

By constructing the recursive skeletonization factorizations  $\mathbf{F}$  of  $\Sigma$  and  $\mathbf{F}_i$  of  $\Sigma_i$  for  $i = 1, \dots, p$ , we see that after the  $O(n^{3/2})$  initial factorization cost each term in the evaluation of the log-likelihood or its gradient can be computed with cost  $O(n \log n)$  except for the product traces  $\text{Tr}(\Sigma^{-1}\Sigma_i)$  for  $i = 1, \dots, p$ . Further, through the approximate generalized square root  $\mathbf{F}^{1/2}$  of  $\Sigma$  we can quickly sample from the distribution  $N(0, \Sigma)$ .

**Remark 4.** *While the recursive skeletonization factorization described here exploits the most well-justified rank assumptions on the covariance kernel  $K(\cdot, \cdot; \theta)$ , each recursive skeletonization factorization in our framework can be replaced by the closely-related hierarchical interpolative factorization (HIF) [24] or strong recursive skeletonization factorization [32], which are observed in practice to exhibit better scaling properties and also admit simple log-determinant computation.*

### 3 Computing the trace terms

There are a number of methods for estimating the term  $\text{Tr}(\Sigma^{-1}\Sigma_i)$  appearing in each gradient component  $g_i$  for  $i = 1, \dots, p$ . Employing the recursive skeletonization factorizations  $\mathbf{F}$  of  $\Sigma$  and  $\mathbf{F}_i$  of  $\Sigma_i$ , the product  $\Sigma^{-1}\Sigma_i \approx \mathbf{F}^{-1}\mathbf{F}_i$  (or a symmetrized form with the same trace) can be applied to a vector with complexity  $O(n \log n)$ . Using  $\mathbf{G}$  to denote this black-box linear operator, the classical statistical approach is the estimator of Hutchinson [25]. Drawing random vectors  $u_i \in \{-1, 1\}^n$  for  $i = 1, \dots, q$  such that the components of  $u_i$  are independent and take value  $\pm 1$  with equal probability, the Hutchinson trace estimator is

$$\text{Tr}(\mathbf{G}) \approx \frac{1}{q} \sum_{i=1}^q u_i^T \mathbf{G} u_i, \quad (8)$$

which is unbiased with variance decaying as  $1/q$  and has cost  $O(qn \log n)$ . For low-accuracy estimates of the trace, the Hutchinson estimator is simple and computationally efficient, but for higher accuracy it proves computationally infeasible to use the Hutchinson approach because of the slow rate of convergence in  $q$ , see Section 5.

When  $\Sigma$  and  $\Sigma_i$  have hierarchical rank structure, it is reasonable to also look for hierarchical rank structure in the product  $\Sigma^{-1}\Sigma_i$ , as matrix inversion and multiplication preserve such rank structure (albeit with different ranks) in many cases [18]. In our framework, we use the matrix peeling algorithm of Lin et al. [28] for constructing an explicit  $\mathcal{H}$ -matrix representation of a fast black-box operator  $\mathbf{G}$ . At a high level, the method proceeds by applying the operator  $\mathbf{G}$  to random vectors drawn with a specific sparsity structure to construct an approximate representation of the off-diagonal blocks at each level. We recursively perform low-rank compression level-by-level, following the same quadtree hierarchy as in the recursive skeletonization factorization, albeit in a top-down traversal rather than bottom-up. Finally, at the bottom level of the tree, the diagonal blocks corresponding to leaf-level subdomains can be extracted and their traces computed. While the full algorithm is applicable to both the strongly-admissible and weakly-admissible setting, the version of the algorithm we detail here is efficient for the simple weakly-admissible case. We point the reader to Lin et al. [28] for more details related to the modifications required for strong admissibility.

The use of a randomized method for computing low-rank representations of matrices, which we review below, is integral to the peeling algorithm.

### 3.1 Randomized low-rank approximations

To begin, suppose that matrix  $\mathbf{A} \in \mathbb{R}^{m_1 \times m_2}$  has (numerical) rank  $r$  and that we wish to construct an explicit rank- $r$  approximation  $\mathbf{A} \approx \mathbf{U}_1 \mathbf{M} \mathbf{U}_2^T$  with  $\mathbf{U}_1 \in \mathbb{R}^{m_1 \times r}$ ,  $\mathbf{U}_2 \in \mathbb{R}^{m_2 \times r}$ , and  $\mathbf{M} \in \mathbb{R}^{r \times r}$ . In the context of approximating the trace terms, for example,  $\mathbf{A}$  will be an off-diagonal block of  $\Sigma^{-1}\Sigma_i$  or perhaps of a related symmetrized form with the same trace. Here we provide an overview of an algorithm that accomplishes this goal.

We begin by constructing approximations to the column space and row space of the matrix  $\mathbf{A}$ . Let  $c$  be a small integer and suppose  $\mathbf{W}_1 \in \mathbb{R}^{m_2 \times (r+c)}$  and  $\mathbf{W}_2 \in \mathbb{R}^{m_1 \times (r+c)}$  are appropriately chosen random matrices, the distribution of which we will discuss later. Following Halko et al. [22], let  $\mathbf{U}_1$  be a well-conditioned basis for the column space of  $\mathbf{A}\mathbf{W}_1$  and  $\mathbf{U}_2$  be a well-conditioned basis for the column space of  $\mathbf{A}^T\mathbf{W}_2$  constructed via, *e.g.*, column-pivoted QR factorizations. Using the Moore-Penrose pseudoinverse, we obtain a low-rank approximation according to the approach summarized by Lin et al. [28, subsection 1.2] via

$$\mathbf{A} \approx \mathbf{U}_1 \left[ (\mathbf{W}_2^T \mathbf{U}_1)^\dagger (\mathbf{W}_2^T \mathbf{A} \mathbf{W}_1) (\mathbf{U}_2^T \mathbf{W}_1)^\dagger \right] \mathbf{U}_2^T = \mathbf{U}_1 \mathbf{M} \mathbf{U}_2^T. \quad (9)$$

Perhaps surprisingly, with an appropriate choice of  $\mathbf{W}_1$  and  $\mathbf{W}_2$  it is the case that with high probability this approximation is *near-optimal*, in the sense that

$$\|\mathbf{A} - \mathbf{U}_1 \mathbf{M} \mathbf{U}_2^T\|_2 \leq \alpha(m_1, m_2, c) \|\mathbf{A} - \mathbf{A}_{r, \text{best}}\|_2,$$

where  $\mathbf{A}_{r,\text{best}}$  is the best rank- $r$  approximation of  $\mathbf{A}$  and  $\alpha(m_1, m_2, c)$  is a small factor dependent on  $c$  and the size of  $\mathbf{A}$ . Further, the approximation process can be monitored and controlled adaptively to ensure a target desired accuracy [22].

It remains to discuss the choice of distribution for  $\mathbf{W}_1$  and  $\mathbf{W}_2$ . The most common and straightforward choice is for both to have i.i.d.  $N(0, 1)$  entries, which guarantees the strongest analytical error bounds and highest success probability. Under this choice, one can show that the algorithm as stated takes  $O(T_{\text{apply}}r + nr^2 + r^3)$ , where  $T_{\text{apply}}$  is the complexity of applying  $\mathbf{A}$  to a vector. This is sufficiently fast for our purposes, though we note that it is possible to accelerate this using other distributions [34, 22, 41].

## 3.2 Matrix peeling for weakly-admissible matrices

For simplicity, we assume a perfect quadtree as in Remark 2. Further, we will assume that the numerical ranks of the off-diagonal blocks to a specified tolerance  $\epsilon_{\text{peel}}$  are known *a priori* at each level, such that off-diagonal blocks of  $\mathbf{G}$  at level  $\ell$  have numerical rank at most  $r_\ell$ . In practice, an adaptive procedure is used to find the ranks. Finally, we assume that  $\mathbf{G}$  is symmetric, since if the trace of nonsymmetric  $\mathbf{G}$  is required we can always instead consider a symmetrized form with the same trace such as  $\frac{1}{2}(\mathbf{G} + \mathbf{G}^T)$ .

### 3.2.1 First level of peeling algorithm

To begin, at level  $\ell = 1$  the domain is partitioned into four subdomains  $\Omega_{1;i}$  with corresponding index sets  $\mathcal{I}_{1;i}$ ,  $i = 1, \dots, 4$  as in Fig. 1. We follow the style of Lin et al. [28] and write the off-diagonal blocks at this level as  $\mathbf{G}_{1;ij} \equiv \mathbf{G}(\mathcal{I}_{1;i}, \mathcal{I}_{1;j})$  to make our notation less cumbersome.

To construct randomized low-rank approximations of  $\mathbf{G}_{1;ij}$  for  $i \neq j$ , we need to find the action of these off-diagonal blocks on random matrices as described in Section 3.1. Define the block-sparse matrices

$$\mathbf{W}_1^{(1)} \equiv \begin{bmatrix} \mathbf{W}_{1;1} \\ 0 \\ 0 \\ 0 \end{bmatrix}, \quad \mathbf{W}_2^{(1)} \equiv \begin{bmatrix} 0 \\ \mathbf{W}_{1;2} \\ 0 \\ 0 \end{bmatrix}, \quad \mathbf{W}_3^{(1)} \equiv \begin{bmatrix} 0 \\ 0 \\ \mathbf{W}_{1;3} \\ 0 \end{bmatrix}, \quad \mathbf{W}_4^{(1)} \equiv \begin{bmatrix} 0 \\ 0 \\ 0 \\ \mathbf{W}_{1;4} \end{bmatrix},$$

where  $\mathbf{W}_{1;j}$  is a random matrix of dimension  $|\mathcal{I}_{1;j}| \times (r_1 + c)$  for  $j = 1, \dots, 4$ . Applying  $\mathbf{G}$  to  $\mathbf{W}_1^{(1)}$  gives the action of  $\mathbf{G}_{1;1j}$  on the random matrix  $\mathbf{W}_{1;1}$  for  $j = 1, \dots, 4$ , since

$$\begin{bmatrix} \mathbf{G}_{1;11} & \mathbf{G}_{1;12} & \mathbf{G}_{1;13} & \mathbf{G}_{1;14} \\ \mathbf{G}_{1;21} & \mathbf{G}_{1;22} & \mathbf{G}_{1;23} & \mathbf{G}_{1;24} \\ \mathbf{G}_{1;31} & \mathbf{G}_{1;32} & \mathbf{G}_{1;33} & \mathbf{G}_{1;34} \\ \mathbf{G}_{1;41} & \mathbf{G}_{1;42} & \mathbf{G}_{1;43} & \mathbf{G}_{1;44} \end{bmatrix} \begin{bmatrix} \mathbf{W}_{1;1} \\ \\ \\ \end{bmatrix} = \begin{bmatrix} \mathbf{G}_{1;11} \mathbf{W}_{1;1} \\ \mathbf{G}_{1;21} \mathbf{W}_{1;1} \\ \mathbf{G}_{1;31} \mathbf{W}_{1;1} \\ \mathbf{G}_{1;41} \mathbf{W}_{1;1} \end{bmatrix}.$$

The top block of the right-hand side vector above is unused as it involves a diagonal block of  $\mathbf{G}$ . However, the remaining blocks are exactly the matrices  $\mathbf{G}_{1;i}\mathbf{W}_{1;1}$  for  $i \neq 1$  as required by the randomized low-rank approximation of Section 3.1. Applying  $\mathbf{G}$  to each  $\mathbf{W}_j^{(1)}$  for  $j = 1, \dots, 4$ , for each block  $\mathbf{G}_{1;j}$  with  $i \neq j$  we obtain a random sampling of its column space  $\mathbf{G}_{1;j}\mathbf{W}_{1;j}$ . Note that by symmetry of  $\mathbf{G}$  we also obtain a random sampling of the row space of each block since  $\mathbf{G}_{1;j}\mathbf{W}_{1;j} = \mathbf{G}_{1;j}^T\mathbf{W}_{1;j}$ .

Using (9) to construct rank- $r_1$  approximations of each of these blocks, we write the approximation of  $\mathbf{G}_{1;j}$  as  $\mathbf{G}_{1;j} \approx \widehat{\mathbf{G}}_{1;j} \equiv \mathbf{U}_{1;j}\mathbf{M}_{1;j}\mathbf{U}_{1;j}^T$ , where the approximation is accurate to the specified tolerance  $\epsilon_{\text{peel}}$  with high probability.

Defining the matrix  $\mathbf{G}^{(1)} \in \mathbb{R}^{n \times n}$  with blocks given by

$$\mathbf{G}^{(1)}(\mathcal{I}_{1;i}, \mathcal{I}_{1;j}) = \begin{cases} \widehat{\mathbf{G}}_{1;j} & i \neq j, \\ 0 & \text{else,} \end{cases}$$

we obtain

$$\mathbf{G} - \mathbf{G}^{(1)} \equiv \mathbf{G} - \begin{bmatrix} & \widehat{\mathbf{G}}_{1;12} & \widehat{\mathbf{G}}_{1;13} & \widehat{\mathbf{G}}_{1;14} \\ \widehat{\mathbf{G}}_{1;21} & & \widehat{\mathbf{G}}_{1;23} & \widehat{\mathbf{G}}_{1;24} \\ \widehat{\mathbf{G}}_{1;31} & \widehat{\mathbf{G}}_{1;32} & & \widehat{\mathbf{G}}_{1;34} \\ \widehat{\mathbf{G}}_{1;41} & \widehat{\mathbf{G}}_{1;42} & \widehat{\mathbf{G}}_{1;43} & \end{bmatrix} \approx \begin{bmatrix} \mathbf{G}_{1;11} & & & \\ & \mathbf{G}_{1;22} & & \\ & & \mathbf{G}_{1;33} & \\ & & & \mathbf{G}_{1;44} \end{bmatrix}.$$

In other words, we have approximated the off-diagonal blocks at this level to a specified accuracy and used the result to obtain a fast operator  $\mathbf{G} - \mathbf{G}^{(1)}$  that is block-diagonal with diagonal blocks the same as those of  $\mathbf{G}$ .

**Remark 5.** *We note that the matrix  $\mathbf{G}^{(1)}$  is not explicitly assembled as a dense matrix inside the peeling algorithm. Instead, we store the non-zero blocks in low-rank form so that  $\mathbf{G}^{(1)}$  may be efficiently applied to vectors.*

### 3.2.2 Second level of peeling algorithm

In the next step of the peeling algorithm, we recurse on the diagonal subblocks  $\mathbf{G}_{1;ii}$  for  $i = 1, \dots, 4$ . Partitioning each subdomain  $\Omega_{1;i}$  at level  $\ell = 1$  into four child subdomains at level  $\ell = 2$  using the quadtree structure and renumbering blocks accordingly, we write the diagonal blocks  $\mathbf{G}_{1;ii}$  for  $i = 1, \dots, 4$  as

$$\mathbf{G}_{1;11} = \begin{bmatrix} \mathbf{G}_{2;11} & \mathbf{G}_{2;12} & \mathbf{G}_{2;13} & \mathbf{G}_{2;14} \\ \mathbf{G}_{2;21} & \mathbf{G}_{2;22} & \mathbf{G}_{2;23} & \mathbf{G}_{2;24} \\ \mathbf{G}_{2;31} & \mathbf{G}_{2;32} & \mathbf{G}_{2;33} & \mathbf{G}_{2;34} \\ \mathbf{G}_{2;41} & \mathbf{G}_{2;42} & \mathbf{G}_{2;43} & \mathbf{G}_{2;44} \end{bmatrix}, \quad \mathbf{G}_{1;22} = \begin{bmatrix} \mathbf{G}_{2;55} & \mathbf{G}_{2;56} & \mathbf{G}_{2;57} & \mathbf{G}_{2;58} \\ \mathbf{G}_{2;65} & \mathbf{G}_{2;66} & \mathbf{G}_{2;67} & \mathbf{G}_{2;68} \\ \mathbf{G}_{2;75} & \mathbf{G}_{2;76} & \mathbf{G}_{2;77} & \mathbf{G}_{2;78} \\ \mathbf{G}_{2;85} & \mathbf{G}_{2;86} & \mathbf{G}_{2;87} & \mathbf{G}_{2;88} \end{bmatrix},$$

and so on for  $\mathbf{G}_{1;33}$  and  $\mathbf{G}_{1;44}$ .

For each  $j = 1, \dots, 16$  we define the random matrix  $W_{2;j} \in \mathbb{R}^{|\mathcal{I}_{2;j}| \times (r_2+c)}$ , which is appropriately sized to give a random sample of the column space of  $G_{2;ij}$  for each  $i = 1, \dots, 16$ ,  $i \neq j$ . We can minimize the number of times we apply the operator  $G - G^{(1)}$  as follows due to its block diagonal structure. For each  $k = 1, \dots, 4$ , we define  $W_k^{(2)} \in \mathbb{R}^{n \times (r_2+c)}$  to have rows divided into 16 blocks according to

$$W_k^{(2)}(\mathcal{I}_{2;j}, \cdot) = \begin{cases} W_{2;j} & j \in \{k, k+4, k+8, k+12\}, \\ 0 & \text{else.} \end{cases}$$

In other words, block  $k$  of  $W_k^{(2)}$  is nonzero, as is every fourth block after  $k$ .

**Definition 3** (Quadtree siblings). *In the context of the quadtree decomposition of  $\Omega$ , we say that  $\Omega_{\ell;i}$  and  $\Omega_{\ell;j}$  are siblings if  $i \neq j$  and both  $\Omega_{\ell;i} \subset \Omega_{\ell-1;k}$  and  $\Omega_{\ell;j} \subset \Omega_{\ell-1;k}$  for some  $1 \leq k \leq 4^{(\ell-1)}$ .*

Let  $B_k \equiv (G - G^{(1)})W_k^{(2)}$  and suppose that  $W_k^{(2)}(\mathcal{I}_{2;j}, \cdot)$  is nonzero. For each  $i$  such that  $\Omega_{2;i}$  and  $\Omega_{2;j}$  are siblings, we have  $B_k(\mathcal{I}_{2;i}, \cdot) \approx G_{2;ij}W_{2;j}$ . For example, in  $W_1^{(2)}$  the nonzero blocks are  $W_1^{(2)}(\mathcal{I}_{2;j}, \cdot) = W_{2;j}$  for  $j \in \{1, 5, 9, 13\}$ , so

$$B_1(\mathcal{I}_{2;i}, \cdot) \approx \begin{cases} G_{2;i1}W_{2;1} & i = 1, \dots, 4, \\ G_{2;i5}W_{2;5} & i = 5, \dots, 8, \\ G_{2;i9}W_{2;9} & i = 9, \dots, 12, \\ G_{2;i13}W_{2;13} & i = 13, \dots, 16. \end{cases}$$

Therefore, applying  $G - G^{(1)}$  to  $W_k^{(2)}$  for  $k = 1, \dots, 4$  gives a random sample of the column space and row space of  $G_{2;ij}$  for each  $i$  and  $j$  such that  $\Omega_{2;i}$  and  $\Omega_{2;j}$  are siblings. For all such  $i$  and  $j$  we use the randomized low-rank approximation algorithm as before to construct

$$\widehat{G}_{2;ij} = U_{2;ij}M_{2;ij}U_{2;ji}^T.$$

Defining  $G^{(2)} \in \mathbb{R}^{n \times n}$  with blocks

$$G^{(2)}(\mathcal{I}_{2;i}, \mathcal{I}_{2;j}) = \begin{cases} \widehat{G}_{2;ij} & \text{if } \Omega_{2;i} \text{ and } \Omega_{2;j} \text{ are siblings,} \\ 0 & \text{else,} \end{cases}$$

we have that  $G - G^{(1)} - G^{(2)}$  is approximately block-diagonal with diagonal blocks  $G_{2;ii}$  for  $i = 1, \dots, 16$ .

### 3.2.3 Subsequent levels of peeling algorithm

In general at level  $\ell > 2$  we see that  $G - \sum_{m=1}^{\ell-1} G^{(m)}$  is approximately block-diagonal with  $4^{(\ell-1)}$  diagonal blocks. For each  $k = 1, \dots, 4$  we define  $W_k^{(\ell)} \in \mathbb{R}^{n \times (r_\ell+c)}$  to have rows divided

into  $4^\ell$  blocks according to

$$\mathbf{W}_k^{(\ell)}(\mathcal{I}_{\ell;j}, \cdot) = \begin{cases} \mathbf{W}_{\ell;j} & j \equiv k \pmod{4}, \\ 0 & \text{else,} \end{cases}$$

where each  $\mathbf{W}_{\ell;j}$  is a random matrix of size  $\mathbb{R}^{|\mathcal{I}_{\ell;j}| \times r_\ell}$ . Using the same logic as in Section 3.2.2, we apply  $\mathbf{G} - \sum_{m=1}^{\ell-1} \mathbf{G}^{(m)}$  to  $\mathbf{W}_k^{(\ell)}$  for each  $k = 1, \dots, 4$  and use the results to construct low rank approximations

$$\widehat{\mathbf{G}}_{\ell;ij} = \mathbf{U}_{\ell;ij} \mathbf{M}_{\ell;ij} \mathbf{U}_{\ell;ji}^T.$$

for each  $i$  and  $j$  such that  $\Omega_{\ell;i}$  and  $\Omega_{\ell;j}$  are siblings. We define

$$\mathbf{G}^{(\ell)}(\mathcal{I}_{\ell;i}, \mathcal{I}_{\ell;j}) = \begin{cases} \widehat{\mathbf{G}}_{\ell;ij} & \text{if } \Omega_{\ell;i} \text{ and } \Omega_{\ell;j} \text{ are siblings,} \\ 0 & \text{else} \end{cases}$$

such that  $\mathbf{G} - \sum_{m=1}^{\ell} \mathbf{G}^{(m)}$  is approximately block-diagonal with  $4^\ell$  diagonal blocks.

### 3.2.4 Extracting the trace

At the bottom level of the quadtree, each diagonal block of  $\mathbf{G} - \sum_{m=1}^L \mathbf{G}^{(m)}$  is of a constant size independent of  $n$  as discussed in Section 2.1. Define  $n_{L;i} \equiv |\mathcal{I}_{L;i}|$  and  $n_L \equiv \max_i n_{L;i}$  such that  $n_L$  is the maximum number of observations in a leaf-level subdomain. We construct a block matrix  $\mathbf{E} \in \mathbb{R}^{n \times n_L}$  such that

$$\mathbf{E}(\mathcal{I}_{L;i}, [n_{L;i}]) = \mathbf{I} \in \mathbb{R}^{|\mathcal{I}_{L;i}| \times |\mathcal{I}_{L;i}|}$$

for each  $i$ , where  $\mathbf{I}$  is an appropriately-sized identity matrix. Letting

$$\mathbf{H} \equiv \left( \mathbf{G} - \sum_{m=1}^L \mathbf{G}^{(m)} \right) \mathbf{E},$$

we find that  $\mathbf{H}(\mathcal{I}_{L;i}, [n_{L;i}]) \approx \mathbf{G}_{L;ii}$  for each  $i = 1, \dots, 4^L$ . We can then approximate the trace of  $\mathbf{G}$  using the relation

$$\text{Tr}(\mathbf{G}) = \sum_{i=1}^{4^L} \text{Tr}(\mathbf{G}_{L;ii}) \approx \sum_{i=1}^{4^L} \text{Tr}(\mathbf{H}(\mathcal{I}_{L;i}, [n_{L;i}])).$$

**Remark 6.** *When using the peeling algorithm to construct an approximate trace of an operator  $\mathbf{G}$  with numerically low-rank off-diagonal blocks, it is important to note that we do not have direct control of the relative error of the trace approximation. This is because a matrix with diagonal entries with large absolute value but mixed signs can have a small trace due to cancellation. In practice, however, our numerical results in Section 5 show excellent agreement between the approximate trace and true trace.*

### 3.3 Computational complexity

For each level of the weak-admissibility-based peeling algorithm described in Section 3.2 there are two key steps: applying the operator  $\mathbf{G} - \sum_{m=1}^{\ell-1} \mathbf{G}^{(m)}$  and forming the low-rank factorizations  $\widehat{\mathbf{G}}_{\ell;ij}$  for each  $i$  and  $j$  such that  $\Omega_{\ell;i}$  and  $\Omega_{\ell;j}$  are siblings. Analyzing the cost of these steps leads to the following complexity result.

**Theorem 2.** *Let the cost of applying  $\mathbf{G} \in \mathbb{R}^{n \times n}$  to a vector be  $T_{\text{apply}}$  and assume that the observations are uniformly distributed in  $\Omega$  such that  $|\mathcal{I}_{\ell;i}| = O(4^{-\ell}n)$  for each  $1 \leq \ell \leq L$  and  $1 \leq i \leq 4^\ell$ . Assuming that the ranks of the off-diagonal blocks  $\mathbf{G}_{\ell;i,j}$  are bounded by  $r_\ell$  for each  $i$  and  $j$  such that  $\Omega_{\ell;i}$  and  $\Omega_{\ell;j}$  are siblings, and define*

$$s_1 \equiv \sum_{\ell=1}^L r_\ell, \quad \text{and} \quad s_2 \equiv \sum_{\ell=1}^L r_\ell^2.$$

Then the complexity of the weak-admissibility-based peeling algorithm is

$$T_{\text{peel}} = O(T_{\text{apply}}s_1 + ns_2 \log n). \quad (10)$$

The storage complexity is  $O(ns_1)$ .

*Proof.* We adapt the proof of Lin et al. [28] to the weak admissibility case. At the first level, applying  $\mathbf{G}$  to each  $\mathbf{W}_k^{(1)}$  costs  $O(T_{\text{apply}}r_1)$  and each randomized factorization costs  $O(nr_1^2)$ , leading to an overall cost for level 1 of  $O(T_{\text{apply}}r_1 + nr_1^2)$ .

At level  $\ell > 1$ , we break the cost of applying  $\mathbf{G} - \sum_{m=1}^{\ell-1} \mathbf{G}^{(m)}$  into two pieces. The cost of applying  $\mathbf{G}$  to each  $\mathbf{W}_k^{(\ell)}$  is  $O(T_{\text{apply}}r_\ell)$ . The matrix  $\sum_{m=1}^{\ell-1} \mathbf{G}^{(m)}$  is a heavily structured matrix with blocks in low-rank form. Applying this to each  $\mathbf{W}_k^{(\ell)}$  costs  $O\left(\sum_{m=1}^{\ell-1} nr_m r_\ell\right)$ , which is  $O(ns_2)$ . We additionally must construct each randomized factorization at this level. Each one costs  $O(4^{-\ell}nr_\ell^2)$  and there are  $O(4^\ell)$  off-diagonal blocks to compress at this level, so the overall cost for level  $\ell$  is  $O(T_{\text{apply}}r_\ell + ns_2 + nr_\ell^2)$ .

Summing the cost of each level from  $\ell = 1, \dots, L$ , we obtain (10). Note that at level  $\ell = L$ , we must additionally extract the diagonal blocks, but by the assumption these blocks are of constant size so this does not increase the asymptotic cost. The storage complexity comes from noting that at level  $\ell$  we must store the  $O(4^\ell)$  matrices of rank  $r_\ell$ , where each has outer dimension that is  $O(4^{-\ell}n)$ .  $\square$

When the underlying matrix has the rank of all off-diagonal blocks bounded by  $r_\ell = O(1)$  for all  $\ell$ , then the computational complexity of weak peeling is  $\tilde{O}(T_{\text{apply}} + n)$ , where we use the so-called ‘‘soft-O’’ notation from theoretical computer science to suppress factors that are polylogarithmic in  $n$ . In this case, peeling  $\mathbf{G} = \Sigma$  itself using its recursive skeletonization factorization results in  $\tilde{O}(n)$  complexity for both time and memory.

Many real matrices of interest, however, do not exhibit off-diagonal blocks with ranks independent of  $n$ . For example, our experiments with the Matérn kernel of (3) show that a constant number of off-diagonal blocks at each level of the hierarchy exhibit ranks bounded

Table 2: The runtime and storage complexity of the peeling algorithm depend on the asymptotic rank  $r_\ell$  of off-diagonal blocks of  $\mathbf{G}$  at level  $\ell$ . The tabulated complexities are based on the assumption that the recursive skeletonization factorization is used to apply  $\mathbf{G}$  as a fast operator.

$r_\ell$	Time	Storage
$O(\log n_\ell)$	$\tilde{O}(n)$	$\tilde{O}(n)$
$O(\sqrt{n_\ell})$	$\tilde{O}(n^2)$	$\tilde{O}(n^{3/2})$

only as  $r_\ell = O(2^{-\ell}\sqrt{n})$ . This coincides with the argument for rank growth in the recursive skeletonization factorization in Section 2.3. Thus, this simplified peeling algorithm in the case of the Matérn kernel has asymptotic time complexity  $\tilde{O}(n^2)$  and storage complexity  $\tilde{O}(n^{3/2})$ , where we pick up at most a polylogarithmic factor in the ranks since we are looking at  $\mathbf{G} = \Sigma^{-1}\Sigma_i$  and not  $\Sigma$  itself. In theory, using the simple peeling algorithm described here is asymptotically no better than extracting the trace by applying  $\mathbf{G}$  to the coordinate vectors  $e_i$  for  $i = 1, \dots, n$ . This necessitates the standard form of peeling for large problems.

**Remark 7.** *In practice, the standard form of the peeling algorithm [28] that uses the full generality of strong admissibility can be employed to remedy such rank growth by explicitly avoiding compression of off-diagonal blocks that are not sufficiently low-rank. Using the modifications described in that paper, the complexity of peeling follows the same bound as Theorem 2 but with the rank bound  $r_\ell$  referring to a bound on the ranks of only those blocks that are compressed in the strongly-admissible hierarchical format. We find in Section 5.1 that using peeling based on strong admissibility is more efficient when  $n$  is large, as expected. However, the implicit constants in the asymptotic runtime lead to weak admissibility being more efficient for moderately-sized problems.*

We summarize our complexity results in Table 2. Note that these results were derived on the assumption that  $n_\ell = O(4^{-\ell}n)$ , *i.e.*, a quasi-uniform distribution of observations and a perfect quadtree decomposition of space. In practice observations that are distributed in a different fashion can actually exhibit better behavior, particularly if the observations are concentrated around a quasi-one-dimensional curve [24].

## 4 Summary of MLE framework

In Algorithm 1 we summarize our complete framework for computing the log-likelihood  $\ell(\theta)$  and gradient  $g(\theta)$  given  $\theta$ , which can be used inside of any first-order optimization routine for Gaussian process maximum likelihood estimation. As mentioned previously, the approach is flexible and does not rely on the specific hierarchical factorization used (*e.g.*, the recursive skeletonization factorization, the hierarchical interpolative factorization, the strong recursive skeletonization factorization) or the form of peeling used (*i.e.*, the peeling based on weak admissibility described in Section 3 or the form by Lin et al. [28] based on

strong admissibility). Rather, the exact components of the framework should be decided on a case-by-case basis depending on the rank properties of the kernel family.

---

**Algorithm 1** Computing the Gaussian process log-likelihood and gradient

---

Given: observation vector  $z \in \mathbb{R}^n$ , observation locations  $\{x_i\}_{i=1}^n \subset \mathbb{R}^2$ , peel tolerance  $\epsilon_{\text{peel}}$ , factorization tolerance  $\epsilon_{\text{fact}} < \epsilon_{\text{peel}}$ , parameter vector  $\theta \in \mathbb{R}^p$ , and covariance kernel  $K(\cdot, \cdot; \theta)$

- 1: // Factor  $\Sigma$  with hierarchical factorization
- 2:  $F \leftarrow$  Recursive skeletonization factorization of  $\Sigma$  with tolerance  $\epsilon_{\text{fact}}$
- 3: // Use fast hierarchical solve and log-determinant
- 4:  $\hat{\ell}(\theta) \leftarrow -\frac{1}{2}z^T F^{-1}z - \frac{1}{2} \log |F| - \frac{1}{2} \log 2\pi \approx \ell(\theta)$
- 5: **for**  $i = 1, \dots, p$  **do**
- 6:   // Factor  $\Sigma_i$  with hierarchical factorization
- 7:    $F_i \leftarrow$  Recursive skeletonization factorization of  $\Sigma_i$  with tolerance  $\epsilon_{\text{fact}}$
- 8:   // Compute trace of  $\Sigma^{-1}\Sigma_i$  with peeling algorithm
- 9:    $t_i \leftarrow$  Trace of operator  $\frac{1}{2}(F^{-1}F_i + F_iF^{-1})$  via peeling algorithm with tolerance  $\epsilon_{\text{peel}}$
- 10:   // Use fast hierarchical apply and solve
- 11:    $\hat{g}_i \leftarrow \frac{1}{2}z^T F^{-1}F_iF^{-1}z - \frac{1}{2}t_i \approx g_i$
- 12: **end for**

Output:  $\hat{\ell}(\theta)$  and  $\hat{g}(\theta)$

---

**Remark 8.** After estimation of the parameter vector  $\theta$ , there remains the question of how to sample from the Gaussian process conditioned on the observed data  $z$ . Assuming  $[z', z]^T$  is jointly distributed according to the original Gaussian process, this conditional distribution is given by

$$z'|z \sim N(\Sigma_{12}\Sigma_{22}^{-1}z, \Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{12}^T),$$

where  $\Sigma_{11}$  is the covariance matrix of  $z'$ ,  $\Sigma_{22}$  is the covariance matrix of  $z$ , and so on. Using the identity

$$\Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{12}^T = \begin{bmatrix} I & -\Sigma_{12}\Sigma_{22}^{-1} \end{bmatrix} \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{12}^T & \Sigma_{22} \end{bmatrix} \begin{bmatrix} I \\ -\Sigma_{22}^{-1}\Sigma_{12}^T \end{bmatrix}$$

and letting  $\Sigma$  denote the two-by-two block matrix in a slight abuse of notation, we can apply a square-root of  $\Sigma_{11} - \Sigma_{12}\Sigma_{22}^{-1}\Sigma_{12}^T$  with skeletonization factorizations  $F_{22} \approx \Sigma_{22}$  and  $F \approx \Sigma$  by using  $F^{1/2}$  to apply a square-root of  $\Sigma$ ,  $F$  to apply  $\Sigma_{12}$  through appropriate padding, and  $F_{22}^{-1}$  to apply  $\Sigma_{22}^{-1}$ . This gives a fast method for sampling from the conditional distribution or computing the conditional mean.

## 5 Numerical results

To demonstrate the effectiveness of our approach to Gaussian process maximum likelihood estimation, we first test the accuracy and runtime of the peeling-based technique for approximating the trace and then test our full method on two examples using synthetic datasets and

one example using a dataset of measurements of ocean surface temperatures. For examples, we take the number of proxy points to be  $n_{\text{prox}} = 256$ , and use a quadtree decomposition of space with a maximum of  $n_{\text{occ}} = 64$  points per leaf subdomain.

In our tests we use the FLAM library (<https://github.com/klho/FLAM/>) for the recursive skeletonization factorization and a custom implementation of matrix peeling as described in Section 3.2. This additional code is available at <https://github.com/asdamle/GPMLE/>. All numerical results shown were run in MATLAB<sup>®</sup> R2015a on a quad-socket Intel<sup>®</sup> Xeon<sup>®</sup> E5-4640 processor clocked at 2.4 GHz using up to 1.5 TB of RAM.

## 5.1 Runtime scaling of the peeling algorithm

To begin, we investigate the numerical performance of the peeling algorithm on synthetic examples. We take the observation locations  $\{x_i\}_{i=1}^n$  to be a  $\sqrt{n} \times \sqrt{n}$  grid of points uniformly discretizing the square  $[0, 100]^2 \subset \mathbb{R}^2$ . We let  $\theta = [\theta_1, \theta_2]$  parameterize the correlation length scale of the process in each coordinate direction, defining the scaled distance

$$\|x - y\|_{\theta}^2 = \frac{(x_1 - y_1)^2}{\theta_1^2} + \frac{(x_2 - y_2)^2}{\theta_2^2},$$

where here  $x_i$  and  $y_i$  are used to denote components of vectors  $x$  and  $y$ . Using this parameterization and incorporating an additive noise term, the two kernels we test are the rational quadratic kernel of (2) with  $\alpha = 1/2$ ,

$$K_{RQ}(x, y; \theta) = (1 + \|x - y\|_{\theta}^2)^{-1/2} + \sigma_N^2 \delta_{xy}, \quad (11)$$

and the Matérn kernel of (3) with parameter  $\nu = 3/2$ ,

$$K_M(x, y; \theta) = (1 + \sqrt{3}\|x - y\|_{\theta}) \exp(-\sqrt{3}\|x - y\|_{\theta}) + \sigma_N^2 \delta_{xy}. \quad (12)$$

Here  $\delta_{xy}$  is the Kronecker delta, which satisfies  $\delta_{xy} = 1$  if  $x = y$  and  $\delta_{xy} = 0$  otherwise.

**Remark 9.** *In both (11) and (12) the additional term  $\sigma_N^2 \delta_{xy}$  can be interpreted as modeling additive white noise with variance  $\sigma_N^2$  on top of the base Gaussian process model. In practice, this so-called “nugget effect” is frequently incorporated to account for measurement error or small-scale variation from other sources [31] and, further, is numerically necessary for many choices of parameter  $\theta$  due to exceedingly poor conditioning of many kernel matrices.*

We compute high-accuracy recursive skeletonization factorizations of the matrices  $\Sigma$  and  $\Sigma_1 \equiv \frac{\partial}{\partial \theta_1} \Sigma$ , which we combine to obtain the fast black-box operator

$$\mathbf{G} = \frac{1}{2}(\Sigma^{-1}\Sigma_1 + \Sigma_1\Sigma^{-1}) \quad (13)$$

for input to the peeling algorithm to compute the trace to specified tolerance  $\epsilon_{\text{peel}} = 1 \times 10^{-6}$ . We choose the parameter vector  $\theta = [10, 7]$  for these examples as in Fig. 4 (left), and set the noise parameter at  $\sigma_N^2 = 1 \times 10^{-4}$ .

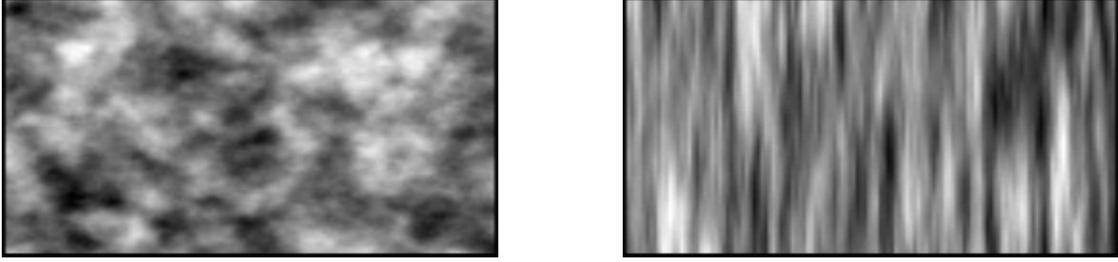


Figure 4: Two different realizations on the domain  $[0, 200] \times [0, 100]$  of the Matérn kernel Gaussian process with covariance seen in (12) and noise parameter  $\sigma_N^2 = 0$ . In the left figure the parameter vector is  $\theta = [10, 7]$  corresponding to a kernel that is relatively close to isotropic. In contrast, in the right figure the parameter vector  $\theta = [3, 30]$  generates strong anisotropy.

Table 3: Runtime  $t_{\text{peel}}$  of the the peeling algorithm with the rational quadratic kernel of (11). Note that we omit the relative error  $e_{\text{peel}}$  in the estimated trace for our largest example, as the operator was too large to determine the true trace using the naïve approach.

$n$	$t_{\text{peel,weak}}$ (s)	$e_{\text{peel,weak}}$	$t_{\text{peel,strong}}$ (s)	$e_{\text{peel,strong}}$
$64^2$	$9.17 \times 10^0$	$5.68 \times 10^{-7}$	$4.85 \times 10^1$	$1.60 \times 10^{-7}$
$128^2$	$9.16 \times 10^1$	$1.02 \times 10^{-5}$	$5.64 \times 10^2$	$2.58 \times 10^{-7}$
$256^2$	$6.63 \times 10^2$	$3.72 \times 10^{-5}$	$3.04 \times 10^3$	$3.32 \times 10^{-6}$
$512^2$	$2.88 \times 10^3$	-	$1.73 \times 10^4$	-

Beginning with the rational quadratic kernel, in Table 3 we give runtime results for both the simplified peeling algorithm described in Section 3.2 (“weak peeling”) as well as the full strong-admissibility-based peeling algorithm of Lin et al. [28] (“strong peeling”). As can be seen in Fig. 5 (left), the runtime of the peeling algorithm with the kernel (11) seems to scale between  $O(n)$  and  $O(n^{3/2})$  with the number of observations  $n$ , regardless of whether weak or strong peeling is used. Further, the relative error in the trace approximation,  $e_{\text{peel}}$  is near the specified tolerance  $\epsilon_{\text{peel}}$ , though the tolerance is not a hard upper bound. Note that we omit the relative error for our largest example, as the operator was too large to determine the true trace using the naïve approach.

In contrast, the results in Table 4 for the Matérn kernel in (12) show different scaling behavior for weak and strong peeling. In Fig. 5 (right), we see that the runtime for weak peeling seems to be close to quadratic in the number of observations, which agrees with our analysis from Section 3. Using strong peeling, however, the complexity of peeling scales considerably better, ultimately following the  $O(n^{3/2})$  trend line. We see again that the relative trace error is well-controlled by  $\epsilon_{\text{peel}}$  in both cases.

Though the observed scaling behavior of strong peeling is as good or better than that for weak peeling for both kernels, in practice we see that for problems with up to a quarter of a million observations weak peeling has a smaller time-to-solution. As such, in the remainder

Table 4: Runtime  $t_{\text{peel}}$  of the the peeling algorithm with the Matérn kernel of (12). Note that we omit the relative error  $e_{\text{peel}}$  in the estimated trace for our largest example, as the operator was too large to determine the true trace using the naïve approach.

$n$	$t_{\text{peel,weak}}$ (s)	$e_{\text{peel,weak}}$	$t_{\text{peel,strong}}$ (s)	$e_{\text{peel,strong}}$
$64^2$	$6.03 \times 10^0$	$5.73 \times 10^{-8}$	$2.06 \times 10^1$	$4.78 \times 10^{-10}$
$128^2$	$5.30 \times 10^1$	$2.46 \times 10^{-7}$	$2.29 \times 10^2$	$3.36 \times 10^{-10}$
$256^2$	$5.37 \times 10^2$	$4.28 \times 10^{-6}$	$1.62 \times 10^3$	$8.14 \times 10^{-10}$
$512^2$	$7.07 \times 10^3$	-	$1.00 \times 10^4$	-

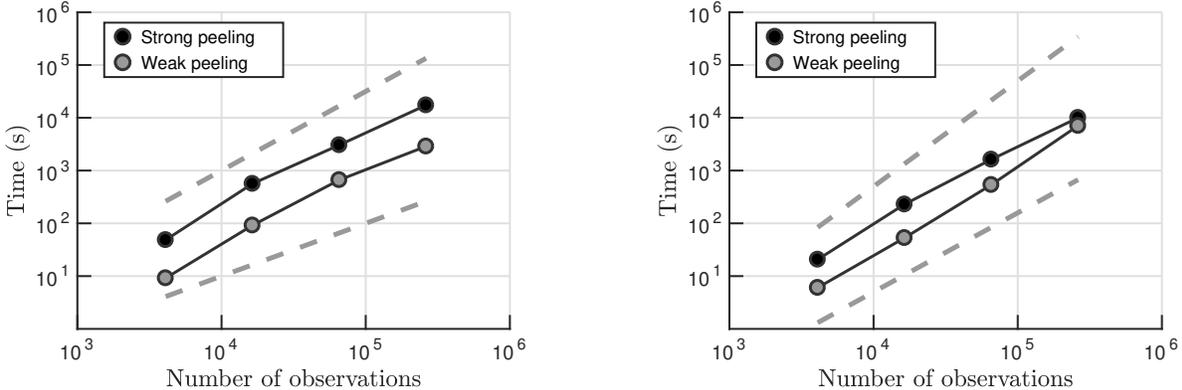


Figure 5: On the left the runtime of peeling for the rational quadratic kernel is plotted along with a  $O(n^{3/2})$  trend line (top) and a  $O(n)$  trend line (bottom), showing subquadratic scaling for weak peeling in this case. In contrast, on the right the runtime of peeling for the Matérn kernel is plotted along with a  $O(n^2)$  trend line (top) and a  $O(n^{3/2})$  trend line (bottom). We see that weak peeling with the Matérn kernel seems to ultimately exhibit quadratic scaling, whereas strong peeling seems to exhibit slightly better than  $O(n^{3/2})$  scaling. The corresponding data are given in Tables 3 and 4.

of our examples we show results using only weak peeling.

## 5.2 Relative efficiency of peeling versus the Hutchinson estimator

As discussed in Section 3, a common alternative statistical approach for approximating the trace of a matrix  $\mathbf{G}$  is the estimator of Hutchinson [25] seen in (8). The aim of this section is to show that for matrices with hierarchical low-rank structure our peeling-based algorithm can be much more efficient when a high-accuracy trace approximation is desired.

As in Section 5.1, we take our observations to be a regular grid discretizing  $[0, 100]^2 \subset \mathbb{R}^2$  using the Matérn kernel of (12) with noise  $\sigma_N^2 = 1 \times 10^{-4}$  and parameter vector  $\theta = [10, 7]$ . We fix the number of observations at  $n = 64^2$  and consider how the accuracy of the trace approximation varies with the number of applications of the black-box operator for both

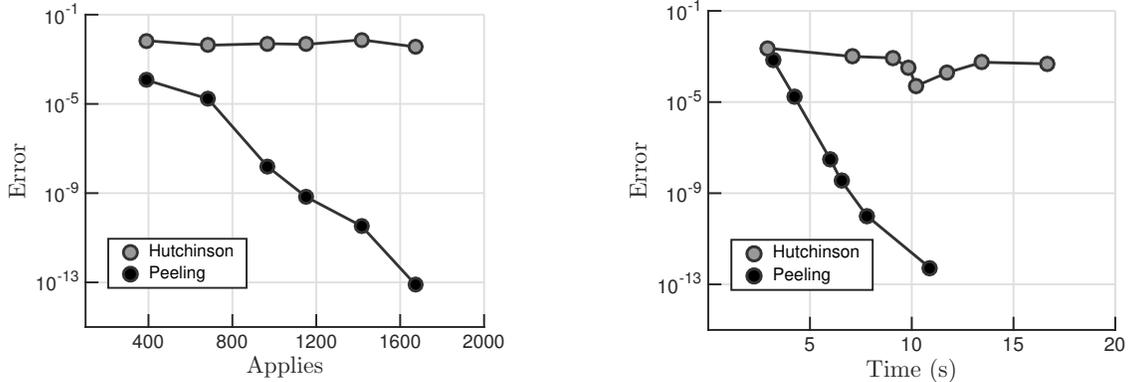


Figure 6: Plotting the relative error in the trace approximation versus the number of applications of the black-box operator, we see in the left figure that the Hutchinson estimator exhibits characteristic inverse square root convergence as dictated by the central limit theorem. In contrast, using the peeling algorithm described in Section 3.2, we see that the same number of black-box applies yields a much improved accuracy, though the rate of convergence depends on the spectra of off-diagonal blocks of the operator. In the right figure, we plot the error of each method versus wall-clock time to establish that the same scaling behavior holds when error is viewed as a function of time-to-solution.

weak peeling and the Hutchinson estimator.

Using a high-accuracy recursive skeletonization factorization to construct the black-box operator in (13) as in Section 5.1, we vary the tolerance  $\epsilon_{\text{peel}}$  in the peeling algorithm and plot in Fig. 6 the relative error in the trace approximation as a function of both the number of black-box applies and total peeling runtime. Additionally, for the Hutchinson estimator we use the same factorizations to construct the unsymmetric operator  $\mathbf{G}' = \Sigma^{-1}\Sigma_1$ . We plot the same quantities for a given instantiation of the estimator for comparison.

For low-accuracy approximations with relative error on the order of  $1 \times 10^{-1}$  to  $1 \times 10^{-3}$ , we see that the Hutchinson estimator is a competitive alternative to the peeling algorithm for finding the trace. When increased accuracy is desired, however, it is clear that in our examples that the peeling algorithm is the more attractive option. While the Hutchinson estimator has a simple form and is easy to compute, the relatively slow inverse square root convergence means that  $M$  in (8) must be taken to be exceedingly large to drive the variance down to reasonable levels, whereas the peeling algorithm is observed to make more economical use of its black-box matrix-vector products. It is worth noting that, for this choice of  $n$ , only 4096 applies are needed to explicitly construct all diagonal entries of the operator via application to the identity, though this is not feasible for larger  $n$ .

### 5.3 Gridded synthetic data example

We now profile a full objective function and gradient evaluation for the MLE problem for  $\theta$ . As before, we consider the Matérn kernel of (12) with noise  $\sigma_N^2 = 1 \times 10^{-4}$ .

We set the parameter vector at  $\theta = [10, 7]$  and again take the observation locations to be a

Table 5: Runtime for one objective function and gradient evaluation (*i.e.*, the work for a single iteration) on a uniform grid of observations.

$\epsilon_{\text{peel}}$	$n$	$t_{\text{fact}}$ (s)	$t_{\text{peel,weak}}$ (s)	$t_{\text{total}}$ (s)
$1 \times 10^{-6}$	$64^2$	$7.34 \times 10^0$	$1.23 \times 10^1$	$1.96 \times 10^1$
	$128^2$	$4.86 \times 10^1$	$1.20 \times 10^2$	$1.68 \times 10^2$
	$256^2$	$2.73 \times 10^2$	$1.22 \times 10^3$	$1.49 \times 10^3$
	$512^2$	$1.41 \times 10^3$	$1.39 \times 10^4$	$1.53 \times 10^4$
$1 \times 10^{-8}$	$64^2$	$9.70 \times 10^0$	$1.55 \times 10^1$	$2.52 \times 10^1$
	$128^2$	$7.08 \times 10^1$	$1.68 \times 10^2$	$2.39 \times 10^2$
	$256^2$	$4.28 \times 10^2$	$1.76 \times 10^3$	$2.29 \times 10^3$
	$512^2$	$2.64 \times 10^3$	$1.54 \times 10^4$	$1.80 \times 10^4$

regular  $\sqrt{n} \times \sqrt{n}$  grid discretizing the square  $[0, 100]^2$ . Evaluating  $\ell(\theta)$  and  $g_i$  for  $i = 1, \dots, 2$  then requires three skeletonization factorizations and two different trace approximations. We investigate the algorithm’s performance for two different peeling tolerances  $\epsilon_{\text{peel}}$ , and in each case take the factorization tolerance to be  $\epsilon_{\text{fact}} = \frac{1}{1000}\epsilon_{\text{peel}}$ . For varying  $n$  between  $64^2$  and  $512^2$ , we measured the runtime of both the factorization portion and peeling portion of Algorithm 1. We note that, given the factorizations and peeled trace estimates, the remaining pieces of Algorithm 1 are several orders of magnitude less costly in terms of runtime.

In Fig. 7 (left), we plot the total runtime for a single objective function and gradient evaluation for the uniform grid of observations (corresponding data in Table 5). We see from the figure that the runtime seems to scale as roughly  $O(n^{3/2})$  with the number of observations; a least-squares fit of the data gives  $O(n^{1.6})$ . As can be seen in the table, the amount of time spent in calculating the recursive skeletonization factorizations is roughly an order of magnitude less than the time spent in the peeling trace approximation, and, further, scales slightly better than peeling for this example.

## 5.4 Scattered synthetic data example

While all examples thus far have used a regular grid of observations, our framework does not rely on this assumption. To complement the examples on gridded observations, we repeat the same experiment from the previous section with real-world observation locations coming from release 2.5 of the International Comprehensive Ocean-Atmosphere Data Set (ICOADS) [45] obtained from the National Center for Atmospheric Research at <http://rda.ucar.edu/datasets/ds540.0/>. We subselect from ICOADS a set of sea surface temperatures measured at varying locations in the North Atlantic ocean between the years 2008 and 2014. Restricting the data to observations made in the month of July across all years and obtain roughly 300,000 unique observation locations and corresponding sea surface temperature

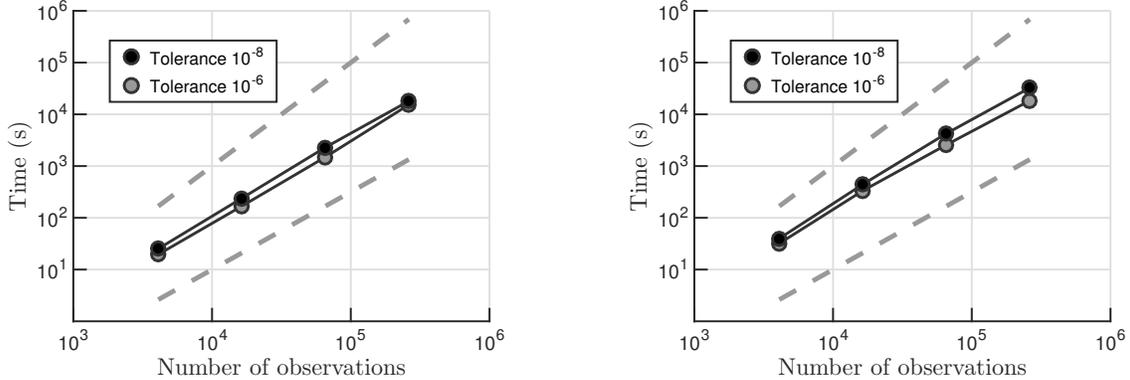


Figure 7: On the left we plot the total runtime of evaluating a single objective function and gradient for the uniform grid example of Section 5.3 as a function of the total number of observations for two different tolerances. The top trend line shows  $O(n^2)$  scaling and the bottom shows  $O(n^{3/2})$  scaling. On the right we plot the corresponding results for the scattered data of Section 5.4 with the same trend lines. We observe that the scaling in all cases looks like  $O(n^{3/2})$ .

measurements, some of which can be seen in Fig. 8.

Because large-scale spatial measurements typically cover a non-trivial range of latitudes and longitudes, the development of valid covariance functions on the entire sphere that respect the proper distance metric has been the subject of much recent work, see, e.g., Gneiting [14] and related work [33, 26, 48]. As the focus of this manuscript is not statistical modeling, we employ a simplified model based on Mercator projection of the observations to two spatial dimensions. Note that the choice of axis scaling in the Mercator projection is arbitrary; in our convention the horizontal axis spans 90 units and the vertical axis spans 70 units.

To perform scaling tests on the cost of an objective function and gradient evaluation according to Algorithm 1, we subselect from our full dataset of unique observation locations by drawing observations uniformly at random without replacement. Figure 7 (right) shows the runtime scaling results as a function of the number of observations, with corresponding data in Table 6. We see that the runtime scaling for the scattered observations follows essentially the same scaling behavior as the gridded observations from Section 5.3, with observed complexity between  $O(n^{1.5})$  and  $O(n^{1.6})$ . Again, the skeletonization factorizations take considerably less time than the trace estimation.

As an illustrative example of the full power of Algorithm 1 in context, we take a subset of  $n = 2^{16}$  scattered observations and realize an instance of a Gaussian process at those locations with true parameter vector  $\theta^* = [10, 7]$  and noise parameter  $\sigma_N^2 = 1 \times 10^{-4}$  to generate the observation vector  $z$ . Setting the peel tolerance to  $\epsilon_{\text{peel}} = 1 \times 10^{-6}$  and the factorization tolerance to  $\epsilon_{\text{fact}} = 1 \times 10^{-9}$ , we plugged our approximate log-likelihood and gradient routines into the MATLAB<sup>®</sup> routine `fminunc` for unconstrained optimization using the quasi-Newton option. Starting from an initial guess of  $\theta_0 = [3, 30]$ , we found that after 13 iterations (14 calls to Algorithm 1) the first-order optimality as measured by the  $\ell_\infty$ -

Table 6: Runtime for one objective function and gradient evaluation (*i.e.*, the work for a single iteration) on scattered observations with locations from ICOADS.

$\epsilon_{\text{peel}}$	$n$	$t_{\text{fact}}$ (s)	$t_{\text{peel,weak}}$ (s)	$t_{\text{total}}$ (s)
$1 \times 10^{-6}$	$2^{12}$	$7.82 \times 10^0$	$2.41 \times 10^1$	$3.19 \times 10^1$
	$2^{14}$	$4.25 \times 10^1$	$2.91 \times 10^2$	$3.34 \times 10^2$
	$2^{16}$	$8.60 \times 10^1$	$2.49 \times 10^3$	$2.57 \times 10^3$
	$2^{18}$	$4.84 \times 10^2$	$1.79 \times 10^4$	$1.84 \times 10^4$
$1 \times 10^{-8}$	$2^{12}$	$1.06 \times 10^1$	$2.92 \times 10^1$	$3.98 \times 10^1$
	$2^{14}$	$5.96 \times 10^1$	$3.83 \times 10^2$	$4.43 \times 10^2$
	$2^{16}$	$2.72 \times 10^2$	$3.98 \times 10^3$	$4.26 \times 10^3$
	$2^{18}$	$1.03 \times 10^3$	$3.15 \times 10^4$	$3.25 \times 10^4$

norm of the gradient had been reduced by three orders of magnitude, yielding an estimate of  $\hat{\theta} = [10.0487, 7.0496]$  after approximately  $4.86 \times 10^4$  seconds.

**Remark 10.** *While a large percentage of this runtime was spent in the peeling algorithm, we find it worthwhile to note that in this example the use of our gradient approximation proved essential—using finite difference approximations to the gradient led to stagnation at the first iteration, even with a factorization tolerance  $\epsilon_{\text{fact}} = 1 \times 10^{-15}$ , *i.e.*, at the limits of machine precision.*

Because the number of iterations to convergence depends on many factors (*e.g.*, the choice of optimization algorithm, how well the data can be modeled by a Gaussian process, and many convergence tolerances depending on the chosen algorithm), we do not find it useful to attempt to profile the full minimization algorithm more extensively than this, but direct the reader instead to the single-iteration results.

## 5.5 Scattered ocean data example

While the factorizations and peeling in Algorithm 1 depend only on the locations of the observations and not their values, the log-likelihood  $\ell(\cdot)$  can have a more complicated shape with real observations  $z$  than with synthetic data, which may impact the required tolerance parameters  $\epsilon_{\text{fact}}$  and  $\epsilon_{\text{peel}}$  and the difficulty of MLE. Further, there are a number of practical considerations relevant for real data not addressed thus far in our synthetic examples.

As a refinement of (1), suppose now that the data are distributed according to  $z \sim N(\mu \mathbf{1}, \sigma^2 \Sigma(\theta))$ , where  $\mathbf{1} \in \mathbb{R}^n$  is the all-ones vector,  $\mu$  and  $\sigma^2$  represent the constant but unknown mean and variance level, and our parameterized Matérn model is given by (for

Table 7: Parameter estimates  $\hat{\theta} = [\hat{\rho}, \hat{\sigma}_N^2]$ , the corresponding mean and variance level  $\mu(\hat{\theta})$  and  $\sigma^2(\hat{\theta})$ , and the log-likelihood values of the fitted parameters for the model with kernel (14). We note that for  $\nu = 1/2$ , the estimate of  $\sigma_N^2$  is at its lower bound.

$\nu$	$\hat{\rho}$	$\hat{\sigma}_N^2$	$\mu(\hat{\theta})$	$\sigma^2(\hat{\theta})$	$\ell(\hat{\theta})$
1/2	$2.12 \times 10^{+0}$	$1.00 \times 10^{-8}$	$-5.65 \times 10^{-1}$	$9.97 \times 10^{-1}$	$9.28 \times 10^{+4}$
3/2	$2.52 \times 10^{-1}$	$4.37 \times 10^{-3}$	$-4.38 \times 10^{-1}$	$8.71 \times 10^{-1}$	$8.42 \times 10^{+4}$
5/2	$1.78 \times 10^{-1}$	$5.67 \times 10^{-3}$	$-4.30 \times 10^{-1}$	$8.57 \times 10^{-1}$	$7.74 \times 10^{+4}$

several different  $\nu$ )

$$[\Sigma(\theta)]_{ij} = \frac{1}{\Gamma(\nu)2^{\nu-1}} \left( \frac{\sqrt{2\nu}r_{ij}}{\rho} \right)^\nu K_\nu \left( \frac{\sqrt{2\nu}r_{ij}}{\rho} \right) + \sigma_N^2 \delta_{ij} \quad (14)$$

$$= \begin{cases} \exp\left(-\frac{r_{ij}}{\rho}\right), & \nu = 1/2, \\ \left(1 + \frac{\sqrt{3}r_{ij}}{\rho}\right) \exp\left(-\frac{\sqrt{3}r_{ij}}{\rho}\right) + \sigma_N^2 \delta_{ij}, & \nu = 3/2, \\ \left(1 + \frac{\sqrt{5}r_{ij}}{\rho} + \frac{5r_{ij}^2}{3\rho^2}\right) \exp\left(-\frac{\sqrt{5}r_{ij}}{\rho}\right) + \sigma_N^2 \delta_{ij}, & \nu = 5/2, \end{cases}$$

with  $r_{ij} = \|x_i - x_j\|$ . In this example, the parameter vector is  $\theta = [\rho, \sigma_N^2]$ , consisting of a single correlation length parameter and the noise level. To optimize the new log-likelihood over  $\mu$ ,  $\sigma^2$ , and  $\theta$ , we note that optimization over  $\mu$  and  $\sigma^2$  results in closed form expressions for these parameters in terms of  $\theta$ , which may then be substituted back into Eq. (4) to obtain the *log profile likelihood* for this model

$$\tilde{\ell}(\theta) \equiv -\frac{1}{2} \log |\Sigma| - \frac{n}{2} \log (z^T (\Sigma + \mathbf{1}\mathbf{1}^T)^{-1} z) + \frac{n}{2} (\log n - 1 - 2\pi) \quad (15)$$

with gradient components given by

$$\tilde{g}_i \equiv -\frac{1}{2} \text{Tr}(\Sigma^{-1} \Sigma_i) + \frac{n}{2} \left( \frac{z^T (\Sigma + \mathbf{1}\mathbf{1}^T)^{-1} \Sigma_i (\Sigma + \mathbf{1}\mathbf{1}^T)^{-1} z}{z^T (\Sigma + \mathbf{1}\mathbf{1}^T)^{-1} z} \right), \quad i = 1, \dots, p.$$

Optimization of this new model fits neatly into the computational framework of Algorithm 1 with trivial modifications. The new model has the advantage of greater plausibility, though it still admits many further improvements.

From the full set of sea surface temperature observations, we subselected  $n = 2^{17}$  unique temperature measurements corresponding to observations between July 2013 and August 2013. Taking  $\epsilon_{\text{fact}} = 1 \times 10^{-9}$  and  $\epsilon_{\text{peel}} = 1 \times 10^{-6}$ , we use the MATLAB<sup>®</sup> optimization routine `fmincon` with the ‘SQP’ option to estimate the correlation length parameter  $\theta_1 = \rho$  and noise parameter  $\theta_2 = \sigma_N^2$  for the standardized temperature measurements. For  $\nu = 3/2$  and  $\nu = 5/2$  this was accomplished by numerically maximizing (15) subject to the lower-bound constraint  $\sigma_N^2 \geq 1 \times 10^{-5}$ , which was necessary to ensure  $\Sigma$  was not numerically

rank-deficient. For  $\nu = 1/2$  the covariance matrix  $\Sigma$  is naturally better conditioned so a looser lower-bound  $\sigma_N^2 \geq 1 \times 10^{-8}$  was used. The tolerances dictating the minimum step-size and minimum change in objective function between successive iterates were both set to  $1 \times 10^{-6}$ .

Due to the non-convex nature of the problem, we tried several choices of starting parameter for each  $\nu$ ; the results we present are for the best initialization in each case. For choices of initial parameters leading to convergent iterates (*e.g.*,  $\theta_0 = [5, 1]$  or  $\theta_0 = [1 \times 10^{-1}, 1 \times 10^{-3}]$ ), the converged solutions all agreed to the specified tolerance and the objective function value at the optimal points agreed to six digits. For some choices of initial parameters, the optimization terminated prematurely due to the relative improvement tolerances used to evaluate convergence (*i.e.*, when the initial parameters are very poor, even a large improvement relative to the initial parameters can be far from the best choice of parameters). We did not observe any evidence of multiple local optima, though the possibility that our reported parameters are globally suboptimal cannot be ruled out.

The results of our numerical optimization for each choice of  $\nu$  can be seen in Table 7, where in each case optimization terminated due to the step-size tolerance. At each corresponding  $\hat{\theta}$ , however, we note that the gradient is small relative to the objective function.

Of the three different models, we find that the fitted model for  $\nu = 1/2$  gives the best fit as measured both by comparative likelihood and qualitatively (see Fig. 8). Since  $\nu$  dictates the smoothness of denoised process, these results imply that the best description of the observed data among our choices is the one with the least assumptions on smoothness. We caution that this does not preclude a much better fit with a more sophisticated model, but this simple example illustrates that our framework is effective for MLE even for real observations.

## 6 Conclusions

The framework for Gaussian process MLE presented in this paper and summarized in Algorithm 1 provides a straightforward method of leveraging hierarchical matrix representations from scientific computing for fast computations with kernelized covariance matrices arising in spatial statistics. The general linear algebraic approach to approximating off-diagonal blocks of the covariance matrix to a specified error tolerance by adaptively determining their ranks gives a flexible way of attaining high-accuracy approximations with reasonable runtime. A further merit to this approach is that it does not rely on having gridded observations or a translation-invariant covariance kernel.

While in this paper we have focused on maximum likelihood estimation for Gaussian processes, these methods are equally viable for the Bayesian setting. For example, computing maximum *a posteriori* estimates follows essentially the same approach with the addition of a term depending on the prior. Further, sampling from the posterior distribution of  $\theta$  in a Bayesian setting can be accomplished using standard Markov chain Monte Carlo methods based on quickly evaluating the likelihood and posterior. This can also be combined with Remark 8 for a fully Bayesian treatment.

Our numerical results in Section 5 show that our framework scales favorably when applied

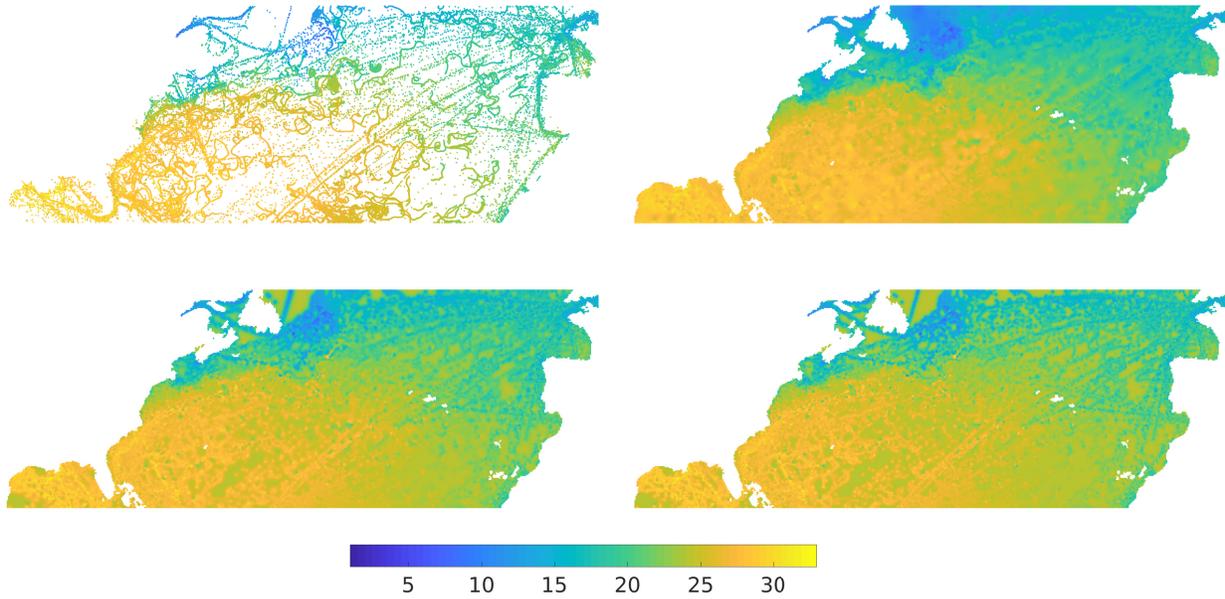


Figure 8: In the top-left plot, we show a subselection of  $n = 2^{17}$  Atlantic ocean surface temperature measurements from ICOADS projected to a 2D plane through Mercator projection and then scattered on top of a white background for visualization. Fitting the model with constant mean and covariance given by (14), we use the estimated parameters to find the conditional mean temperatures throughout this region of the Atlantic for  $\nu = 1/2$  (top right),  $\nu = 3/2$  (bottom left), and  $\nu = 5/2$  (bottom right). The color bar shows the estimated sea surface temperature in Celsius.

to our two test cases (the rational quadratic and Matérn family kernels), leading to runtimes scaling approximately as  $O(n^{3/2})$  with  $n$  the number of observations. Further, we see that the tolerance parameter  $\epsilon_{\text{peel}}$  controlling the rank of off-diagonal block approximations in the peeling algorithm serves as a good estimate of the order of the error in the ultimate trace approximation as well. In practice, the tolerances  $\epsilon_{\text{peel}}$  and  $\epsilon_{\text{fact}}$  can be dynamically modified during the course of the maximum likelihood process for performance, *e.g.*, one could use relatively low-accuracy approximations during initial iterations of the optimization routine and slowly decrease the tolerance as the optimization progresses.

While the methods and complexity estimates discussed in this paper relate to the case of two spatial dimensions, they trivially extend to one-dimensional (time-series) data or quasi-two-dimensional data, *e.g.*, observations in three dimensions where the sampling density in one dimension is much smaller than in the other two. While the same methods apply in principle to truly three-dimensional data, the corresponding computational complexity is bottlenecked by the cost of using peeling to obtain a high-accuracy trace estimate of the matrices  $\Sigma^{-1}\Sigma_i$  for  $i = 1, \dots, p$  due to increased rank growth. In fact, even in the two-dimensional case it is clear from Table 5 and Table 6 that the most expensive piece of our framework in practice is determining these traces. One solution is to instead use the hierarchical matrix representations inside of an estimator such as that of Stein et al. [39], which obviates the need for the trace. For the true MLE, however, future work on efficiently computing this trace to high accuracy is necessary. Given a method for efficiently computing this trace for 3D data, we expect that related factorizations based on more sophisticated use of skeletonization should give complexities for computing the log-likelihood and gradient that are as good as or better than those we obtain with recursive skeletonization in the 2D case. For example, the hierarchical interpolative factorization [24] (which uses further levels of compression to mitigate rank-growth of off-diagonal blocks) may be used in our framework as an efficient method of applying  $\Sigma_i$  and  $\Sigma^{-1}$  and computing the log-determinant of  $\Sigma$  for 3D problems.

While Gaussian process regression is widely used for data in  $\mathbb{R}^d$  with  $d$  much larger than three, the methods of this paper are designed with spatial data in mind. In particular, in the high-dimensional setting the geometry of the observations becomes very important for efficiency. If the data can be well-approximated according to an intrinsic low-dimensional embedding that is efficient to identify, there is hope for efficient approximations using hierarchical rank structure (see, for example, Yu et al. [47]). However, in general we expect that rank-structured factorizations will continue to be most effective for low-dimensional spatial applications.

## Acknowledgments

The authors thank Matthias Cremon, Eileen Martin, Sven Schmit, and Austin Benson for useful discussion on Gaussian process regression, the anonymous reviewers for thoughtful comments that improved the presentation of this paper, and Stanford University and the Stanford Research Computing Center for providing computational resources and support

that have contributed to these research results.

## References

- [1] S. AMBIKASARAN, D. FOREMAN-MACKEY, L. GREENGARD, D. W. HOGG, AND M. O'NEIL, *Fast direct methods for Gaussian processes*, IEEE Transactions on Pattern Analysis and Machine Intelligence, 38 (2016), pp. 252–265.
- [2] M. ANITESCU, J. CHEN, AND L. WANG, *A matrix-free approach for solving the parametric Gaussian process maximum likelihood problem*, SIAM Journal on Scientific Computing, 34 (2012), pp. A240–A262.
- [3] E. AUNE, D. P. SIMPSON, AND J. EIDSVIK, *Parameter estimation in high dimensional Gaussian distributions*, Statistics and Computing, 24 (2014), pp. 247–263, doi:10.1007/s11222-012-9368-y, <http://dx.doi.org/10.1007/s11222-012-9368-y>.
- [4] S. BÖRM AND J. GARCKE, *Approximating Gaussian processes with  $\mathcal{H}^2$ -matrices*, in Proceedings of the 18th European Conference on Machine Learning, Springer, 2007, pp. 42–53.
- [5] J. E. CASTRILLÓN-CANDÁS, M. G. GENTON, AND R. YOKOTA, *Multi-level restricted maximum likelihood covariance estimation and Kriging for large non-gridded spatial datasets*, Spatial Statistics, (2015).
- [6] S. CHANDRASEKARAN, P. DEWILDE, M. GU, W. LYONS, AND T. PALS, *A fast solver for HSS representations via sparse matrices*, SIAM Journal on Matrix Analysis and Applications, 29 (2007), pp. 67–81.
- [7] S. CHANDRASEKARAN, M. GU, AND T. PALS, *A fast ULV decomposition solver for hierarchically semiseparable representations*, SIAM Journal on Matrix Analysis and Applications, 28 (2006), pp. 603–622.
- [8] H. CHENG, Z. GIMBUTAS, P.-G. MARTINSSON, AND V. ROKHLIN, *On the compression of low rank matrices*, SIAM Journal on Scientific Computing, 26 (2005), pp. 1389–1404.
- [9] E. CORONA, P.-G. MARTINSSON, AND D. ZORIN, *An  $O(N)$  direct solver for integral equations on the plane*, Applied and Computational Harmonic Analysis, 38 (2015), pp. 284 – 317.
- [10] N. CRESSIE AND G. JOHANNESSON, *Fixed rank Kriging for very large spatial data sets*, Journal of the Royal Statistical Society, Series B, 70 (2008), pp. 209–226.
- [11] J. EIDSVIK, B. A. SHABY, B. J. REICH, M. WHEELER, AND J. NIEMI, *Estimation and prediction in spatial models with block composite likelihoods*,

- Journal of Computational and Graphical Statistics, 23 (2014), pp. 295–315, doi:10.1080/10618600.2012.760460, <http://dx.doi.org/10.1080/10618600.2012.760460>, arXiv:<http://dx.doi.org/10.1080/10618600.2012.760460>.
- [12] R. FURRER, M. G. GENTON, AND D. NYCHKA, *Covariance tapering for interpolation of large spatial datasets*, Journal of Computational and Graphical Statistics, 15 (2006), pp. 502–523.
- [13] A. GILLMAN, P. M. YOUNG, AND P.-G. MARTINSSON, *A direct solver with  $O(N)$  complexity for integral equations on one-dimensional domains*, Frontiers of Mathematics in China, 7 (2012), pp. 217–247.
- [14] T. GNEITING, *Strictly and non-strictly positive definite functions on spheres*, Bernoulli, 19 (2013), pp. 1327–1349, doi:10.3150/12-BEJSP06, <http://dx.doi.org/10.3150/12-BEJSP06>.
- [15] L. GREENGARD, D. GUEYFFIER, P.-G. MARTINSSON, AND V. ROKHLIN, *Fast direct solvers for integral equations in complex three-dimensional domains*, Acta Numerica, 18 (2009), pp. 243–275.
- [16] L. GREENGARD AND V. ROKHLIN, *On the numerical solution of two-point boundary value problems*, Communications on Pure and Applied Mathematics, 44 (1991), pp. 419–452.
- [17] W. HACKBUSCH, *A sparse matrix arithmetic based on  $\mathcal{H}$ -matrices. Part I: Introduction to  $\mathcal{H}$ -matrices*, Computing, 62 (1999), pp. 89–108.
- [18] W. HACKBUSCH, *Hierarchical Matrices: Algorithms and Analysis*, Springer Series in Computational Mathematics, Springer-Verlag Berlin Heidelberg, 2015.
- [19] W. HACKBUSCH AND S. BÖRM, *Data-sparse approximation by adaptive  $\mathcal{H}^2$ -matrices*, Computing, 69 (2002), pp. 1–35.
- [20] W. HACKBUSCH AND B. N. KHOROMSKIJ, *A sparse  $\mathcal{H}$ -matrix arithmetic. Part II: Application to multi-dimensional problems*, Computing, 64 (2000), pp. 21–47.
- [21] W. HACKBUSCH, B. N. KHOROMSKIJ, AND R. KRIEMANN, *Hierarchical matrices based on a weak admissibility criterion*, Computing, 73 (2004), pp. 207–243.
- [22] N. HALKO, P.-G. MARTINSSON, AND J. A. TROPP, *Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions*, SIAM Review, 53 (2011), pp. 217–288.
- [23] K. L. HO AND L. GREENGARD, *A fast direct solver for structured linear systems by recursive skeletonization*, SIAM Journal on Scientific Computing, 34 (2012), pp. A2507–A2532.

- [24] K. L. HO AND L. YING, *Hierarchical interpolative factorization for elliptic operators: Integral equations*, Communications on Pure and Applied Mathematics, (2015).
- [25] M. F. HUTCHINSON, *A stochastic estimator of the trace of the influence matrix for Laplacian smoothing splines*, Communications in Statistics - Simulation and Computation, 19 (1990), pp. 433–450.
- [26] M. JUN AND M. L. STEIN, *Nonstationary covariance models for global data*, Ann. Appl. Stat., 2 (2008), pp. 1271–1289, doi:10.1214/08-AOAS183, <http://dx.doi.org/10.1214/08-AOAS183>.
- [27] B. N. KHOROMSKIJ, A. LITVINENKO, AND H. G. MATTHIES, *Application of hierarchical matrices for computing the Karhunen–Loève expansion*, Computing, 84 (2008), pp. 49–67.
- [28] L. LIN, J. LU, AND L. YING, *Fast construction of hierarchical matrix representation from matrix-vector multiplication*, Journal of Computational Physics, 230 (2011), pp. 4071–4087, arXiv:1001.0149.
- [29] F. LINDGREN, H. RUE, AND J. LINDSTRÖM, *An explicit link between Gaussian fields and Gaussian Markov random fields: the stochastic partial differential equation approach*, Journal of the Royal Statistical Society: Series B (Statistical Methodology), 73 (2011), pp. 423–498, doi:10.1111/j.1467-9868.2011.00777.x, <http://dx.doi.org/10.1111/j.1467-9868.2011.00777.x>.
- [30] P.-G. MARTINSSON AND V. ROKHLIN, *A fast direct solver for boundary integral equations in two dimensions*, Journal of Computational Physics, 205 (2005), pp. 1–23.
- [31] G. MATHERON, *Principles of geostatistics*, Economic geology, 58 (1963), pp. 1246–1266.
- [32] V. MINDEN, K. L. HO, A. DAMLE, AND L. YING, *A recursive skeletonization factorization based on strong admissibility*, Multiscale Modeling & Simulation, 15 (2017), pp. 768–796, doi:10.1137/16M1095949, <http://dx.doi.org/10.1137/16M1095949>, arXiv:<http://dx.doi.org/10.1137/16M1095949>.
- [33] E. PORCU, M. BEVILACQUA, AND M. G. GENTON, *Spatio-temporal covariance and cross-covariance functions of the great circle distance on a sphere*, Journal of the American Statistical Association, 111 (2016), pp. 888–898, doi:10.1080/01621459.2015.1072541, <http://dx.doi.org/10.1080/01621459.2015.1072541>, arXiv:<http://dx.doi.org/10.1080/01621459.2015.1072541>.
- [34] V. ROKHLIN AND M. TYGERT, *A fast randomized algorithm for overdetermined linear least-squares regression*, Proceedings of the National Academy of Sciences, 105 (2008), pp. 13212–13217.

- [35] H. SANG AND J. Z. HUANG, *A full-scale approximation of covariance functions for large spatial data sets*, Journal of the Royal Statistical Society, Series B, 74 (2012), pp. 111–132.
- [36] P. STARR AND V. ROKHLIN, *On the numerical solution of two-point boundary value problems II*, Communications on Pure and Applied Mathematics, 47 (1994), pp. 1117–1159.
- [37] M. L. STEIN, *Interpolation of Spatial Data: Some Theory for Kriging*, Springer Series in Statistics, Springer New York, 1999.
- [38] M. L. STEIN, J. CHEN, AND M. ANITESCU, *Difference filter preconditioning for large covariance matrices*, SIAM Journal on Matrix Analysis and Applications, 33 (2012), pp. 52–72.
- [39] M. L. STEIN, J. CHEN, AND M. ANITESCU, *Stochastic approximation of score functions for Gaussian processes*, Annals of Applied Statistics, 7 (2013), pp. 1162–1191, arXiv:1312.2687.
- [40] M. L. STEIN, Z. CHI, AND L. J. WELTY, *Approximating likelihoods for large spatial data sets*, Journal of the Royal Statistical Society, Series B, 66 (2004), pp. 275–296.
- [41] J. A. TROPP, *Improved analysis of the subsampled randomized Hadamard transform*, Advances in Adaptive Data Analysis, 03 (2011), pp. 115–126.
- [42] J. VANHATALO, V. PIETILÄINEN, AND A. VEHTARI, *Approximate inference for disease mapping with sparse Gaussian processes*, Statistics in Medicine, 29 (2010), pp. 1580–1607.
- [43] A. V. VECCHIA, *Estimation and model identification for continuous spatial processes*, Journal of the Royal Statistical Society, Series B, 50 (1988), pp. pp. 297–312.
- [44] P. WHITTLE, *On stationary processes in the plane*, Biometrika, 41 (1954), pp. pp. 434–449.
- [45] S. D. WOODRUFF, S. J. WORLEY, S. J. LUBKER, Z. JI, J. E. FREEMAN, D. I. BERRY, P. BROHAN, E. C. KENT, R. W. REYNOLDS, S. R. SMITH, AND C. WILKINSON, *ICOADS release 2.5: extensions and enhancements to the surface marine meteorological archive*, International Journal of Climatology, 31 (2011), pp. 951–967.
- [46] J. XIA, S. CHANDRASEKARAN, M. GU, AND X. S. LI, *Fast algorithms for hierarchically semiseparable matrices*, Numerical Linear Algebra With Applications, 17 (2010), p. 953–976.

- [47] C. D. YU, W. B. MARCH, AND G. BIROS, *An  $n \log n$  parallel fast direct solver for kernel matrices*, in 2017 IEEE International Parallel and Distributed Processing Symposium (IPDPS), May 2017, pp. 886–896, doi:10.1109/IPDPS.2017.10.
- [48] M. HEATON, M. KATZFUSS, C. BERRETT, AND D. NYCHKA, *Constructing valid spatial processes on the sphere using kernel convolutions*, *Environmetrics*, 25 (2014), pp. 2–15, doi:10.1002/env.2251, <http://dx.doi.org/10.1002/env.2251>.