

HLIBCov: Likelihood Approximation With Parallel Hierarchical Matrices For Large Spatial Datasets

Alexander Litvinenko

5 *King Abdullah University of Science and Technology, Thuwal 23955-6900, Saudi Arabia*

Abstract. The main goal of this article is to introduce the parallel hierarchical matrix library HLIBpro to statistical community. We describe the HLIBCov package, which is - an extension of the HLIBpro for approximating large covariance matrices and maximizing likelihood functions. We show that an approximate Cholesky factorization of a dense matrix of size $2Mi \times 2Mi$ can be computed on a modern multi-core desktop in few minutes. Further, HLIBCov is used for estimating the unknown parameters such as the covariance length, variance and smoothness parameter of a Matérn covariance function by maximizing the joint Gaussian log-likelihood function. The computational bottleneck here is expensive linear algebra arithmetics due to large and dense covariance matrices. Therefore covariance matrices are approximated in the hierarchical (\mathcal{H} -) matrix format with computational cost $\mathcal{O}(k^2 n \log^2 n / p)$ and storage $\mathcal{O}(kn \log n)$, where the rank k is a small integer (typically $k < 25$), p the number of cores and n the number of locations on a fairly general mesh. We demonstrate a synthetic example, where the true values of known parameters are known. For reproducibility we provide the C++ code, the documentation, and the synthetic data.

Key words: Computational statistics; parallel Hierarchical matrices; Large datasets; Matérn covariance; Random Fields; Spatial statistics, HLIB, HLIBpro, Cholesky, matrix determinant

10

HLIBCov software library: Likelihood Approximation With Parallel Hierarchical Matrices For Large Spatial Datasets

Program title: HLIBCov

<http://www.global-sci.com/>

Global Science Preprint

Nature of problem: To estimate the unknown parameters (variance, smoothness parameter, and covariance length) of a covariance function by maximizing the joint Gaussian log-likelihood function with a log-linear computational cost and storage.

Software licence: GPL 2.0

CiCP scientific software URL:

Distribution format: *.cc files via github

Programming language(s): C++

Computer platform: 32 or 64

Operating system: Unix-like operating systems

Compilers: clang

RAM: 4GB and more (depending on the matrix size)

External routines/libraries: HLIBpro requires:

- BLAS/LAPACK,
- Threading Building Blocks (<http://www.threadingbuildingblocks.org/>),
- Boost (<http://www.boost.org/>),
- zlib (<http://zlib.net>),
- FFTW3 (system dependent, <http://www.fftw.org>),
- Scotch (system dependent, <http://gforge.inria.fr/projects/scotch/>),
- METIS (system dependent, <http://glaros.dtc.umn.edu/gkhome/views/metis>)
- SCons (<http://www.scons.org>)

HLIBCov requires: GNU Scientific Library (<https://www.gnu.org/software/gsl/>)

Running time: $\mathcal{O}(k^2 n \log^2 n) / p$ with p number of cores

Restrictions: 1. For multiple-core architectures with shared memory systems
2. Compiling with GCC 5 will generate linking errors.

Supplementary material and references: www.HLIBpro.com and references therein.

Additional Comments: HLIBpro is a software library implementing parallel algorithms for Hierarchical matrices. It is freely available in binary form for academic purposes. HLIBpro algorithms design for 1, 2 and 3 - dimensional problems.

1 Introduction

This paper is complementary to the paper [Litvinenko, Ying, Genton, Keyes, *Likelihood Approximation With Hierarchical Matrices For Large Spatial Datasets*, submitted to J. Computational Statistics and Data Analysis, Sept. 2017].

Novelty of this work. In this paper we use parallel hierarchical (\mathcal{H} -) matrices for approximating dense covariance matrices and the joint Gaussian log-likelihood with computational complexity $\mathcal{O}(k^\alpha n \log^\alpha n / p)$, where p is the number of cores, n is the number of measurements; $k \ll n$ is the maximal rank, used in the hierarchical matrix, which defines the quality of the approximation; and $\alpha = 1$ (for a matrix-vector product) or 2 (for the Cholesky). Additionally, we estimate unknown parameters of the covariance function.

Parameter estimation. Problem settings. Let n be the number of spatial measurements \mathbf{Z} located irregularly across a given geographical region at locations $\mathbf{s} := \{\mathbf{s}_1, \dots, \mathbf{s}_n\} \in \mathbb{R}^d$, $d \geq 1$. These data are frequently modeled as a realization from a stationary Gaussian spatial random field. Specifically, we let $\mathbf{Z} = \{Z(\mathbf{s}_1), \dots, Z(\mathbf{s}_n)\}^\top$, where $Z(\mathbf{s})$ is a Gaussian random field. Then, we assume that \mathbf{Z} has mean zero and a stationary parametric covariance function $C(\mathbf{h}; \boldsymbol{\theta}) = \text{cov}\{Z(\mathbf{s}), Z(\mathbf{s} + \mathbf{h})\}$, where $\mathbf{h} \in \mathbb{R}^d$ is a spatial lag vector and $\boldsymbol{\theta} \in \mathbb{R}^q$ is the unknown parameter vector of interest. To infer unknown parameters $\boldsymbol{\theta}$ we maximize the Gaussian log-likelihood function:

$$\mathcal{L}(\boldsymbol{\theta}) = -\frac{n}{2} \log(2\pi) - \frac{1}{2} \log |\mathbf{C}(\boldsymbol{\theta})| - \frac{1}{2} \mathbf{Z}^\top \mathbf{C}(\boldsymbol{\theta})^{-1} \mathbf{Z}, \quad (1.1)$$

where $\mathbf{C}(\boldsymbol{\theta})_{ij} = C(\mathbf{s}_i - \mathbf{s}_j; \boldsymbol{\theta})$, $i, j = 1, \dots, n$. The maximum likelihood estimator of $\boldsymbol{\theta}$ is the value $\hat{\boldsymbol{\theta}}$ that maximizes (1.1). When the sample size n is large, the evaluation of (1.1) becomes challenging, due to the computation of the quadratic form and log-determinant of the n -by- n dense covariance matrix $\mathbf{C}(\boldsymbol{\theta})$. Indeed, this requires $\mathcal{O}(n^2)$ memory and $\mathcal{O}(n^3)$ computational steps. Hence, scalable and efficient methods that can process larger sample sizes are needed.

Definition 1.1. An \mathcal{H} -matrix approximation with the maximal rank k of the exact log-likelihood $\mathcal{L}(\boldsymbol{\theta})$ is defined by $\tilde{\mathcal{L}}(\boldsymbol{\theta}; k)$:

$$\tilde{\mathcal{L}}(\boldsymbol{\theta}; k) = -\frac{n}{2} \log(2\pi) - \sum_{i=1}^n \log \{\tilde{\mathbf{L}}_{ii}(\boldsymbol{\theta}; k)\} - \frac{1}{2} \mathbf{v}(\boldsymbol{\theta})^\top \mathbf{v}(\boldsymbol{\theta}), \quad (1.2)$$

where $\tilde{\mathbf{L}}(\boldsymbol{\theta}; k)$ is an \mathcal{H} -matrix approximation of the Cholesky factor $\mathbf{L}(\boldsymbol{\theta})$ with the maximal rank k in sub-blocks, $\mathbf{C}(\boldsymbol{\theta}) = \mathbf{L}(\boldsymbol{\theta})\mathbf{L}(\boldsymbol{\theta})^\top$, and $\mathbf{v}(\boldsymbol{\theta})$ is the solution of the linear system $\tilde{\mathbf{L}}(\boldsymbol{\theta}; k)\mathbf{v}(\boldsymbol{\theta}) = \mathbf{Z}$.

The argument k in $\tilde{\mathcal{L}}(\boldsymbol{\theta}; k)$ indicates that the fixed rank strategy, i.e., each sub-block in the covariance matrix has a rank equal to k or smaller, is used. In this case we cannot say anything about the accuracy in each sub-block (see more in Remark 2.2).

To maximize $\tilde{\mathcal{L}}(\boldsymbol{\theta}; k)$ in (1.2), we use Brent-Dekker method [14, 54]. We note that maximization of the log-likelihood function is an ill-posed problem, since even very small perturbations in the covariance matrix $\mathbf{C}(\boldsymbol{\theta})$ may result in large perturbations in the log-determinant and in the log-likelihood. A possible remedy is to take a higher rank k .

Features of the \mathcal{H} -matrix approximation. Other advantages for applying the \mathcal{H} -matrix technique are the following:

1. The class of \mathcal{H} -matrices is wide, for instance, it includes the classes of low-rank and sparse matrices;
2. $\mathbf{C}(\boldsymbol{\theta})^{-1}$, $\mathbf{C}(\boldsymbol{\theta})^{1/2}$, $\det \mathbf{C}(\boldsymbol{\theta})$, Cholesky decomposition, the Schur complement and, many others can be computed in the \mathcal{H} -matrix format [28];
3. The \mathcal{H} -matrix technique is well-studied and has solid theory, many examples, and multiple sequential and parallel implementations, no specific MPI or OpenMP knowledge is needed;
4. The \mathcal{H} -matrix accuracy is controllable by the rank, k . The full rank gives an exact representation;
5. The \mathcal{H} -matrix format preserves the structure (in contrast, for instance, to sparse matrices) after the Cholesky decomposition and the inverse have been computed, although with possibly a larger rank;

Related work. Stationary covariance functions, which have block Toeplitz or block circulant structure can be resolved with the Fast Fourier Transform (FFT) with the computing cost $\mathcal{O}(n \log n)$ [22]. However, this approach either does not work for data measured at irregularly spaced locations or requires expensive, non-trivial modifications.

Recently, a large amount of research has been devoted to approximating large covariance matrices: for example, covariance tapering [21, 36, 61], likelihood approximations in both the spatial [65] and spectral [19] domains, latent processes such as Gaussian predictive processes [5] and fixed-rank kriging [16], and Gaussian Markov random-field approximations [20, 43, 58, 59]; see [66] for a review. Each of these methods has its strengths and drawbacks [63, 64, 67]. Some other ideas include the nearest-neighbor Gaussian process models [17], multiresolution Gaussian process models [52], equivalent

kriging [39], multi-level restricted Gaussian maximum likelihood estimators [15], and hierarchical low-rank approximations [35]. Bayesian approaches to identify unknown or uncertain parameters could be also applied [47, 50, 50, 53, 56, 57].

Previous results [28] show that the \mathcal{H} -matrix technique is very stable when approximating the covariance matrix itself [3, 4, 11, 32, 37, 60], its inverse [3, 9], its Cholesky decomposition [7, 8], and the conditional covariance matrix [28, 44, 45].

Recently, the maximum likelihood estimator for parameter-fitting Gaussian observations with a Matérn covariance matrix was computed via a framework for unstructured observations in two spatial dimensions, which allowed the evaluation of the log-likelihood and its gradient with computational complexity $\mathcal{O}(n^{3/2})$, where n was the number of observations; the method relied on the recursive skeletonization factorization procedure [33, 48]. However, the consequences of the approximation on the maximum likelihood estimator were not studied.

In [11] authors computed the exact solution of Gaussian process regression via replacing the kernel matrix by a data-sparse approximation, called an \mathcal{H}^2 -matrix technique. An \mathcal{H}^2 -matrix approximation has $\mathcal{O}(kn)$ computational complexity and storage cost, where k is a parameter controlling the accuracy of the approximation.

Memory efficient kernel approximation was offered by [62]. The authors first make the observation that the structure of shift-invariant kernels changes from low-rank to block-diagonal (without any low-rank structure) when varying the scale parameter. Based on this observation, the authors propose a new kernel approximation algorithm – Memory Efficient Kernel Approximation, which considers both low-rank and clustering structure of the kernel matrix. They show that the resulting algorithm outperforms state-of-the-art low-rank kernel approximation methods regarding speed, approximation error, and memory usage.

In [4] authors estimated the covariance matrix of a set of normally distributed random vectors. To overcome large numerical issues in high-dimensional regime, they computed (dense) inverses of sparse covariance matrices by the usage of \mathcal{H} -matrices.

In [32] authors offered methods for rapid computation of separable expansions for the approximation of random fields. In particular, authors suggested the pivoted Cholesky method and provided a-posteriori error estimate in the trace norm. Low-rank tensor techniques in kriging are introduced in [51]. BigQUIC method for sparse inverse covariance estimation for a million variables was introduced in [34]. This method can solve 1 million dimensional ℓ_1 regularized Gaussian MLE problems. In [60] authors applied \mathcal{H} -matrices to linear inverse problems in large-scale geostatistics. In [3] authors used the

\mathcal{H}^2 matrix technique for solving large-scale stochastic linear inverse problems with applications in the subsurface modeling. In [2] authors showed that for the most commonly used covariance functions can be hierarchically factored into a product of block low-rank updates of the identity matrix, yielding an $\mathcal{O}(n \log n)$ algorithm for inversion. Additionally, they demonstrate that this factorization enables the evaluation of the determinant, permitting the direct calculation of probabilities in high dimensions.

The structure of the paper. In Section 2 we introduce the methodology and algorithms. We review the \mathcal{H} -matrix technique, the \mathcal{H} -matrix approximations of Matérn covariance functions and Gaussian likelihood functions. In Section 3 we estimate the memory storage and computing costs. In Section 4 we describe software installation details, procedures of the HLIBCov code and the algorithm for parameter estimation. Estimation of unknown parameters is reported in Section 5. The best practices are listed in Section 6. We end the paper with a conclusion in Section 7. The auxiliary \mathcal{H} -matrix details are provided in the Appendix A.

2 Methodology and algorithms

2.1 Matérn covariance functions

The class of Matérn covariance functions [26] is very widely used in spatial statistics, geostatistics, machine learning, image analysis, and other areas. The Matérn form of spatial correlations was introduced into statistics as a flexible parametric class [49], with one parameter determining the smoothness of the underlying spatial random field [31].

Let \mathbf{s} and \mathbf{s}' are any two spatial locations, and $\mathbf{h} := \|\mathbf{s} - \mathbf{s}'\|$. The Matérn class of covariance functions is defined as

$$C(h; \boldsymbol{\theta}) = \frac{\sigma^2}{2^{\nu-1} \Gamma(\nu)} \left(\frac{\mathbf{h}}{\ell} \right)^\nu K_\nu \left(\frac{\mathbf{h}}{\ell} \right), \quad (2.1)$$

with $\boldsymbol{\theta} = (\sigma^2, \ell, \nu)^\top$; where σ^2 is the variance; $\nu > 0$ controls the smoothness of the random field, with larger values of ν corresponding to smoother fields; and $\ell > 0$ is a spatial range parameter. Here K_ν denotes the modified Bessel function of the second kind of order ν . When $\nu = 1/2$, the Matérn covariance function reduces to the exponential covariance model and describes a rough field. The value $\nu = \infty$ corresponds to a Gaussian covariance model that describes a very smooth, infinitely differentiable field. Random fields with a Matérn covariance function are $\lfloor \nu - 1 \rfloor$ times mean square differentiable.

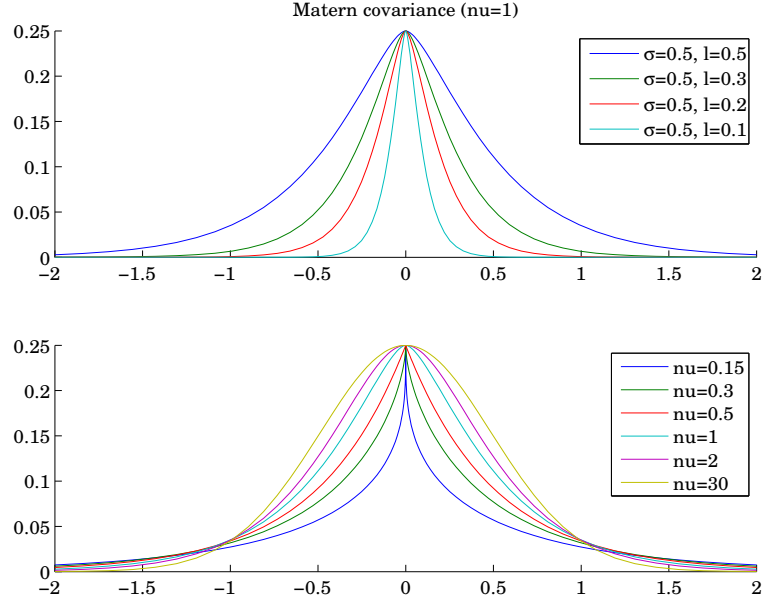


Figure 1: Matérn function for different parameters (computed in sglib [68]).

The Matérn covariance functions become especially simple when ν is half-integer [55]: $\nu = p + 1/2$, where p is a non-negative integer. In this case the covariance function is a product of an exponential and a polynomial of order p [1]:

$$C_{\nu=3/2}(\mathbf{h}) = \left(1 + \frac{\sqrt{3}\mathbf{h}}{\ell}\right) \exp\left(-\frac{\sqrt{3}\mathbf{h}}{\ell}\right), \quad C_{\nu=5/2}(\mathbf{h}) = \left(1 + \frac{\sqrt{5}\mathbf{h}}{\ell} + \frac{5\mathbf{h}^2}{3\ell^2}\right) \exp\left(-\frac{\sqrt{5}\mathbf{h}}{\ell}\right).$$

2.2 Hierarchical matrices

Hierarchical matrices have been described in detail [24, 27–30, 44]. Applications of the \mathcal{H} -matrix technique to the covariance matrices can be found in [2, 3, 11, 32, 37, 38, 46].

170 Originally the \mathcal{H} -matrix technique was developed for the approximation of stiffness
matrices coming from partial differential and integral equations [13, 24, 27]. The idea be-
hind \mathcal{H} -matrices is to approximate blocks far from diagonal (i.e., far from the singularity)
by low-rank matrices. The admissibility condition (criteria) is used to divide a given ma-
trix into sub-blocks (Fig. 4) and define which sub-blocks can be approximated well by
175 low-rank matrices and which not.

In general, covariance matrices are dense and, therefore, require $\mathcal{O}(n^2)$ units of mem-
ory for the storage and $\mathcal{O}(n^2)$ FLOPS for the matrix-vector multiplication. The \mathcal{H} -matrix

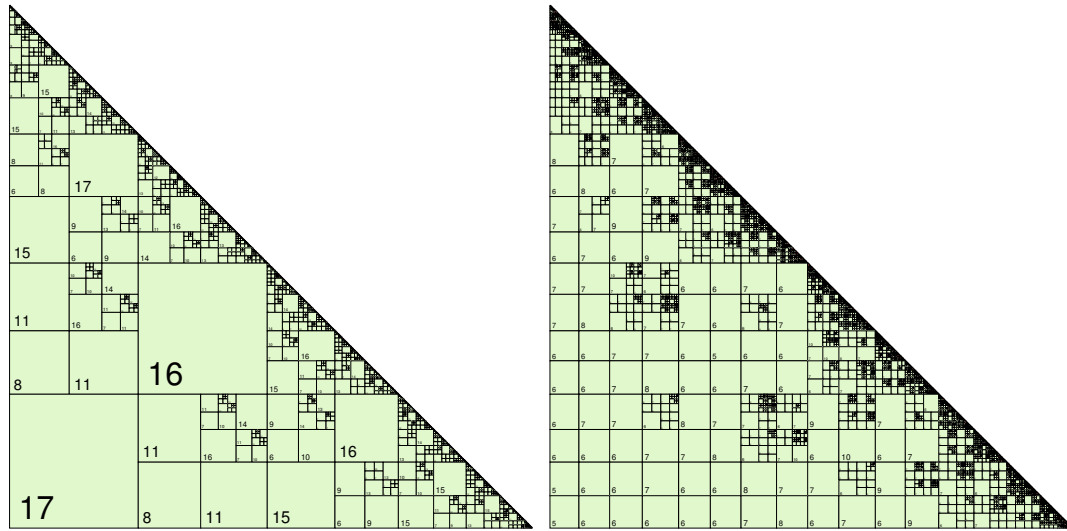


Figure 2: An example of \mathcal{H} -matrix $\tilde{\mathbf{C}}$ with coarse sub-blocks (left) and fine (right). The adaptive rank arithmetic is used.

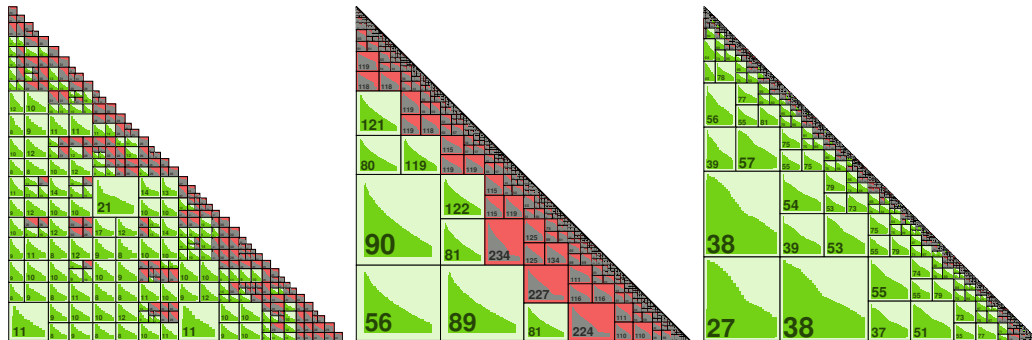


Figure 3: Examples of \mathcal{H} -matrix approximations of the Cholesky factor $\tilde{\mathbf{L}}$. The adaptive rank arithmetic is used.

technique is defined as a hierarchical partitioning of a given matrix into sub-blocks (Fig. 4, right) followed by the further approximation of the majority of these sub-block by low-rank matrices (Fig. 3). After the decomposition of the matrix into sub-blocks has been completed, an important question is how to compute the low-rank approximations. For this purpose the Adaptive Cross Approximation (or similar techniques like ACA+, HACA) algorithm [6,10,12,13,23] is used, which performs the approximations with a linear complexity $\mathcal{O}(kn)$, where n is the size of the sub-block.

Remark 2.1. The \mathcal{H} -matrix approximation error may destroy the symmetry, therefore we do the following trick $\mathbf{C} := \frac{1}{2}(\mathbf{C} + \mathbf{C}^\top)$.

Definition 2.1. We let I be an index set. A hierarchical decomposition of I we will call the cluster tree T_I (see Fig. 4 and in Appendix). A hierarchical division of the index set product, $I \times I$, into sub-blocks (Fig. 4, right) is called block cluster tree $T_{I \times I}$. The set of \mathcal{H} -matrices is

$$\mathcal{H}(T_{I \times I}, k) := \{\mathbf{C} \in \mathbb{R}^{I \times I} \mid \text{rank}(\mathbf{C}|_b) \leq k \text{ for all admissible blocks } b \text{ of } T_{I \times I}\},$$

where k is the maximum rank. Here, $\mathbf{C}|_b = (c_{ij})_{(i,j) \in b}$ denotes the matrix block of $\mathbf{C} = (c_{ij})_{i,j \in I}$ corresponding to sub-block $b \in T_{I \times I}$ (see Appendix).

Blocks that satisfy the admissibility condition can be approximated by low-rank matrices; see [27]. An \mathcal{H} -matrix approximation of \mathbf{C} is denoted by $\tilde{\mathbf{C}}$.

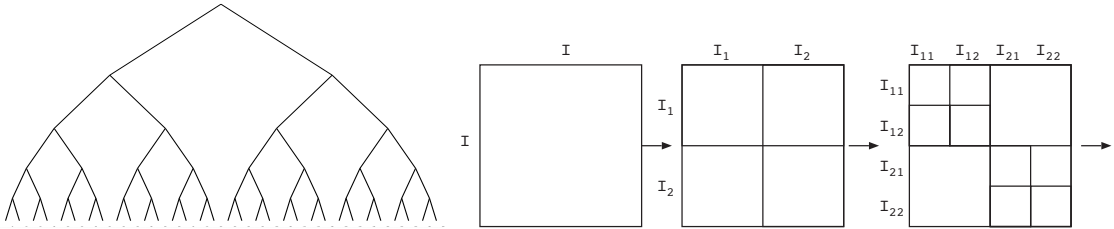


Figure 4: (left) A cluster tree T_I and (right) a block cluster tree $T_{I \times I}$ (matrix decomposition into sub-blocks).

Remark 2.2. (Fixed and adaptive rank strategies) For each sub-block we can either fix the accuracy ε or the maximal approximative rank k . In the first case we speak about the *adaptive rank strategy*, i.e., where accuracy in the spectral norm for each sub-block is ε . In the second variant we speak about the *fixed rank strategy*, i.e., each sub-block has maximal rank k . The last variant does not provide us with the reliable estimate, but make it easier

to estimate the total computing cost and memory storage. We write $\mathcal{H}(T_{I \times I}; \varepsilon)$ and $\tilde{\mathcal{L}}(\theta; \varepsilon)$ for the adaptive ranks, and $\mathcal{H}(T_{I \times I}; k)$ and $\tilde{\mathcal{L}}(\theta; k)$ for the fixed. The fixed rank strategy is useful for a priori evaluations of the computational resources and storage memory. The adaptive rank strategy is preferable for practical approximations and is useful when the accuracy in each sub-block is crucial.

In Fig. 5 (left) we demonstrate boxplots for different ranks $k \in \{3, 7, 9, 12\}$. In the experiment we run 100 replicates. One can see that boxplots for $k \in \{7, 9, 12\}$ are not changing, therefore there is no point to increase the rank, i.e. rank 7 is sufficient. Fig. 5 shows that the larger is the number of observations the better is estimation of the unknown parameter.

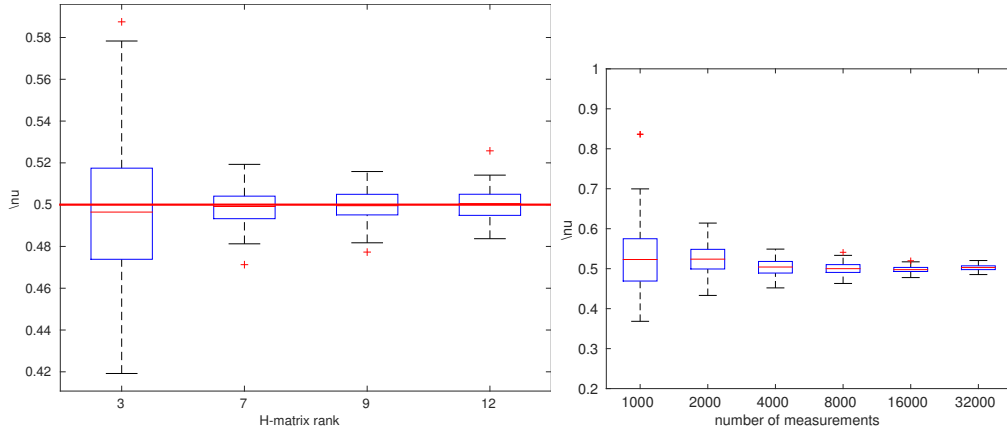


Figure 5: (left) Dependence of the boxplots for ν on the \mathcal{H} -matrix rank, $n = 16,000$; (right) Convergence of boxplots for ν with increasing n , 100 replicates.

In Fig. 6 (left) we demonstrate the shape of the negative log-likelihood, of the quadratic form and of the log-determinant $\log(|C|)$. The true value of the parameter $\theta^* = \ell = 12$. In Fig. 6 (right) one can see three log-likelihood functions - the exact one and two others computed with different ranks 7 and 17. One can see, that all three plots look very similar.

2.3 Parallel hierarchical matrix technique

We used the parallel \mathcal{H} -matrix library HLIBpro [25, 40–42] to approximate the Matérn covariance matrix, to compute a Cholesky factorization, solve a linear system, calculate the determinant and a quadratic form. HLIBpro is fast, robust and efficient; see theoretical

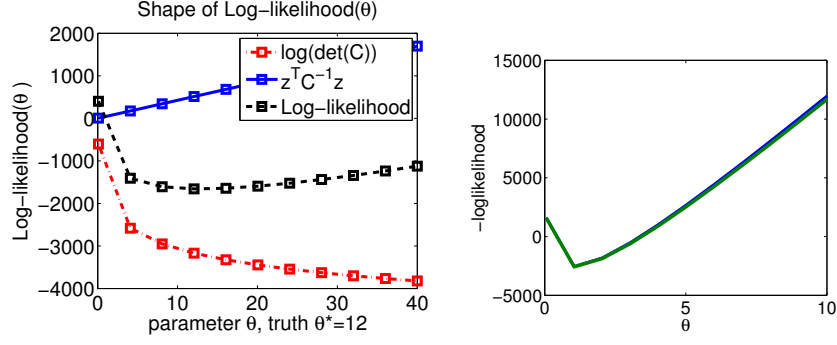


Figure 6: (left) Shape of the log-likelihood function; (right) Three negative log-likelihood functions: the exact, and computed with \mathcal{H} -matrix ranks 7 and 17. One can see that even with rank 7 one can achieve very accurate results.

215 parallel complexity in Table 1. Here $|V(T)|$ denotes the number of vertices, $|L(T)|$ the number of leaves in the block-cluster tree $T = T_I \times I$, and n_{min} the size of a block when we stop further division into sub-blocks (see Section 2.2). Usually $n_{min} = 32$ or 64 , since a very deep hierarchy slows down computations.

Table 1: Parallel complexity of the main linear operations in HLIBpro on p cores.

Operation	Parallel Complexity [40] (Shared Memory)
build \tilde{C}	$\frac{\mathcal{O}(n \log n)}{p} + \mathcal{O}(V(T) \setminus L(T))$
store \tilde{C}	$\mathcal{O}(kn \log n)$
$\tilde{C} \cdot z$	$\frac{\mathcal{O}(kn \log n)}{p} + \frac{n}{\sqrt{p}}$
$\alpha \tilde{A} \oplus \beta \tilde{B}$	$\frac{\mathcal{O}(n \log n)}{p}$
$\alpha \tilde{A} \odot \tilde{B} \oplus \beta \tilde{C}$	$\frac{\mathcal{O}(n \log n)}{p} + \mathcal{O}(C_{sp}(T) V(T))$
\tilde{C}^{-1}	$\frac{\mathcal{O}(n \log n)}{p} + \mathcal{O}(nn_{min}^2), n_{min} \approx 64$
\mathcal{H} -Cholesky \tilde{L}	$\frac{\mathcal{O}(n \log n)}{p} + \mathcal{O}(\frac{k^2 n \log^2 n}{n^{1/d}}), d = 1, 2, 3$
determinant $ \tilde{C} $	$\frac{\mathcal{O}(n \log n)}{p} + \mathcal{O}(\frac{k^2 n \log^2 n}{n^{1/d}}), d = 1, 2, 3$

3 Memory storage, computing time and convergence

The Kullback-Leibler divergence (KLD) $D_{KL}(P\|Q)$ is a measure of information lost when distribution Q is used to approximate P . For multivariate normal distributions (μ_0, \mathbf{C}) and $(\mu_1, \tilde{\mathbf{C}})$ it is defined as follow

$$D_{KL}(\mathbf{C}, \tilde{\mathbf{C}}) = 0.5 \left(\text{tr}(\tilde{\mathbf{C}}^{-1} \mathbf{C}) + (\mu_1 - \mu_0)^\top \tilde{\mathbf{C}}^{-1} (\mu_1 - \mu_0) - n - \ln \left(\frac{|\mathbf{C}|}{|\tilde{\mathbf{C}}|} \right) \right)$$

220

In Tables 2 and 3 we show dependence of KLD and two matrix errors on the \mathcal{H} -matrix rank k for Matérn covariance function with parameters $\ell = \{0.25, 0.75\}$ and $\nu = \{0.5, 1.5\}$, computed on the domain $\mathcal{G} = [0, 1]^2$. All errors are under control, besides the last column. The ranks $k = 10, 12$ are not sufficient to approximate the inverse, the error $\|\mathbf{C}(\tilde{\mathbf{C}})^{-1} - \mathbf{I}\|_2$ is large. A remedy is to take a larger rank.

k	KLD		$\ \mathbf{C} - \tilde{\mathbf{C}}\ _2$		$\ \mathbf{C}(\tilde{\mathbf{C}})^{-1} - \mathbf{I}\ _2$	
	$\ell = 0.25$	$\ell = 0.75$	$\ell = 0.25$	$\ell = 0.75$	$\ell = 0.25$	$\ell = 0.75$
10	2.6e-3	0.2	7.7e-4	7.0e-4	6.0e-2	3.1
12	5.0e-4	2e-2	9.7e-5	5.6e-5	1.6e-2	0.5
15	1.0e-5	9e-4	2.0e-5	1.1e-5	8.0e-4	0.02
20	4.5e-7	4.8e-5	6.5e-7	2.8e-7	2.1e-5	1.2e-3
50	3.4e-13	5e-12	2.0e-13	2.4e-13	4e-11	2.7e-9

Table 2: Convergence of the \mathcal{H} -matrix approximation error vs the \mathcal{H} -matrix rank k , Matérn covariance with parameters $\ell = \{0.25, 0.75\}$ and $\nu = 0.5$, domain $\mathcal{G} = [0, 1]^2$, $\|\mathbf{C}_{(\ell=0.25, 0.75)}\|_2 = \{212, 568\}$.

k	KLD		$\ C - \tilde{C}\ _2$		$\ C(\tilde{C})^{-1} - I\ _2$	
	$\ell = 0.25$	$\ell = 0.75$	$\ell = 0.25$	$\ell = 0.75$	$\ell = 0.25$	$\ell = 0.75$
20	0.12	2.7	5.3e-7	2e-7	4.5	72
30	3.2e-5	0.4	1.3e-9	5e-10	4.8e-3	20
40	6.5e-8	1e-2	1.5e-11	8e-12	7.4e-6	0.5
50	8.3e-10	3e-3	2.0e-13	1.5e-13	1.5e-7	0.1

Table 3: Convergence of the \mathcal{H} -matrix approximation error vs the \mathcal{H} -matrix rank k , Matérn covariance with parameters $\ell = \{0.25, 0.75\}$ and $\nu = 1.5$, domain $\mathcal{G} = [0, 1]^2$, $\|C_{(\ell=0.25, 0.75)}\|_2 = \{720, 1068\}$.

225

Figure 7 shows that the \mathcal{H} -matrix storage cost is almost not changing for different parameters $\ell = \{0.15, \dots, 2.2\}$ (on the left) and $\nu = \{0.3, \dots, 1.3\}$ on the right. The computational domain is $[32.4, 43.4] \times [-84.8 - 72.9]$ with $n = 2,000$.

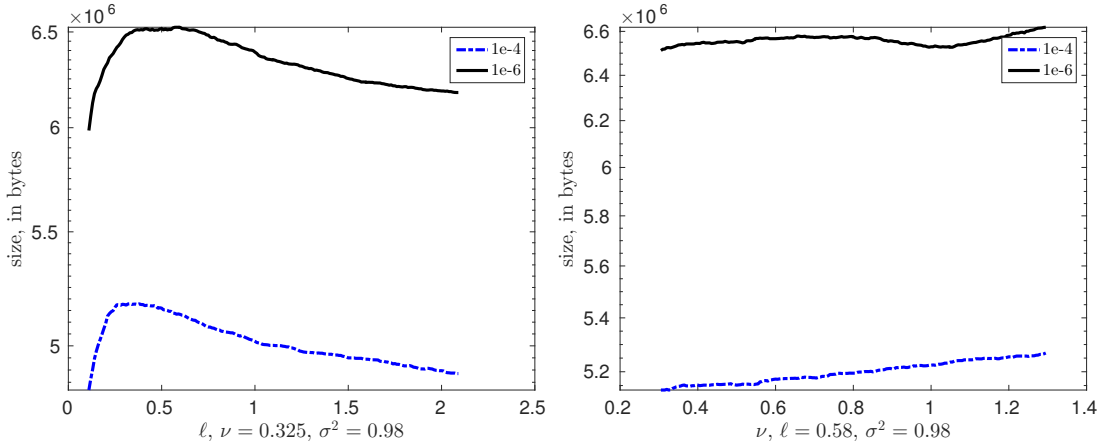


Figure 7: (left) Dependence of the matrix size on (left) the covariance length ℓ ; (right) the smoothness ν for two different accuracies in \mathcal{H} -matrix sub-blocks $\varepsilon = \{1e-4, 1e-6\}$, $n = 2,000$ locations in the domain $[32.4, 43.4] \times [-84.8 - 72.9]$.

230

In Fig. 8 we plotted convergence of $\|C - \tilde{C}\|$ in the Frobenius and spectral norms vs. the rank k for different smoothness parameters and covariance lengths. In Fig. 9 we plotted $\|C - \tilde{C}\|_2$ vs. rank k for different ν and covariance lengths.

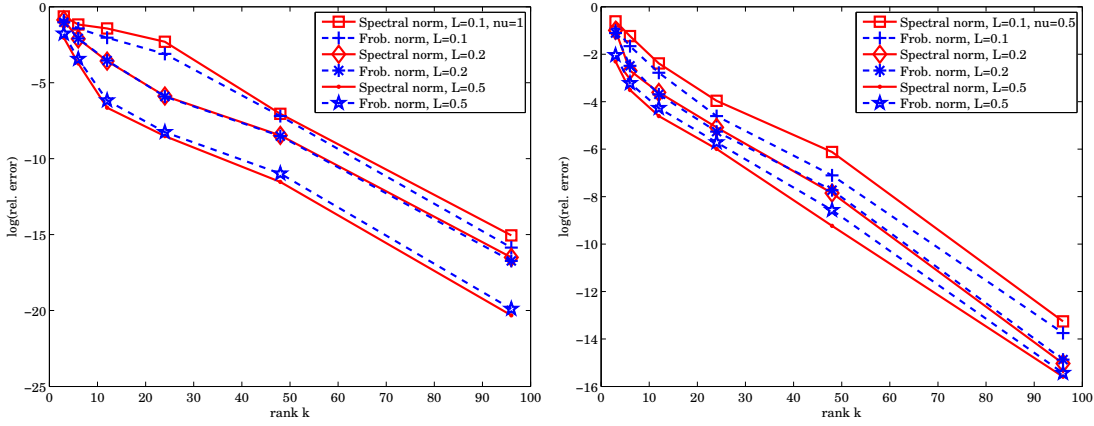


Figure 8: Convergence of the \mathcal{H} -matrix approximation errors for covariance lengths $\{0.1, 0.2, 0.5\}$; (left) $\nu = 1$ and (right) $\nu = 0.5$.

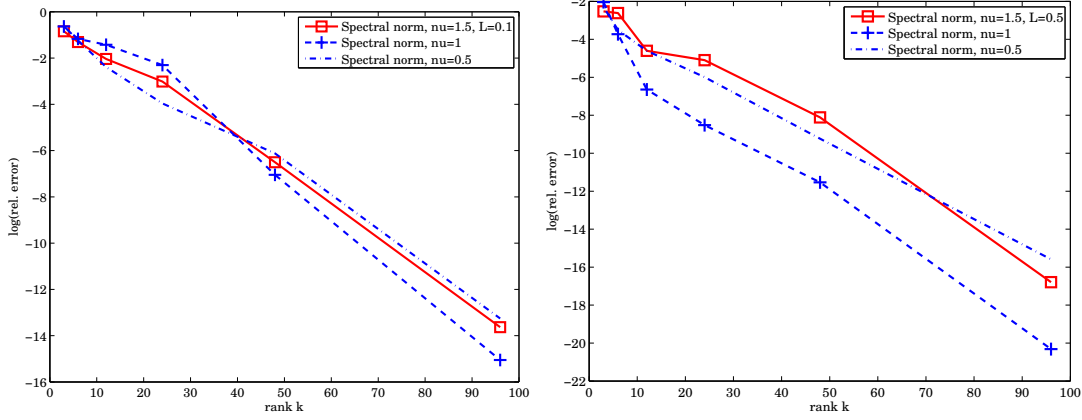


Figure 9: Convergence of the \mathcal{H} -matrix approximation errors for $\nu = \{0.5, 1, 1.5\}$; (left) covariance length 0.1 and (right) covariance length 0.5.

4 Software Installation and Numerical Examples

4.1 Introduction to HLIBpro

This Section contains shortened information taken from www.hlib.com. The software library HLIBpro is a robust parallel C++ implementation of the \mathcal{H} -matrix technique and \mathcal{H} -matrix arithmetics, developed by Ronald Kriemann [40]. HLIBpro supports both the shared and distributed memory architectures, but in this work we use only the shared

memory version. Threading Building Blocks are used for parallelisation. HLIBpro is free for academic purposes, is distributed as a compiled code (no source code is available). Originally HLIBpro was developed for solving FEM and BEM problems [25, 42].
 240 In this work we extend the applicability of HLIBpro to dense covariance matrices and log-likelihood functions.

4.2 Installation of HLIBpro

HLIBCov uses functionality of HLIBpro, therefore, first, the reader needs to install HLIBpro.
Hardware. All of the numerical experiments herein are performed on a Dell workstation
 245 with 20 processors (40 cores) and total 128 GB RAM.

Operation system. All the following steps are described for Linux systems (in our case Ubuntu 14.04). We did not tested it for MacOS and for Windows operation systems. Since HLIBcov is just few c++ procedures (only one file loglikelihood.cc), the applicability of HLIBCov is limited only by the applicability of HLIBpro. To install HLIBpro on MacOS
 250 and on Windows we refer to www.HLIBpro.com for further details.

Software. To use the binary distribution of HLIBpro, the following libraries have to be available together with HLIBpro: LAPACK/BLAS, Threading Building Blocks, Boost, zlib, FFTW3, Scotch, METIS, SCons (see details on www.HLIBpro.com). Additionally, the HLIBCov requires the GNU Scientific Library (GSL). We use GSL-1.16, the installation
 255 steps are described on <https://www.gnu.org/software/gsl/>. The GSL library provides maximization algorithms and Bessel functions. The reader can easily replace GSL on his own optimization library. The Bessel functions are also available in other packages.

To install LAPACK, zlib, FFTW3, Scotch and SCons using the standard package tools, you may run

```
260 1 $ aptitude install liblapack-dev zlib1g-dev libfftw3-dev libscotch-dev
    2 $ aptitude install scons
```

A manual installation of TBB, Boost, Scotch and METIS requires appropriate modification
 265 of variables `tbb_dir`, `boost_dir`, `scotch_dir` and `metis_dir` in the `SConstruct` file. After that, compile everything within the HLIBpro installation directory with

```
1 $ scons
270 $ export LD_LIBRARY_PATH=$LD_LIBRARY_PATH:<path_to_HLIBpro>/lib
```

The last terminal command is usually not necessary when using the provided SCons makefiles.

An alternative - Using the Auxiliary Library Package. As an alternative to the manual

installing of all required packages (listed above), the HLIBpro provides a *tar.gz* file with all auxiliary library packages. Download and unpack it into your HLIBpro directory:

```
1 $ cd HLIBpro-2.6
2 $ tar -xzf HLIBpro-2.6-Linux-aux.tgz
```

Create a copy of *SConstruct* file, since the original *SConstruct* file will be overwritten. Afterwards, by calling *scons* you may compile the HLIBpro examples using the provided libraries.

Adding HLIBCov to HLIBpro. For default settings/paths, there is no need to change *SConstruct* file. To compile HLIBCov, add the following line to */examples/SConscript*

```
285 $ examples.append(cxenv.Program('loglikelihood.cc'))
1 $ cd ..
2 $ scons
3 $ cd examples/
4 $ ./loglikelihood
290
```

Input of the HLIBCov

The input contains in the first line the total number of locations N . Lines $2, \dots, N+1$ contain the coordinates x_i, y_i , and the measurement value. An example

```
295 3
1 0.1 0.2 88.1
2 0.1 0.3 87.2
3 0.2 0.4 86.0
300
```

The HLIBpro does not require neither a list of finite elements nor a list of edges. On the github we provide few input files examples of different size.

Output of the HLIBCov

The main output is the three identified parameter values $\theta = (\ell, \nu, \sigma^2)$. The auxiliary output may include a lot of details: \mathcal{H} -matrix details (the maximal rank k , the maximal accuracy in each sub-block, Frobenius and spectral norms of $\tilde{\mathbf{C}}, \tilde{\mathbf{L}}, \tilde{\mathbf{L}}^{-1}, \|\mathbf{I} - \tilde{\mathbf{L}}\tilde{\mathbf{L}}^\top^{-1}\|$). Additionally, you can also print out iterations of the maximization algorithm. An example an output file below contains two iterations: index, $\nu, \ell, \sigma^2, \tilde{\mathcal{L}}$, and the residual of the iterative method

```
310 1 0.27 2.4 1.30 L = 1762.1 TOL= 0.007
1 2 0.276 2.41 1.29 L = 1757.2 TOL= 0.009
2
```

If the iterative process is converging, then the last row will contain the solution $\theta^* = (\ell^*, \nu^*, \sigma^{*2})$. While computing error boxes, the output file will contain M solutions $(n, \ell^*, \nu^*, \sigma^{*2})$, where M is the number of replicants:

```
1 4000 5.4e-1 8.2e-2 1.01
```

```

2    4000 5.3e-1 8.3e-2 1.02
320  4000 5.5e-1 8.1e-2 1.02

```

4.3 The preprocessing C++ function

Below we list the C++ code, which reads the data and initializes all required C++ objects (*loglikelihood.cc*).

```

325 1      vector< double * > vertices;
2      TScalarVector    rhs;
3      INIT();
4      in >> N;
330 cout << " reading " << N << " coordinates" << endl;
6      vertices.resize( N );
7      rhs.set_size( N );
8      double x, y, v=0.0;
9      for ( int i = 0; i < N; ++i ){
335         in >> x >> y >> v;
11         vertices[i] = new double[ dim ];
12         vertices[i][0] = x; vertices[i][1] = y;
13         rhs.set_entry( i, v ); } // for
14 TCoordinate coord( vertices, dim );
340 TAutoBSPPartStrat part_strat( adaptive_split_axis );
16 TBSPCTBuilder      ct_builder( & part_strat );
17 auto               ct = ct_builder.build( & coord );
18 TStdGeomAdmCond    adm_cond( 2.0, use_min_diam );
19 TBCTBuilder        bct_builder( std::log2( 16 ) );
345 auto              bct = bct_builder.build( ct.get(), ct.get(), & adm_cond );
21 // bring RHS into H-ordering
22 ct->perm_e2i()->permute( & rhs );
23 double output[3];
24 call_compute_max_likelihood(rhs, nu, length, sigma2, bct.get(), ct.get(), vertices, output);
350 DONE();

```

The C++ code computing the maximum of the log-likelihood function (*loglikelihood.cc*).

```

1 double call_compute_max_likelihood(TScalarVector Z, double nu, double covlength, double sigma2, TBlockClusterTree* bct,
325 TClusterTree* ct, std::vector<double*> vertices, double output[3])
2 { gsl_function F;
3   int status; iter = 0, max_iter = 200; smy_f_params params ;
4   FILE* f1; double size;
5   const gsl_multimin_fminimizer_type *T = gsl_multimin_fminimizer_nmsimplex2;
360 gsl_multimin_fminimizer *s = NULL; gsl_vector *ss, *x;
7   gsl_multimin_function minex_func;
8   params.bct = bct; params.ct = ct; params.Z = Z; params.nu = nu;
9   params.covlength=covlength; params.sigma2=sigma2; params.vertices=vertices;
10  /* Starting point */
365 x = gsl_vector_alloc(3); gsl_vector_set (x, 0, nu);
12  gsl_vector_set (x, 1, covlength); gsl_vector_set (x, 2, sigma2);
13  /* Set initial step sizes to 0.1 */
14  ss = gsl_vector_alloc (3);
15  gsl_vector_set (ss, 0, 0.02); //for nu
370  gsl_vector_set (ss, 1, 0.04); //for theta
17  gsl_vector_set (ss, 2, 0.01); //for sigma2
18  /* Initialize method and iterate */
19  minex_func.n = 3; //dimension
20  minex_func.f = &eval_logli;
375 minex_func.params = &params;
22  s = gsl_multimin_fminimizer_alloc (T, 3); /* or 2? like in the example */
23  gsl_multimin_fminimizer_set (s, &minex_func, x, ss);
24  do{ iter++;

```

```

25     status = gsl_multimin_fminimizer_iterate(s);
380    if (status) break;
27     size = gsl_multimin_fminimizer_size (s); //for stopping criteria
28     status = gsl_multimin_test_size (size, 1e-5);
29     if (status == GSL_SUCCESS) printf ("converged to minimum at \n");}}
30 while (status == GSL_CONTINUE && iter < max_iter);
385 output[0]= gsl_vector_get(s->x, 0); //nu
32 output[1]= gsl_vector_get(s->x, 1); //theta
33 output[2]= gsl_vector_get(s->x, 2); //sigma2
34 gsl_vector_free(x); gsl_vector_free(ss); gsl_multimin_fminimizer_free (s);
390 return status; }

```

Below we list the C++ code, which computes the value of the log-likelihood for given parameters (*loglikelihood.cc*).

```

1 double eval_logli (const gsl_vector *sol, void* p)
395 { pmy_f_params params ;
3     double nu = gsl_vector_get(sol, 0);
4     double length = gsl_vector_get(sol, 1);
5     double sigma2 = gsl_vector_get(sol, 2);
6     unique_ptr< TProgressBar > progress( verbose(2) ? new TConsoleProgressBar : nullptr );
400     params = (pmy_f_params)p;
8     TScalarVector rhs= (params->Z);
9     TBlockClusterTree* bct = (params->bct); TClusterTree* ct = (params->ct);
10    vector< double * > vertices= (params->vertices);
11    double err2=0.0, nugget = 1.0e-4, s = 0.0;
405    auto acc = fixed_prec( 1e-5 ); int dim = 2, N = 0;
13    TCovCoeffFn coefffn(length,nu,sigma2,nugget,vertices,ct->perm_i2e(),ct->perm_i2e());
14    TACAPlus< real_t > aca( & coefffn );
15    TDenseMatBuilder< real_t > h_builder( & coefffn, & aca );
16    // enable coarsening during construction
410    h_builder.set_coarsening( false );
18    auto A = h_builder.build( bct, acc, progress.get() );
19    N=A->cols();
20    auto A_copy = A->copy();
21    auto options = fac_options_t( progress.get() );
415    options.eval = point_wise; //! Extreme important
23    auto A_inv = ldl_inv( A_copy.get(), acc, options );
24    for ( int i = 0; i < N; ++i ) {
25        const auto v = A_copy->entry( i, i );
26        s = s + log(v);} // for
420    TStopCriterion sstop( 150, 1e-6, 0.0 );
28    TCG solver( sstop );
29    TSolverInfo sinfo( false, verbose( 4 ) );
30    auto solu = A->row_vector();
31    solver.solve( A.get(), solu.get(), & rhs, A_inv.get(), & sinfo );
425    auto dotp = re( rhs.dot( solu.get() ) );
33    auto LL = 0.5*N*log(2*Math::pi<double>())+0.5*s+0.5*dotp;}

```

HLIBpro License

To receive the license file **HLIBpro.lic** send an email with your user login name and the machine name, on which you plan to run the code, to Ronald Kriemann rok@mis.mpg.de and request a license file. The machine name can be received with the terminal command *hostname*. The received license file must be placed either in the directory where the final program will be executed, e.g. in the examples/ subdirectory, or in the directory pointed to by the environment variable HLIBpro_LIC_PATH. You can add this variable to the system paths

```
1 $ export HLIBpro_LIC_PATH=<path_to_HLIBpro.lic>
```

5 Numerical experiments with synthetic data

440 We performed numerical experiments with simulated data to recover the true values of the parameters of the Matérn covariance matrix, known to be $(\ell^*, \nu^*, \sigma^*) = (1.0, 0.5, 1.0)$.

5.1 Generation of the synthetic data

To build M various data sets (M replicates) with $n \in \{128, 64, \dots, 4, 2\} \times 1000$ locations, we generate a large vector \mathbf{Z}_0 with $n_0 = 2 \cdot 10^6$ locations and randomly sample n points from it.
 445 it. Note, that if locations are very close to each other, then the covariance matrix could be singular or it will be difficult to compute the Cholesky factorization.

To generate the random data $\mathbf{Z}_0 \in \mathbb{R}^{n_0}$ we compute the \mathcal{H} -Cholesky factorization of $\mathbf{C}(1.0, 0.5, 1.0) = \tilde{\mathbf{L}}\tilde{\mathbf{L}}^\top$. Then we evaluate $\mathbf{Z}_0 = \tilde{\mathbf{L}}\boldsymbol{\xi}$, where $\boldsymbol{\xi} \in \mathbb{R}^{n_0}$ is a normal vector with zero mean and unit variance. We generate \mathbf{Z}_0 only once. After that we run our optimization algorithm and try to identify (recover) the “unknown” parameters $(\ell^*, \nu^*, \sigma^*)$. The resulting boxplots for ν and σ^2 over $M = 100$ replicates are illustrated in Fig. 10.

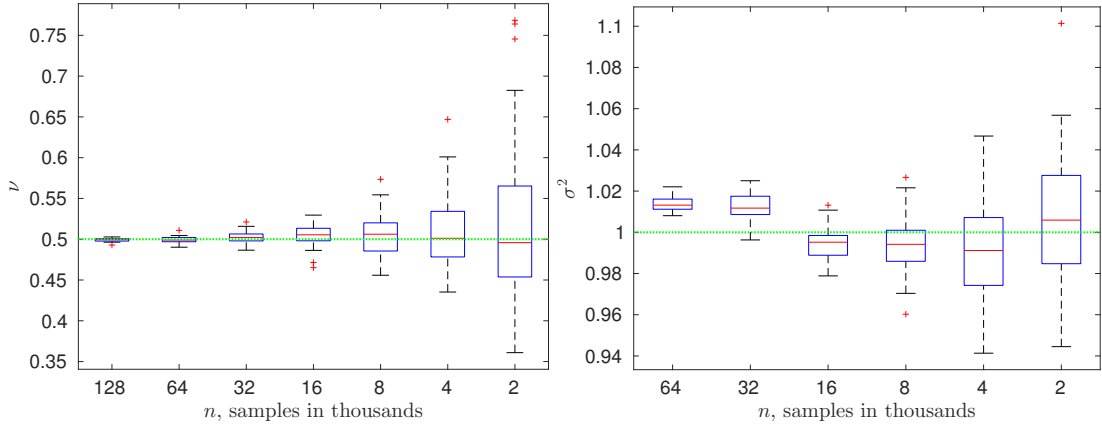


Figure 10: Synthetic data with known parameters $(\ell^*, \nu^*, \sigma^{2*}) = (0.5, 1, 0.5)$. Boxplots for ν, σ^2 for $n = 1,000 \times \{64, 32, \dots, 4, 2\}$, 100 replicates.

To identify all three parameters simultaneously a three-dimensional optimization problem is solved. The maximal number of iterations is set to 200, the residual to 10^{-6} . The

behavior and accuracy of the boxplots depend on the hierarchical matrix rank, the maximum number of iterations to achieve a certain threshold, the threshold (or residual) itself, the initial guess, the step size in each parameter in the maximization algorithm, and the maximization algorithm. All replicates of \mathbf{Z} are sampled from the same generated vector of size $n_0 = 2 \cdot 10^6$.

In Table 4 we demonstrate the almost linear storage cost (columns 3 and 6) and computing time (columns 2 and 5).

Table 4: Computing time and storage cost for parallel \mathcal{H} -matrix approximation; number of cores is 40, $\hat{\nu} = 0.33$, $\hat{\ell} = 0.65$, $\hat{\sigma}^2 = 1.0$. \mathcal{H} -matrix accuracy in each sub-block for both $\tilde{\mathbf{C}}$ and $\tilde{\mathbf{L}}$ is 10^{-5} .

n	$\tilde{\mathbf{C}}$			$\tilde{\mathbf{L}}\tilde{\mathbf{L}}^\top$		
	comp. time sec.	size MB	kB/dof	comp. time sec.	size MB	$\ \mathbf{I} - (\tilde{\mathbf{L}}\tilde{\mathbf{L}}^\top)^{-1}\tilde{\mathbf{C}}\ _2$
32,000	3.3	162	5.1	2.4	172.7	$2.4 \cdot 10^{-3}$
128,000	13.3	776	6.1	13.9	881.2	$1.1 \cdot 10^{-2}$
512,000	52.8	3420	6.7	77.6	4150	$3.5 \cdot 10^{-2}$
2,000,000	229	14790	7.4	473	18970	$1.4 \cdot 10^{-1}$

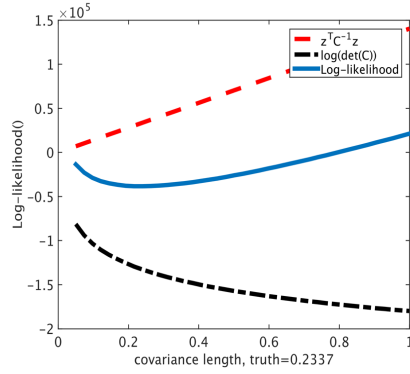


Figure 11: Dependence of the negative log-likelihood and its ingredients on parameter ℓ , parameter $\nu = 0.5$ is fixed, the true value of covariance length is 0.2337. Exponential covariance, $n = 66049$, rank $k = 16$, $\sigma^2 = 1$.

To get an idea about the shape of the log-likelihood function, and about its components we draw Fig. 11. This picture helps us to understand the behavior of the iterative

465 optimization method, contribution of the log-determinant, of the quadratic functional. One can see that the log-likelihood is almost flat and it may be necessary a lot of iteration steps to find the minimum.

Table 5: Comparison of three log-likelihood functions computed with three different \mathcal{H} -matrix accuracies $\{10^{-7}, 10^{-9}, 10^{-11}\}$. Exponential covariance function discretized in the domain $[32.4, 43.4] \times [-84.8 - 72.9]$, $n=32,000$ locations. Columns corresponds to different covariance lengths $\{0.001, \dots, 0.1\}$.

ℓ	0.001	0.005	0.01	0.02	0.03	0.05	0.07	0.1
$-\tilde{\mathcal{L}}(\ell; 10^{-7})$	44657	36157	36427	40522	45398	68450	70467	90649
$-\tilde{\mathcal{L}}(\ell; 10^{-9})$	44585	36352	36113	41748	47443	60286	70688	90615
$-\tilde{\mathcal{L}}(\ell; 10^{-11})$	44529	37655	36390	42020	47954	60371	72785	90639

Nugget τ^2 . Assume $|\tilde{\mathbf{C}}| \neq 0$. It is known [18], that for a small perturbation matrix \mathbf{E} it is hold

$$\frac{\|(\tilde{\mathbf{C}} + \mathbf{E})^{-1} - (\tilde{\mathbf{C}})^{-1}\|}{\|\tilde{\mathbf{C}}^{-1}\|} \leq \kappa(\mathbf{C}) \cdot \frac{\|\mathbf{E}\|}{\|\tilde{\mathbf{C}}\|} = \frac{\kappa(\tilde{\mathbf{C}})\tau^2}{\|\tilde{\mathbf{C}}\|},$$

where $\kappa(\tilde{\mathbf{C}})$ is the condition number of $\tilde{\mathbf{C}}$, and $\mathbf{E} = \tau^2 \mathbf{I}$. Or, substituting $\kappa(\tilde{\mathbf{C}}) := \|\tilde{\mathbf{C}}^{-1}\| \cdot \|\tilde{\mathbf{C}}\|$, obtain

$$\frac{\|(\tilde{\mathbf{C}} + \tau^2 \mathbf{I})^{-1} - (\tilde{\mathbf{C}})^{-1}\|}{\|\tilde{\mathbf{C}}^{-1}\|} \leq \tau^2 \|\tilde{\mathbf{C}}^{-1}\|.$$

From the last equation one can see that if $a \leq \|\tilde{\mathbf{C}}^{-1}\| \leq b$, where a, b are some positive constants, then the relative error will be of order τ^2 . Figure 12 (left) demonstrates three
 470 negative log-likelihood functions. One is computed with the nugget 0.01, another with 0.005 and the third with 0.001. As one can see, for this particular example, the behavior of likelihood is preserved, and the minimum is not changing (or changing very slightly). Figure 12 (right) is just a zoomed version of the left picture.

6 Best practices

475 In this section we list our recommendations and warnings.

1. Initialize variables and the license with `INIT()` ;
2. For practical computations use the adaptive rank arithmetics.

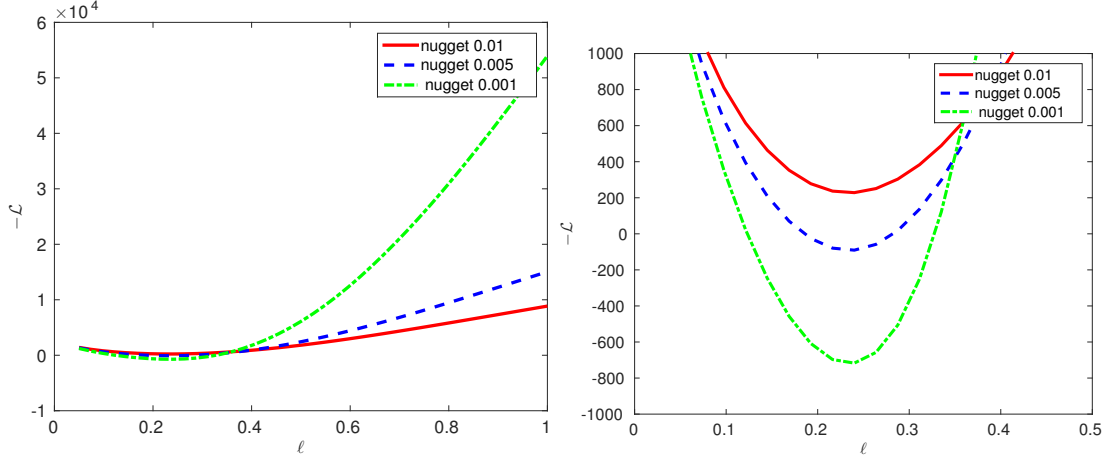


Figure 12: (left) Dependence of the log-likelihood on parameter ℓ with nuggets $(\{0.01, 0.005, 0.001\})$ for Gaussian covariance. (right) Zoom of the left figure near minimum; $n = 2000$ random points from moisture example, rank $k = 14$, $\sigma^2 = 1$.

3. For the input it is enough to define a file with 3 columns: location coordinates (x, y) and the observed value, no triangles, no edges are required;
- 480 4. If two locations are coincide or very close to each other, then the matrix will be close to singular or singular. As a result, it will be hard to compute the Cholesky factorization. A remedy is to improve the quality of locations.
5. Construction of \mathcal{H} -matrices permutes the indices of the degrees of freedom. Therefore, there are two permutation mappings: `ct->perm_e2i()` and `ct->perm_i2e()`.
- 485 For example, to prepare the observation vector \mathbf{Z} , computed on a grid, for multiplication with the hierarchical matrix $\tilde{\mathbf{L}}^{-1}$ (or $\tilde{\mathbf{C}}^{-1}$) use `ct->perm_e2i()->permute(&Z)`. To find to which grid node corresponds the i -th row (column), use `ct->perm_i2e()`.
6. HLIBpro uses smart pointers, e.g., `std::unique_ptr`. If the user leaves a code block those smart pointers will automatically delete the object they manage. This simplifies programming and makes it much more secure.
- 490 7. By default, the \mathcal{H} -Cholesky or \mathcal{H} -LU factorizations are always DAG based. You can turn this off by setting `HLIB::CFG::Arith::use_dag = false`;
8. By default, HLIBpro uses all available computing cores. To perform computations on 16 cores, use `HLIB::CFG::set_nthreads(16)` at the beginning of the program

495 (after command `INIT()`).

9. Since HLIBpro is working for 2D and 3D, there are only very minor changes in HLIBCov to move from 2D locations to 3D. Replace `dim=2` to `dim=3` in

```
TCoordinate coord(vertices, dim);
Then add >> z to in >> x >> y>> z >> v;
```

- 500 10. Finalize all computations with `DONE()`;

7 Conclusion

We extended functionality of the parallel \mathcal{H} -matrix library HLIBpro for application in spatial statistics and in parameters inference. This new extension allows us to work with large covariance matrices. The first novelty is that we approximated the joint multivariate Gaussian likelihood function and have found its maxima in the \mathcal{H} -matrix format. This maxima estimated unknown parameters (ℓ , ν , and σ^2) of the covariance model. The second novelty is that the new code is parallel, highly efficient and written in C++ language. With the \mathcal{H} -matrix technique we reduced the storage cost and the computing cost (Tables 3, 4) of the log-likelihood function dramatically, from cubic to almost linear. We demonstrated a synthetic example, where we were able to identify the true parameters of the covariance model. With the suggested method we were able to compute the log-likelihood function for 2,000,000 locations just in a few minutes on a powerful desktop machine (Table 4). The \mathcal{H} -matrix technique allows us to increase the spatial resolution, to handle more measurements, consider larger regions and identify more parameters simultaneously.

Acknowledgments

The author would like to thank Ronald Kriemann from the Max Planck Institute for Mathematics in the Sciences in Leipzig for his assistance in the HLIBpro code. This work would be impossible without assistance from Marc Genton, Ying Sun and David Keyes. The research reported in this publication was supported by funding from King Abdullah University of Science and Technology (KAUST).

References

- [1] M. Abramowitz and I. A. Stegun. *Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables*. Number 55 in National Bureau of Standards Applied Mathematics Series. Superintendent of Documents, U.S. Government Printing Office, Washington, DC, 1964.
- [2] S. Ambikasaran, D. Foreman-Mackey, L. Greengard, D. W. Hogg, and M. O’Neil. Fast direct methods for Gaussian processes and the analysis of NASA Kepler mission data. *arXiv preprint arXiv:1403.6015*, 2014.
- [3] S. Ambikasaran, J. Y. Li, P. K. Kitanidis, and E. Darve. Large-scale stochastic linear inversion using hierarchical matrices. *Computational Geosciences*, 17(6):913–927, 2013.
- [4] J. Ballani and D. Kressner. Sparse inverse covariance estimation with hierarchical matrices. http://sma.epfl.ch/~anchpcommon/publications/quic_ballani_kressner_2014.pdf, 2015.
- [5] S. Banerjee, A. E. Gelfand, A. O. Finley, and H. Sang. Gaussian predictive process models for large spatial data sets. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 70(4):825–848, 2008.
- [6] M. Bebendorf. Approximation of boundary element matrices. *Numerical Mathematics*, 86(4):565–589, 2000.
- [7] M. Bebendorf. Why approximate LU decompositions of finite element discretizations of elliptic operators can be computed with almost linear complexity. *SIAM J. Numerical Analysis*, 45:1472–1494, 2007.
- [8] M. Bebendorf and T. Fischer. On the purely algebraic data-sparse approximation of the inverse and the triangular factors of sparse matrices. *Numerical Linear Algebra with Applications*, 18(1):105–122, 2011.
- [9] M. Bebendorf and W. Hackbusch. Existence of H-matrix approximants to the inverse FE-matrix of elliptic operators with L^∞ -coefficients. *Numerische Mathematik*, 95(1):1–28, 2003.
- [10] M. Bebendorf and S. Rjasanow. Adaptive low-rank approximation of collocation matrices. *Computing*, 70(1):1–24, 2003.
- [11] S. Börm and J. Garcke. Approximating Gaussian processes with H^2 -matrices. In J. N. Kok, J. Koronacki, R. L. de Mantaras, S. Matwin, D. Mladen, and A. Skowron, editors, *Proceedings of 18th European Conference on Machine Learning, Warsaw, Poland, September 17-21, 2007. ECML 2007*, volume 4701, pages 42–53, 2007.
- [12] S. Börm and L. Grasedyck. Hybrid cross approximation of integral operators. *Numer. Math.*, 101(2):221–249, 2005.
- [13] S. Börm, L. Grasedyck, and W. Hackbusch. *Hierarchical Matrices*, volume 21 of *Lecture Note*. Max-Planck Institute for Mathematics, Leipzig, 2003. www.mis.mpg.de.
- [14] R. P. Brent. Chapter 4: An algorithm with guaranteed convergence for finding a zero of a function, algorithms for minimization without derivatives. *Englewood Cliffs, NJ: Prentice-Hall*, 1973.
- [15] J. E. Castrillon, M. G. Genton, and R. Yokota. Multi-Level Restricted Maximum Likelihood

- Covariance Estimation and Kriging for Large Non-Gridded Spatial Datasets. *Spatial Statistics*, 18:105–124, 2016.
- [16] N. Cressie and G. Johannesson. Fixed rank kriging for very large spatial data sets. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 70(1):209–226, 2008.
- 565 [17] A. Datta, S. Banerjee, A. O. Finley, and A. E. Gelfand. Hierarchical nearest-neighbor Gaussian process models for large geostatistical datasets. *Journal of the American Statistical Association*, 0(ja):00–00, 2015.
- [18] J. Demmel. The componentwise distance to the nearest singular matrix. *SIAM Journal on Matrix Analysis and Applications*, 13(1):10–19, 1992.
- 570 [19] M. Fuentes. Approximate likelihood for large irregularly spaced spatial data. *Journal of the American Statistical Association*, 102:321–331, 2007.
- [20] G.-A. Fuglstad, D. Simpson, F. Lindgren, and H. Rue. Does non-stationary spatial data always require non-stationary random fields? *Spatial Statistics*, 14:505–531, 2015.
- [21] R. Furrer, M. G. Genton, and D. Nychka. Covariance tapering for interpolation of large spatial datasets. *Journal of Computational and Graphical Statistics*, 15(3):502–523, 2006.
- 575 [22] G. H. Golub and C. F. Van Loan. Matrix computations, 2012.
- [23] S. A. Goreinov, E. E. Tyrtyshnikov, and N. L. Zamarashkin. A theory of pseudoskeleton approximations. *Linear Algebra Appl.*, 261:1–21, 1997.
- [24] L. Grasedyck and W. Hackbusch. Construction and arithmetics of \mathcal{H} -matrices. *Computing*, 580 70(4):295–334, 2003.
- [25] L. Grasedyck, R. Kriemann, and S. LeBorne. Parallel black box H-LU preconditioning for elliptic boundary value problems. *Computing and visualization in science*, 11(4-6):273–291, 2008.
- [26] P. Guttorp and T. Gneiting. Studies in the history of probability and statistics XLIX: On the Matérn correlation family. *Biometrika*, 93:989–995, 2006.
- 585 [27] W. Hackbusch. A sparse matrix arithmetic based on \mathcal{H} -matrices. I. Introduction to \mathcal{H} -matrices. *Computing*, 62(2):89–108, 1999.
- [28] W. Hackbusch. *Hierarchical matrices: Algorithms and Analysis*, volume 49 of *Springer Series in Comp. Math.* Springer, 2015.
- 590 [29] W. Hackbusch and B. N. Khoromskij. A sparse \mathcal{H} -matrix arithmetic. II. Application to multi-dimensional problems. *Computing*, 64(1):21–47, 2000.
- [30] W. Hackbusch, B. N. Khoromskij, and R. Kriemann. Hierarchical matrices based on a weak admissibility criterion. *Computing*, 73(3):207–243, 2004.
- [31] M. S. Handcock and M. L. Stein. A Bayesian analysis of kriging. *Technometrics*, 35:403–410, 595 1993.
- [32] H. Harbrecht, M. Peters, and M. Siebenmorgen. Efficient approximation of random fields for numerical applications. *Numerical Linear Algebra with Applications*, 22(4):596–617, 2015.
- [33] K. L. Ho and L. Ying. Hierarchical interpolative factorization for elliptic operators: differential equations. *Communications on Pure and Applied Mathematics*, 2015.

- [34] C.-J. Hsieh, M. A. Sustik, I. S. Dhillon, P. K. Ravikumar, and R. Poldrack. Big & QUIC: Sparse inverse covariance estimation for a million variables. In C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Weinberger, editors, *Advances in Neural Information Processing Systems 26*, pages 3165–3173. Curran Associates, Inc., 2013.
- [35] H. Huang and Y. Sun. Hierarchical low rank approximation of likelihoods for large spatial datasets. *Journal of Computational and Graphical Statistics*, to appear. *ArXiv e-prints* 1605.08898, 2017.
- [36] C. G. Kaufman, M. J. Schervish, and D. W. Nychka. Covariance tapering for likelihood-based estimation in large spatial datasets. *Journal of the American Statistical Association*, 103(484):1545–1555, 2008.
- [37] B. N. Khoromskij and A. Litvinenko. Data sparse computation of the Karhunen-Loève expansion. *AIP Conference Proceedings*, 1048(1):311, 2008.
- [38] B. N. Khoromskij, A. Litvinenko, and H. G. Matthies. Application of hierarchical matrices for computing the Karhunen-Loève expansion. *Computing*, 84(1-2):49–67, 2009.
- [39] W. Kleiber and D. W. Nychka. Equivalent kriging. *Spatial Statistics*, 12:31–49, 2015.
- [40] R. Kriemann. Parallel H-matrix arithmetics on shared memory systems. *Computing*, 74(3):273–297, 2005.
- [41] R. Kriemann. HLIBpro user manual. Technical report, Max Planck Institute for Mathematics in the Sciences, 2008.
- [42] R. Kriemann. H-LU factorization on many-core systems. Preprint, Max Planck Institute for Mathematics in the Sciences, 2014.
- [43] 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(4):423–498, 2011.
- [44] A. Litvinenko. Application of hierarchical matrices for solving multiscale problems. *PhD Dissertation, Leipzig University*, 2006.
- [45] A. Litvinenko. Partial inversion of elliptic operator to speed up computation of likelihood in bayesian inference. *arXiv preprint arXiv:1708.02207*, 2017.
- [46] A. Litvinenko and H. G. Matthies. Sparse data representation of random fields. *PAMM*, 9(1):587–588, 2009.
- [47] A. Litvinenko and H. G. Matthies. Inverse problems and uncertainty quantification. *arXiv preprint arXiv:1312.5048*, 2013, 2013.
- [48] P.-G. Martinsson and V. Rokhlin. A fast direct solver for boundary integral equations in two dimensions. *Journal of Computational Physics*, 205(1):1–23, 2005.
- [49] B. Matérn. *Spatial Variation*, volume 36 of *Lecture Notes in Statistics*. Springer-Verlag, Berlin; New York, second edition edition, 1986.
- [50] H. Matthies, A. Litvinenko, O. Pajonk, B. V. Rosić, and E. Zander. Parametric and uncertainty computations with tensor product representations. In A. M. Dienstfrey and R. F. Boisvert, editors, *Uncertainty Quantification in Scientific Computing*, volume 377 of *IFIP Advances in*

- Information and Communication Technology*, pages 139–150. Springer Berlin Heidelberg, 2012.
- 640 [51] W. Nowak and A. Litvinenko. Kriging and spatial design accelerated by orders of magnitude: combining low-rank covariance approximations with FFT-techniques. *Mathematical Geosciences*, 45(4):411–435, 2013.
- [52] D. Nychka, S. Bandyopadhyay, D. Hammerling, F. Lindgren, and S. Sain. A multiresolution Gaussian process model for the analysis of large spatial datasets. *Journal of Computational and Graphical Statistics*, 24(2):579–599, 2015.
- 645 [53] O. Pajonk, B. V. Rosić, A. Litvinenko, and H. G. Matthies. A deterministic filter for non-Gaussian Bayesian estimation — applications to dynamical system estimation with noisy measurements. *Physica D: Nonlinear Phenomena*, 241(7):775–788, 2012.
- [54] W. H. Press, S. A. Teukolsky, W. T. Vetterling, and B. P. Flannery. Section 9.3. *Van Wijngaarden-Dekker-Brent Method. Numerical Recipes: The Art of Scientific Computing*, volume 3rd ed. New York: Cambridge University Press., 2007.
- 650 [55] C. E. Rasmussen and C. K. I. Williams. *Gaussian Processes for Machine Learning: www.gaussianprocess.org/gpml/chapters/*, volume 497. MIT Press, 2006.
- [56] B. V. Rosić, A. Kučerová, J. Šykora, O. Pajonk, A. Litvinenko, and H. G. Matthies. Parameter identification in a probabilistic setting. *Engineering Structures*, 50:179–196, 2013.
- 655 [57] B. V. Rosić, A. Litvinenko, O. Pajonk, and H. G. Matthies. Sampling-free linear bayesian update of polynomial chaos representations. *Journal of Computational Physics*, 231(17):5761–5787, 2012.
- [58] H. Rue and L. Held. *Gaussian Markov random fields: theory and applications*, 2005.
- 660 [59] H. Rue and H. Tjelmeland. Fitting Gaussian Markov random fields to Gaussian fields. *Scandinavian Journal of Statistics*, 29(1):31–49, 2002.
- [60] A. Saibaba, S. Ambikasaran, J. Yue Li, P. Kitanidis, and E. Darve. Application of hierarchical matrices to linear inverse problems in geostatistics. *Oil & Gas Science and Technology—Rev. IFP Energies Nouvelles*, 67(5):857–875, 2012.
- 665 [61] 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 (Statistical Methodology)*, 74(1):111–132, 2012.
- [62] S. Si, C.-J. Hsieh, and I. S. Dhillon. Memory efficient kernel approximation. In *International Conference on Machine Learning (ICML)*, jun 2014.
- 670 [63] M. L. Stein. Statistical properties of covariance tapers. *Journal of Computational and Graphical Statistics*, 22(4):866–885, 2013.
- [64] M. L. Stein. Limitations on low rank approximations for covariance matrices of spatial data. *Spatial Statistics*, 8:1–19, 2014.
- [65] M. L. Stein, Z. Chi, and L. J. Welty. Approximating likelihoods for large spatial data sets. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 66(2):275–296, 2004.
- 675 [66] Y. Sun, B. Li, and M. G. Genton. Geostatistics for large datasets. In M. Porcu, J. M. Montero, and M. Schlather, editors, *Space-Time Processes and Challenges Related to Environmental*

Problems, pages 55–77. Springer, 2012.

[67] Y. Sun and M. L. Stein. Statistically and computationally efficient estimating equations for large spatial datasets. *Journal of Computational and Graphical Statistics*, 25(1):187–208, 2016.

[68] E. Zander. SGLib - A Matlab Toolbox for Stochastic Galerkin Methods, Feb. 2010.

A (Block) Cluster trees and admissibility condition

Let I be an index set of all degrees of freedom. Denote for each index $i \in I$ corresponding to a basis function b_i the support $\mathcal{G}_i := \text{supp } b_i \subset \mathbb{R}^d$. Now we define two trees which are necessary for the definition of hierarchical matrices. These trees are labeled trees where the label of a vertex t is denoted by \hat{t} .

Definition A.1. (Cluster Tree $T_{I \times I}$)

A finite tree T_I is a cluster tree over the index set I if the following conditions hold:

- I is the root of T_I and a subset $\hat{t} \subseteq I$ holds for all $t \in T_I$.
- If $t \in T_I$ is not a leaf, then the set of sons $\text{sons}(t)$ contains disjoint subsets of I and the subset \hat{t} is the disjoint union of its sons, $\hat{t} = \bigcup_{s \in \text{sons}(t)} \hat{s}$.
- If $t \in T_I$ is a leaf, then $|\hat{t}| \leq n_{\min}$ for a fixed number n_{\min} .

Definition A.2. (Block Cluster Tree $T_{I \times I}$) Let T_I be a cluster tree over the index set I . A finite tree $T_{I \times I}$ is a block cluster tree based on T_I if the following conditions hold:

- $\text{root}(T_{I \times I}) = I \times I$.
- Each vertex b of $T_{I \times I}$ has the form $b = (\tau, \sigma)$ with clusters $\tau, \sigma \in T_I$.
- For each vertex (τ, σ) with $\text{sons}(\tau, \sigma) \neq \emptyset$, we have

$$\text{sons}(\tau, \sigma) = \begin{cases} (\tau, \sigma') : \sigma' \in \text{sons}(\sigma), & \text{if } \text{sons}(\tau) = \emptyset \wedge \text{sons}(\sigma) \neq \emptyset \\ (\tau', \sigma) : \tau' \in \text{sons}(\tau), & \text{if } \text{sons}(\tau) \neq \emptyset \wedge \text{sons}(\sigma) = \emptyset \\ (\tau', \sigma') : \tau' \in \text{sons}(\tau), \sigma' \in \text{sons}(\sigma), & \text{otherwise} \end{cases}$$

- The label of a vertex (τ, σ) is given by $\widehat{(\tau, \sigma)} = \widehat{\tau} \times \widehat{\sigma} \subseteq I \times I$.

We can see that $\widehat{\text{root}(T_{I \times I})} = I \times I$. This implies that the set of leaves of $T_{I \times I}$ is a partition of $I \times I$.

B Rank- k Adaptive Cross Approximation (ACA)

An \mathcal{H} -matrix contains many sub-blocks, which can be well approximated by low-rank matrices. How to compute these low-rank approximations? The truncated singular value decomposition is accurate, but very slow. HLIBpro uses the Adaptive Cross Approximation method (ACA) [23] and its improved modifications such as ACA+ and HACA [6, 10, 12].

Definition B.1. Let $\mathbf{R} \in \mathbb{R}^{p \times q}$, then factorization

$$\mathbf{R} = \mathbf{A}\mathbf{B}^\top, \quad \text{where } \mathbf{A} \in \mathbb{R}^{p \times k}, \quad \mathbf{B} \in \mathbb{R}^{q \times k}, \quad k \ll \min\{p, q\} \in \mathbb{N}. \quad (\text{B.1})$$

we will call a *low-rank representation*.

Note that any matrix of rank k can be represented in the form (B.1).

Suppose that b is a sub-block in the block-cluster tree, and $\mathbf{R} := \mathbf{C}|_b$. Suppose it is known that \mathbf{R} could be approximated by a rank- k matrix. We explain below how to compute \mathbf{R} in the form (B.1). ACA is especially effective for assembling low-rank matrices. It requires to compute only k columns and k rows of the matrix under consideration and, thus, has the computational cost $k(p+q)$. In [23] it is proved that if there exists a sufficiently good low-rank approximation, then there also exists a cross approximation with almost the same accuracy in the sense of the spectral norm. The ACA algorithm computes vectors a_ℓ and b_ℓ which form $\tilde{\mathbf{R}} = \sum_{\ell=1}^k a_\ell b_\ell^\top$ such that $\|\mathbf{R} - \tilde{\mathbf{R}}\| \leq \varepsilon$, where ε is the desired accuracy [10, 13]. In [12] the reader can also find different counterexamples when the standard ACA algorithm does not work. Here we present the standard version of the ACA algorithm.

Algorithm B.1. ACA algorithm

begin

/* input is a required accuracy ε and a function to compute \mathbf{R}_{ij} */;

/* output is matrix $\tilde{\mathbf{R}}$ */;

$k = 0$; $\tilde{\mathbf{R}} = \mathbf{0}$;

$S = \emptyset$; $T = \emptyset$; /* sets of row and column indices */

do

Take a row $i^* \notin S$;

Subtract $\mathbf{R}_{i^*j} := \mathbf{R}_{i^*j} - \tilde{\mathbf{R}}_{i^*j}$, $j = 1..q$;

Find $\max_j |a_{i^*j}| \neq 0$, $j < q$. Suppose it lies in column j^* ;

Compute all elements b_{ij^*} in column j^* , $i < p$;

Subtract $\mathbf{R}_{ij^*} := \mathbf{R}_{ij^*} - \tilde{\mathbf{R}}_{ij^*}$, $i = 1..p$;

```

     $k := k + 1; S := S \cup \{i^*\}; T := T \cup \{j^*\};$ 
    Compute  $\tilde{\mathbf{R}} = \tilde{\mathbf{R}} + a_{i^*} \cdot b_{j^*}^\top$ ; /* it is rank  $k$  approximation */
    if ( $\|a_{i^*} \cdot b_{j^*}^\top\|_2 \leq \varepsilon \cdot \|a_1 \cdot b_1^\top\|_2$ ) return  $\tilde{\mathbf{R}}$ ;
735   Find  $\max_i |b_{ij^*}|, i < p, i \neq i^*$ . The row where it lies is a new row  $i^*$ ;
    until ( $k < k_{\max}$ )
    return  $\tilde{\mathbf{R}}$ ;
end;
```

740 Note that the algorithm does not compute the whole matrix \mathbf{R} . The subtraction is done only from the elements under consideration, i.e. row a_ℓ and column b_ℓ , $\ell = 1, \dots, k$.

Remark B.1. Further optimisation of the ACA algorithm can be done by the truncated SVD. Suppose that a factorisation of matrix $\mathbf{R} = \mathbf{A}\mathbf{B}^\top$, $\mathbf{A} \in \mathbb{R}^{p \times K}$, $\mathbf{B} \in \mathbb{R}^{q \times K}$, is found by ACA. Suppose also that the rank of \mathbf{R} is k , $k < K$. Then one can apply the truncated SVD algorithm to compute $\mathbf{R} = \mathbf{U}\Sigma\mathbf{V}^\top$ requiring $\mathcal{O}((p+q)K^2 + K^3)$ FLOPS.