

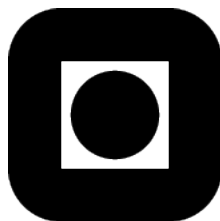
NORGES TEKNISK-NATURVITENSKAPELIGE
UNIVERSITET

**Think continuous: Markovian Gaussian models in
spatial statistics**

by

Daniel Simpson, Finn Lindgren and Håvard Rue

PREPRINT
STATISTICS NO. 9/2011



NORWEGIAN UNIVERSITY OF SCIENCE AND
TECHNOLOGY
TRONDHEIM, NORWAY

This preprint has URL
<http://www.math.ntnu.no/preprint/statistics/2011/S9-2011.pdf>
Daniel Simpson has homepage: <http://www.math.ntnu.no/~daniel>
E-mail: daniel@math.ntnu.no
Address: Department of Mathematical Sciences, Norwegian University of Science and
Technology, N-7491 Trondheim, Norway.

Think continuous: Markovian Gaussian models in spatial statistics

Daniel Simpson*, Finn Lindgren & Håvard Rue
Department of Mathematical Sciences
Norwegian University of Science and Technology
N-7491 Trondheim, Norway

October 31, 2011

Abstract

Gaussian Markov random fields (GMRFs) are frequently used as computationally efficient models in spatial statistics. Unfortunately, it has traditionally been difficult to link GMRFs with the more traditional Gaussian random field models as the Markov property is difficult to deploy in continuous space. Following the pioneering work of Lindgren et al. (2011), we expound on the link between Markovian Gaussian random fields and GMRFs. In particular, we discuss the theoretical and practical aspects of fast computation with continuously specified Markovian Gaussian random fields, as well as the clear advantages they offer in terms of clear, parsimonious and interpretable models of anisotropy and non-stationarity.

1 Introduction

From a practical viewpoint, the primary difficulty with spatial Gaussian models in applied statistics is dimension, which typically scales with the number of observations. Computationally speaking, this is a disaster! It is, however, not a disaster unique to spatial statistics. Time series models, for example, can suffer from the same problems. In the temporal case, the ballooning dimensionality is typically tamed by adding a conditional independence, or *Markovian*, structure to the model. The key advantage of the Markov property for time series models is that the computational burden then grows only linearly (rather than cubically) in the dimension, which makes inference on these models feasible for long time series.

Despite its success in time series modelling, the Markov property has had a less exalted role in spatial statistics. Almost all instances where the Markov property has been used in spatial modelling has been in the form of Markov random fields defined over a set of discrete locations connected by a graph. The most common Markov random field models are *Gaussian* Markov random fields (GMRFs), in which the value of the random field at the nodes is jointly Gaussian (Rue and Held, 2005). GMRFs are typically written as

$$\mathbf{x} \sim N(\boldsymbol{\mu}, \mathbf{Q}^{-1}),$$

where \mathbf{Q} is the *precision matrix* and the Markov property is equivalent to requiring that \mathbf{Q} is sparse, that is $Q_{ij} = 0$ iff x_i and x_j are conditionally independent (Rue and Held, 2005).

As problems in spatial statistics are usually concerned with inferring a spatially continuous effect over a domain of interest, it is difficult to directly apply the fundamentally discrete GMRFs. For this reason, it is commonly stated that there are two essential fields in spatial statistics: the one that uses GMRFs and the one that uses continuously indexed Gaussian random fields. In a recent read paper, Lindgren et al. (2011) showed that these two approaches are not distinct. By carefully utilising the

*Corresponding author. Email: Daniel.Simpson@math.ntnu.no

continuous space Markov property, it is possible to construct Gaussian random fields for which all quantities of interest can be computed using GMRFs!

The most exciting aspect of the Markovian models of Lindgren et al. (2011) is their flexibility. There is no barrier—conceptual or computational—to extending them to construct non-stationary, anisotropic Gaussian random fields. Furthermore, it is even possible to construct them on the sphere and other manifolds. In fact, Simpson et al. (2011a) showed that there is essentially no computational difference between inferring a log-Gaussian Cox process on a rectangular observation window and inferring one on a non-convex, multiply connected region on the sphere! This type of flexibility is not found in any other method for constructing Gaussian random field models.

In this paper we carefully review the connections between GMRFs, Gaussian random fields, the spatial Markov property and deterministic approximation theory. It is hoped that this will give the interested reader some insight into the theory and practice of Markovian Gaussian random fields. In Section 2 we briefly review the practical computational properties of GMRFs. Section 3 we take a detailed tour of the theory of Markovian Gaussian random fields. We begin with a discussion of the spatial Markov property and show how it naturally leads to differential operators. We then present a practical method for approximating Markovian Gaussian random fields and discuss what is meant by a continuous approximation. In particular, we show that deterministic approximation theory can provide essential insights into the behaviour of these approximations. We then discuss some practical issues with choosing sets of basis functions before discussing extensions of the models. Finally we mention some further extensions of the method.

2 Practical computing with Gaussian Markov random fields

Gaussian Markov random fields possess two pleasant properties that make them useful for spatial problems: they facilitate fast computation for large problems, and they are quite stable with respect to conditioning. In this section we will explore these two properties in the context of spatial statistics.

2.1 Fast computations with Gaussian Markov random fields

As in the temporal setting, the Markovian property allows for fast computation of samples, likelihoods and other quantities of interest (Rue and Held, 2005). This allows the investigation of much larger models than would be available using general multivariate Gaussian models. The situation is not, however, as good as it is in the one dimensional case, where all of these quantities can be computed using $\mathcal{O}(n)$ operations, where n is the dimension of the GMRF. Instead, for the two dimensional spatial models, samples and likelihoods can be computed in $\mathcal{O}(n^{3/2})$ operations, which is still a significant saving on the $\mathcal{O}(n^3)$ operations required for a general Gaussian model. A quick order calculation shows that computing a sample from an \tilde{n} -dimensional Gaussian random vector without any special structure takes the same amount of time as computing a sample from GMRF of dimension $n = \tilde{n}^2$!

The key object when computing with GMRFs is the Cholesky decomposition $\mathbf{Q} = \mathbf{L}\mathbf{L}^T$, where \mathbf{L} is a lower triangular matrix. When \mathbf{Q} is sparse, its Cholesky decomposition can be computed very efficiently (see, for instance, Davis, 2006). Once the Cholesky triangle has been computed, it is easy to show that $\mathbf{x} = \boldsymbol{\mu} + \mathbf{L}^{-T}\mathbf{z}$ is a sample from the GMRF $\mathbf{x} \sim N(\boldsymbol{\mu}, \mathbf{Q}^{-1})$ where $\mathbf{z} \sim N(\mathbf{0}, \mathbf{I})$. Similarly, the log density for a GMRF can be computed as

$$\log \pi(\mathbf{x}) = -\frac{n}{2} \log(2\pi) + \sum_{i=1}^n \log L_{ii} - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{Q}(\mathbf{x} - \boldsymbol{\mu}),$$

where L_{ii} is the i th diagonal element of \mathbf{L} and the inner product $(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{Q}(\mathbf{x} - \boldsymbol{\mu})$ can be computed in $\mathcal{O}(n)$ calculations using the sparsity of \mathbf{Q} . It is also possible to use the Cholesky triangle \mathbf{L} to compute $\text{diag}(\mathbf{Q}^{-1})$, which are the marginal variances of the GMRF (Rue and Martino, 2007).

Furthermore, it is possible to sample from \mathbf{x} conditioned on a *small number* of linear constraints $\mathbf{x}|\mathbf{B}\mathbf{x} = \mathbf{b}$, where $\mathbf{B} \in \mathbb{R}^{k \times n}$ is usually a dense matrix and the number of constraints, k , is very small. This occurs, for instance, when the GMRF is constrained to sum to zero. However, if one wishes to sample conditional on data, which usually corresponds to a *large number* of linear constraints, the methods in the next section are almost always significantly more efficient. While direct calculation of the conditional density is possible, when \mathbf{B} is a dense matrix conditioning destroys the Markov structure of the problem. It is still, however, possible to sample efficiently using a technique known as *conditioning by Kriging* (Rue and Held, 2005), whereby an unconditional sample \mathbf{x} is drawn from $N(\boldsymbol{\mu}, \mathbf{Q}^{-1})$ and then corrected using the equation

$$\mathbf{x}^* = \mathbf{x} - \mathbf{Q}^{-1}\mathbf{B}^T(\mathbf{B}\mathbf{Q}^{-1}\mathbf{B}^T)^{-1}(\mathbf{B}\mathbf{x} - \mathbf{b}).$$

When k is small, the conditioning by Kriging update can be computed efficiently from the Cholesky factorisation. We reiterate, however, that when there are a large number of constraints, the conditioning by Kriging method will be inefficient, and, if \mathbf{B} is sparse (as is the case when conditioning on data), it is usually better use the methods in the next subsection.

An alternative method for conditional sampling can be constructed by noting that the conditioning by Kriging update $\mathbf{x}^* = \mathbf{x} - \boldsymbol{\delta}\mathbf{x}$ can be computed by solving the augmented system

$$\begin{pmatrix} \mathbf{Q} & \mathbf{B}^T \\ \mathbf{B} & \mathbf{0} \end{pmatrix} \begin{pmatrix} \boldsymbol{\delta}\mathbf{x} \\ \mathbf{y} \end{pmatrix} = \begin{pmatrix} \mathbf{0} \\ \mathbf{B}\mathbf{x} - \mathbf{b} \end{pmatrix}, \quad (1)$$

where \mathbf{y} is an auxiliary variable. A standard application of Sylvester's inertia theorem (Theorem 8.1.17 in Golub and van Loan, 1996) shows that the matrix in (1) is not positive definite, however the system can still be solved using ordinary direct (or iterative) methods (Simpson et al., 2008). The augmented system (1) reveals the geometric structure of conditioning by Kriging: augmented systems of equations of this form arise from the Karush-Kuhn-Tucker conditions in constrained optimisation (Benzi et al., 2005). In particular, the conditional sample solves the minimisation problem

$$\begin{aligned} \mathbf{x}^* &= \arg \min_{\mathbf{x}^* \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{x}^* - \mathbf{x}\|_{\mathbf{Q}}^2 \\ &\text{subject to } \mathbf{B}\mathbf{x}^* = \mathbf{b}, \end{aligned}$$

where \mathbf{x} is the unconditional sample and $\|\mathbf{y}\|_{\mathbf{Q}}^2 = \mathbf{y}^T \mathbf{Q} \mathbf{y}$. That is, the conditional sample is the closest vector that satisfies the linear constraint to the unconditional sample when the distance is measured in the natural norm induced by the GMRF.

The methods in this section are all predicated on the computation of a Cholesky factorisation. However, for genuinely large problems, it may be impossible to compute or store the Cholesky factor. With this in mind, a suite of methods were developed by Simpson et al. (2007) based on modern iterative methods for solving sparse linear systems. These methods have shown some promise for large problems (Strickland et al., 2011; Aune et al., 2011), although more work is needed on approximating the log density (Simpson, 2009; Aune and Simpson, 2011).

2.2 The effect of conditioning: fast Bayesian inference

The second appealing property of GMRFs is that they behave well under conditioning. The discussion in this section is intimately tied to discretely specified GMRFs, however we will see in Section 3.6 that the formulation below, and especially the matrix \mathbf{A} , has an important role to play in the continuous setting. Consider the simple Bayesian hierarchical model

$$\mathbf{y}|\mathbf{x} \sim N(\mathbf{A}\mathbf{x}, \mathbf{Q}_y^{-1}) \quad (2a)$$

$$\mathbf{x} \sim N(\boldsymbol{\mu}, \mathbf{Q}_x^{-1}), \quad (2b)$$

where \mathbf{A} , \mathbf{Q}_x and \mathbf{Q}_y are sparse matrices. A simple manipulation shows that $(\mathbf{x}^T, \mathbf{y}^T)^T$ is jointly a Gaussian Markov random field with joint precision matrix

$$\mathbf{Q}_{xy} = \begin{pmatrix} \mathbf{Q}_x + \mathbf{A}^T \mathbf{Q}_y \mathbf{A} & -\mathbf{A}^T \mathbf{Q}_y \\ -\mathbf{Q}_y \mathbf{A} & \mathbf{Q}_y \end{pmatrix} \quad (3)$$

and the mean defined implicitly through the equation

$$\mathbf{Q}_{xy} \boldsymbol{\mu}_{xy} = \begin{pmatrix} \mathbf{Q}_x \boldsymbol{\mu} \\ \mathbf{0} \end{pmatrix}.$$

As $(\mathbf{x}^T, \mathbf{y}^T)^T$ is jointly a GMRF, it is easy to see (see, for example, Rue and Held, 2005) that

$$\mathbf{x}|\mathbf{y} \sim N(\boldsymbol{\mu} + (\mathbf{Q}_x + \mathbf{A}^T \mathbf{Q}_y \mathbf{A})^{-1} \mathbf{A}^T \mathbf{Q}_y (\mathbf{y} - \mathbf{A} \boldsymbol{\mu}), (\mathbf{Q}_x + \mathbf{A}^T \mathbf{Q}_y \mathbf{A})^{-1}). \quad (4)$$

It is important to note that the precision matrices for the joint (3) and the conditional (4) distributions are only sparse—and the corresponding fields are only GMRFs—if \mathbf{A} is sparse. This observation directly links the structure of the matrix \mathbf{A} to the availability of efficient inference methods and will be important in the coming sections.

For practical problems in spatial statistics, the model (2) is not enough: there will be unknown parameters in both the likelihood and the latent field. If we group the unknown parameters into a vector $\boldsymbol{\theta}$, we get the following hierarchical model

$$\mathbf{y}|\mathbf{x}, \boldsymbol{\theta} \sim N(\mathbf{A}\mathbf{x}, \mathbf{Q}_y(\boldsymbol{\theta})^{-1}) \quad (5a)$$

$$\mathbf{x}|\boldsymbol{\theta} \sim N(\boldsymbol{\mu}, \mathbf{Q}_x(\boldsymbol{\theta})^{-1}) \quad (5b)$$

$$\boldsymbol{\theta} \sim \pi(\boldsymbol{\theta}). \quad (5c)$$

In order to perform inference on (5), it is common to use Markov chain Monte Carlo (MCMC) methods for sampling from the posterior $\pi(\mathbf{x}, \boldsymbol{\theta}|\mathbf{y})$, however, this is not necessary. It's an easy exercise in Gaussian density manipulation to show that the marginal posterior for the parameters, denoted $\pi(\boldsymbol{\theta}|\mathbf{y})$ can be computed without integration and is given by

$$\pi(\boldsymbol{\theta}|\mathbf{y}) \propto \frac{\pi(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta})\pi(\boldsymbol{\theta})}{\pi(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta})} \Big|_{\mathbf{x}=\mathbf{x}^*},$$

where \mathbf{x}^* can be any point, but is typically taken to be the conditional mode $E(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta})$, and the corresponding marginals $\pi(\theta_j|\mathbf{y})$ can be computed using numerical integration. Similarly, the marginals $\pi(x_i|\mathbf{y})$ can be computed using numerical integration and the observation that, for every $\boldsymbol{\theta}$, $\pi(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta})$ is a GMRF (Rue and Martino, 2007; Rue et al., 2009).

It follows that, for models with Gaussian observations, it is possible to perform *deterministic* inference that is exact up to the error in the numerical integration. In particular, if there are only a moderate number of parameters, this will be extremely fast. For non-Gaussian observation processes, exact deterministic inference is no longer possible, however, Rue et al. (2009) showed that it is possible to construct extremely accurate approximate inference schemes by cleverly deploying a series of Laplace approximation. The integrated nested Laplace approximation (INLA) has been used successfully on a large number of spatial problems (see, for example, Fong et al., 2010; Akerkar et al., 2010; Schrödle and Held, 2011; Riebler et al., 2011; Cameletti et al., 2011; Illian et al., 2011) and a user friendly R interface is available from <http://r-inla.org>.

3 Continuously specified, Markovian Gaussian random fields

One of the primary aims of spatial statistics is to infer a *spatially continuous* surface $x(s)$ over the region of interest. It is, therefore, necessary to build probability distributions over the space of

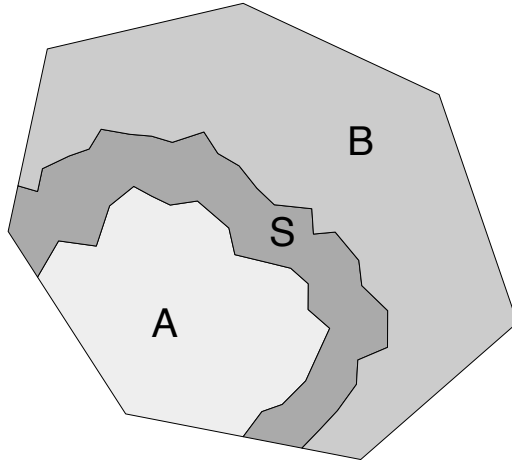


Figure 1: An illustration of the spatial Markov property. If, for any appropriate set S , $\{x(s) : s \in A\}$ is independent of $\{x(s) : s \in B\}$ given $\{x(s) : s \in S\}$, then the field $x(s)$ has the spatial Markov property.

functions, and the standard way of doing this is to construct Gaussian random fields, which are the generalisation to functions of multivariate Gaussian distributions in the sense that for any collection of points $(s_1, s_2, \dots, s_p)^T$, the field evaluated at those points is jointly Gaussian. In particular $\mathbf{x} \equiv (x(s_1), x(s_2), \dots, x(s_p))^T \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, where the covariance matrix is given by $\Sigma_{ij} = c(s_i, s_j)$ for some positive definite covariance function $c(\cdot, \cdot)$. In most commonly used cases, the covariance function is non-zero everywhere and, as a result, $\boldsymbol{\Sigma}$ is a dense matrix.

It is clear that we would like to transfer some of the pleasant computational properties of GMRFs, which are outlined above, to the Gaussian random field setting. The obvious barrier to this is that classical GMRF models are strongly tied to discrete sets of points, such as graphs (Rue and Held, 2005) and lattices (Besag, 1974). Throughout this section, we will discuss the recent work of Lindgren et al. (2011) that has broken down the barrier between GMRFs and spatially continuous Gaussian random field models.

3.1 The spatial Markov property

For temporal processes, defining the Markov property is greatly simplified by the structure of time: its directional nature and the clear distinction between past, present and future allow for a very natural discussion of neighbourhoods. Unfortunately, space is far less structured and, as such the Markov property is harder to define exactly. Intuitively, however, the definition is generalised in an obvious way and is demonstrated by Figure 1. Informally, a Gaussian random field $x(s)$ has the spatial Markov property if, for every appropriate set S separating A and B , the values of $x(s)$ in A are conditionally independent of the values in B given the values in S . A formal definition of the spatial Markov property can be found in Rozanov (1977).

It is not immediately obvious how the spatial Markov property can be used for computational inference. However, in an almost completely ignored paper, Rozanov (1977) provided the vital characterisation of Markovian Gaussian random fields in terms of their power spectra. The power spectrum of a stationary Gaussian random field is defined as the Fourier transform of its covariance function $c(\mathbf{h})$, that is

$$R(\mathbf{k}) \equiv \frac{1}{(2\pi)^d} \int_{\mathbb{R}^d} \exp(-i\mathbf{k}^T \mathbf{h}) c(\mathbf{h}) d\mathbf{h}.$$

Rozanov showed that a stationary field is Markovian if and only if $R(\mathbf{k}) = 1/p(\mathbf{k})$, where $p(\mathbf{k})$ is a positive, symmetric polynomial.

3.2 From spectra to differential operators

While Rozanov's characterisation of stationary Markovian random fields in terms of their power spectra is elegant, it is not obvious how we can turn it into a useful computational tool. The link comes from Fourier theory: there is a one-to-one correspondence between polynomials in the frequency space and differential operators. To see this, define the covariance operator C of the Markovian Gaussian random field as the convolution

$$\begin{aligned} C[f(\cdot)](\mathbf{h}) &= \int_{\mathbb{R}^d} c(\mathbf{h}' - \mathbf{s}) f(\mathbf{h}') d\mathbf{h}' \\ &= \int_{\mathbb{R}^d} \exp(i\mathbf{k}^T \mathbf{h}) \frac{\hat{f}(\mathbf{k})}{p(\mathbf{k})} d\mathbf{k}, \end{aligned}$$

where $p(\mathbf{k}) \equiv 1/R(\mathbf{k}) = \sum_{|\mathbf{i}| \leq \ell} a_{\mathbf{i}} \mathbf{k}^{\mathbf{i}}$ is a d -variate, positive, symmetric polynomial of degree ℓ ; $\mathbf{i} = (i_1, \dots, i_d)^T \in \mathbb{N}^d$ is a multi-index, meaning that $\mathbf{k}^{\mathbf{i}} = \prod_{l=1}^d k_l^{i_l}$ and $|\mathbf{i}| = \sum_{l=1}^d i_l$; $f(\cdot)$ is a smooth function that goes to zero rapidly at infinity; and $\hat{f}(\mathbf{k})$ is the Fourier transform of $f(\mathbf{h})$. It follows from Fourier theory that the covariance operator C has an inverse, which we will call the *precision* operator Q , and it is given by

$$\begin{aligned} Q[f(\cdot)](\mathbf{h}) &\equiv C^{-1}[f(\cdot)](\mathbf{h}) = \int_{\mathbb{R}^d} \exp(i\mathbf{k}^T \mathbf{h}) \hat{f}(\mathbf{k}) p(\mathbf{k}) d\mathbf{k} \\ &= \sum_{|\mathbf{i}| \leq \ell} a_{\mathbf{i}} D^{\mathbf{i}} f(\mathbf{h}), \end{aligned}$$

where $D^{\mathbf{i}} = i^{|\mathbf{i}|} \frac{\partial^{|\mathbf{i}|}}{\partial h_1^{i_1} \partial h_2^{i_2} \dots \partial h_d^{i_d}}$ are the appropriate multivariate derivatives.

The following proposition summarises the previous discussion and specialises the result to isotropic fields.

Proposition 1. *A stationary Gaussian random field $x(s)$ defined on \mathbb{R}^d is Markovian if and only if its covariance operator $C[f(\cdot)]$ has an inverse of the form*

$$Q[f(\cdot)](\mathbf{h}) = \sum_{|\mathbf{i}| \leq \ell} a_{\mathbf{i}} D^{\mathbf{i}} f(\mathbf{h}),$$

where $a_{\mathbf{i}}$ are the coefficients of a real, symmetric polynomial $p(\mathbf{k}) = \sum_{|\mathbf{i}| \leq \ell} a_{\mathbf{i}} \mathbf{k}^{\mathbf{i}}$. Furthermore, $x(s)$ is isotropic if and only if

$$Q[f(\cdot)](\mathbf{h}) = \sum_{i=0}^{\ell} \tilde{a}_i (-\Delta)^i f(\mathbf{h}),$$

where $\Delta = \sum_{i=1}^d \frac{\partial^2}{\partial h_i^2}$ is the d -dimensional Laplacian and $\tilde{p}(t) = \sum_{i=0}^{\ell} \tilde{a}_i t^i$ is a real, symmetric univariate polynomial.

The critical point of the above paragraph is that for Markovian Gaussian random fields, the covariance is the inverse of a *local* operator, in the sense that the value of $Q[f(\cdot)](\mathbf{h})$ only depends on the value of $f(\mathbf{h})$ in an infinitesimal neighbourhood of \mathbf{h} . This is in stark contrast to the covariance operator, which is an integral operator and therefore depends on the value of $f(\mathbf{h})$ everywhere that the covariance function is non-zero. It is this locality that will lead us to GMRFs and the promised land of sparse precision matrices and fast computations.

3.3 Stochastic differential equations and Matérn fields

The discussion in the previous subsection gave a very general form of Markovian Gaussian random fields. In this section we will, for the sake of sanity, simplify the setting greatly. Consider the stochastic partial differential (SPDE)

$$(\kappa^2 - \Delta)^{\alpha/2} x(s) = W(s), \quad (6)$$

where $W(s)$ is Gaussian white noise and $\alpha > 0$. The solution to the SPDE will be a Gaussian random field and some formal manipulations (made precise by Rozanov, 1977) show that the precision operator is

$$\begin{aligned} Q &= [E(xx^*)]^{-1} = \left[(\kappa^2 - \Delta)^{-\alpha/2} E(WW^*) (\kappa^2 - \Delta)^{-\alpha/2} \right]^{-1} \\ &= (\kappa^2 - \Delta)^\alpha. \end{aligned}$$

It follows that Q is a local differential operator when α is an integer and, by Proposition 1, the stationary solution to (6) for integer α is a Markovian Gaussian random field.

Somewhat surprisingly, we have now ventured back into the more common parts of spatial statistics: Whittle (1954, 1963) showed that, for any $\alpha > d/2$, the solution to (6) has a Matérn covariance function. That Matérn family of covariance functions, which is given by

$$c_{\sigma^2, \nu, \kappa}(\mathbf{h}) = \frac{\sigma^2}{2^{\nu-1} \Gamma(\nu)} (\kappa \|\mathbf{h}\|)^\nu K_\nu(\kappa \|\mathbf{h}\|),$$

where σ^2 is the variance parameter, κ controls the range, and $\nu = \alpha - d/2 > 0$ is the shape parameter, is one of the most widely used families of stationary, isotropic covariance functions. The results of the previous section show that, when $\nu + d/2$ is an integer, the Matérn fields are Markovian.

3.4 Approximating Gaussian random fields: the finite element method

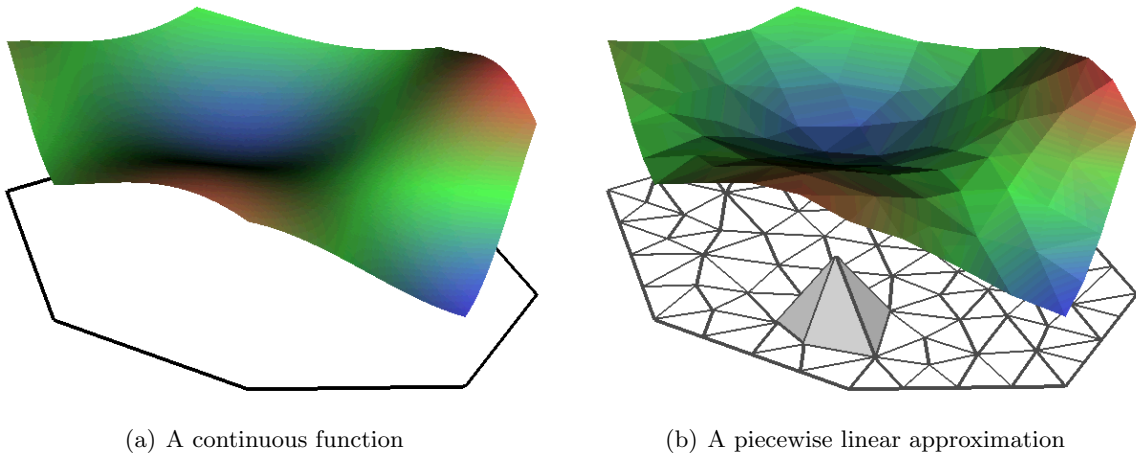


Figure 2: Piecewise linear approximation of a function over a triangulated mesh.

Having now laid all of the theoretical groundwork, we can consider the fundamental question in this paper: how do we use GMRFs to approximate a Gaussian random field? In order to construct a good approximation, it is vital to remember that every realisation of a Gaussian random field is a *function*. It follows that any sensible method for approximating Gaussian random fields is

necessarily connected to a method for approximating an appropriate class of deterministic functions. Therefore, following Lindgren et al. (2011), we are on the hunt for simple methods for approximating functions.

A reasonably simple method for approximating continuous functions is demonstrated in Figure 2, which shows a piecewise linear approximation to a deterministic function $f(s)$ that is defined over a triangulation of the domain. This approximation is of the form

$$f(s) \approx f_h(s) = \sum_{i=1}^n w_i \phi_i(s),$$

where the basis function $\phi_i(s)$ is the piecewise linear function that is equal to one at the i th vertex of the mesh and is zero at all other vertices and the subscript h denotes the largest triangle edge in the mesh and is used to differentiate between piecewise linear functions over a mesh and general functions. The grey pyramid in Figure 2 is an example of a basis function. With this in hand, we can define a piecewise linear Gaussian random field as

$$x_h(s) = \sum_{i=1}^n w_i \phi_i(s), \quad (7)$$

where $\phi_i(s)$ are as before and the weights \mathbf{w} are now jointly a Gaussian random vector.

Our aim is now to find the statistical properties of \mathbf{w} that make $x_h(s)$ approximate the Markovian Matérn fields. In order to do this, we use the characterisation of a Matérn field as the stationary solution to (6). For simplicity, let us consider $\alpha = 2$ on a domain $\Omega \subset \mathbb{R}^2$. Any solution of (6) also satisfies, for any suitable function $\psi(s)$,

$$\int_{\Omega} \psi(s) (\kappa^2 - \Delta) x(s) ds = \int_{\Omega} \psi(s) W(ds)$$

and an application of Green's formula leads to

$$\int_{\Omega} \kappa^2 \psi(s) x(s) + \nabla \psi(s) \cdot \nabla x(s) dx = \int_{\Omega} \psi(s) W(ds), \quad (8)$$

where the integrals on the right hand sides are integrals with respect to white noise (see Chapter 5 of Adler and Taylor, 2007) and the second line follows from Green's formula and the (new) condition that the normal derivative of $x(s)$ vanishes on the boundary of Ω . Furthermore, if we find $x(s)$ such that (8) holds for *any* sensible $\psi(s)$, then $x(s)$ is the weak solution to (6) and it can be shown that it is a Gaussian random field with the appropriate Matérn covariance function.

Unfortunately, we cannot test (8) against every function $\psi(s)$, so we will instead chose a finite set $\{\psi_j(s)\}_{j=1}^n$ to test against. Substituting $x_h(s)$ into (8) and testing against this set of ψ_j s, we get the system of linear equations

$$\sum_{i=1}^n \left(\kappa^2 \int_{\Omega} \psi_j(s) \phi_i(s) ds + \int_{\Omega} \nabla \psi_j(s) \cdot \nabla \phi_i(s) ds \right) w_i = \int_{\Omega} \psi_j(s) dW(s), \quad j = 1, \dots, n. \quad (9)$$

Finally, we chose our test functions $\psi_j(s)$ to be the same as our basis functions $\phi_i(s)$ and arrive at the Galerkin finite element method. For piecewise linear functions over a triangular mesh, it is easy to compute all of the integrals on the left hand side of (9). The white noise integral on the right hand side can be computed and it is Gaussian with mean zero and

$$\text{Cov} \left(\int_{\Omega} \phi_i(s) dW(s), \int_{\Omega} \phi_j(s) dW(s) \right) = \int_{\Omega} \phi_i(s) \phi_j(s) ds. \quad (10)$$

We can, therefore write (9) in matrix form as

$$\mathbf{K}\tilde{\mathbf{w}} \sim N(\mathbf{0}, \tilde{\mathbf{C}}),$$

where $\mathbf{K} = \kappa\tilde{\mathbf{C}} + \mathbf{G}$ and the matrices are given by $\tilde{C}_{ij} = \int_{\Omega} \phi_i(s)\phi_j(s) ds$ and $G_{ij} = \int_{\Omega} \nabla\psi_j(s) \cdot \nabla\phi_i(s) ds$.

A quick look at the definitions of $\tilde{\mathbf{C}}$ and \mathbf{G} shows that, due to the highly local nature of the basis functions, these matrices are sparse. Unfortunately, $\tilde{\mathbf{C}}^{-1}$ is a dense matrix and, therefore $\tilde{\mathbf{w}}$ will not be a GMRF. However, replacing $\tilde{\mathbf{C}}$ by the diagonal matrix $\mathbf{C} = \text{diag}(\int_{\Omega} \phi_i(s) ds, i = 1, \dots, n)$ gives essentially the same result numerically (Bolin and Lindgren, 2009) and can be shown not to increase the rate of convergence of the approximation (Appendix C.5 of Lindgren et al., 2011). Using \mathbf{C} , it follows that the solution to

$$\mathbf{K}\mathbf{w} \sim N(\mathbf{0}, \tilde{\mathbf{C}})$$

is a GMRF with zero mean and sparse precision matrix $\mathbf{Q} = \mathbf{K}^T \mathbf{C} \mathbf{K}$. The GMRFs that correspond to other integer values of α can be found in Lindgren et al. (2011).

3.5 A continuous approximation to a continuous random field

We have now derived a GMRF \mathbf{w} from the continuous Matérn field $x(s)$. It is, therefore, reasonable to wonder how well \mathbf{w} approximates $x(s)$. *It doesn't!* The GMRF \mathbf{w} was defined as the weights of a basis function expansion (7) and only make sense in this context. In particular, \mathbf{w} is not trying to approximate $x(s)$ at the mesh vertices. Instead, the piecewise linear Gaussian random field $x_h(s) = \sum_{i=1}^n w_i \phi_i(s)$ tries to approximate the true random field *everywhere*. Because of the nature of this continuous approximation, $x_h(s)$ will necessarily overestimate the variance at the vertices and underestimate in the centre of the triangles.

One of the real advantages of using piecewise linear basis functions is that a lot of work has gone into working out their approximation properties (Brenner and Scott, 2007). We can leverage this information to get hard bounds on the convergence of $x_h(s)$ to $x(s)$. The following theorem is a simple example of this type of result. It shows that functionals of both the field and its derivative converge to the true functionals and the error is $\mathcal{O}(h)$, where h is the length of the largest edge in the mesh. The theorem also links the convergence of functionals of the random field to how well the piecewise linear basis functions can approximate functions in the Sobolev space H^1 , which consists of square integrable functions $f(s)$ for which $\|f\|_{H^1}^2 = \kappa^2 \int_{\Omega} f(s)^2 ds + \int_{\Omega} \nabla f(s) \cdot \nabla f(s) ds$ is finite.

Theorem 1. *Let $L = \kappa^2 - \Delta$. Then, for any $f \in H^1$,*

$$E \left(\int_{\Omega} f(s) L(x(s) - x_h(s)) ds \right)^2 \leq ch^2 \|f\|_{H^1}^2,$$

where c is a constant and h is the size of the largest triangle edge in the mesh.

Proof. Let $f \in H^1$, $f_h(s)$ be the H^1 -projection of f onto the finite element space $V_h = \text{span}\{\phi_1(s), \dots, \phi_n(s)\}$, that is let f_h be the solution to

$$\min_{f_h \in V_h} \|f - f_h\|_{H^1}.$$

It follows that

$$\begin{aligned} \int_{\Omega} f(s) Lx_h(s) ds &= \int_{\Omega} (f(s) - f_h(s)) Lx_h(s) ds + \int_{\Omega} f_h(s) Lx_h(s) ds \\ &= \int_{\Omega} f_h(s) Lx_h(s) dx \\ &= \int_{\Omega} f_h(s) dW(s), \end{aligned}$$

where the second equality follows from the Galerkin property of $x_h(s)$ (Brenner and Scott, 2007), which states that $\int_{\Omega} g(s) Lx_h(s) ds = 0$ whenever $g(s)$ is in the orthogonal complement of V_h with respect to the H^1 inner product. It follows directly that

$$\int_{\Omega} f(s) L(x(s) - x_h(s)) ds = \int_{\Omega} (f(s) - f_h(s)) dW(s).$$

Therefore, it follows from the properties of white noise integrals that

$$\begin{aligned} E \left(\int_{\Omega} f(s) L(x(s) - x_h(s)) ds \right)^2 &= E \left(\int_{\Omega} (f(s) - f_h(s)) dW(s) \right)^2 \\ &= \int_{\Omega} (f(s) - f_h(s))^2 ds. \end{aligned}$$

It follows from Theorem 4.4.20 in Brenner and Scott (2007) that (under some suitable assumptions on the triangulation)

$$\|f_n - f\|_{L^2(\Omega)} \leq ch \|f\|_{H^1}.$$

□

The key lesson from the proof of Theorem 1 is that the convergence of functionals of $x_h(s)$ depends solely on how well the basis functions can approximate a fixed H^1 function. Therefore, it is vital to consider the approximation properties of your basis functions!

3.6 Choosing basis functions: don't forget about \mathbf{A}

While we have computed everything with respect to piecewise linear functions, the methods considered above work for any set of test and basis functions for which all of the computations make sense. However, we strongly warn against using any of the more esoteric choices. There are two main issues that can appear. The first issue is that the wrong choice of basis functions will destroy the Markov structure of the posterior model, which will annihilate the computational gains we have worked so hard for. The second issue, which is related to Theorem 1 is much more problematic: not all sets of basis functions will provide good approximations to $x(s)$.

In Section 2.2, we looked at the GMRF computations for hierarchical Gaussian models. Consider a simple Gaussian observation process of the form

$$y_i \sim N(x_h(s_i), \sigma^2), \quad i = 1, \dots, N$$

where the number of data points N does not need to be related to the number of mesh vertices n . It follows that the datapoint (s_i, y_i) requires the computation of the sum $\sum_{j=1}^n w_j \phi_j(s_i)$. When the basis functions are local, there are only a few non-zero terms in this sum and the corresponding matrix \mathbf{A} , which has entries $A_{ij} = \phi_j(s_i)$, is sparse. For piecewise linear functions on a triangular mesh, each row of \mathbf{A} has at most 3 non-zero entries.

On the other hand, consider a basis consisting of the first n functions from the Karhunen-Lo  ve expansion of $x(s)$. In this case, it's easy to show that the precision matrix for \mathbf{w} will be diagonal. Unfortunately, these basis functions are usually non-zero everywhere, so \mathbf{A} will be a completely dense $N \times n$ matrix, which, for large datasets will become the dominant computational cost.

The approximation properties of sets of basis functions are typically very hard to determine. There are, however, some good guiding principles. The first principle is that your basis functions should do a good job approximating simple functions. You usually want to be able to at least approximate constant and linear functions well. A second guiding principle is that you should be very careful when your basis functions depend on a parameter that is being estimated. It is important to check that the approximation is still sensible over the entire parameter range. Figure 3 is an example of this problem, taken from Simpson et al. (2011b), shows the best approximation to a constant function

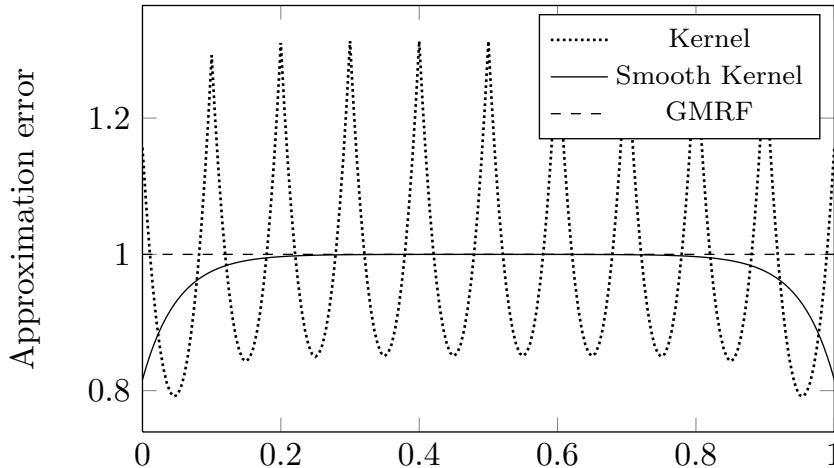


Figure 3: Taken from (Simpson et al., 2011b). The error in the best approximation to $x(s) = 1$, $s \in [0, 1]$. The dotted line is the naïve kernel basis, while the solid line is an integrated kernel basis derived by Simpson et al. (2011b). The piecewise linear basis (dashed line) reproduces the function exactly.

using the basis derived from a convolution kernel approximation to a one dimensional Matérn field with $\nu = 3/2$ and $\kappa = 20$. The basis functions are computed with respect to a fixed grid of 11 equally spaced knots s_i , $i = 1, \dots, 11$ and are given by

$$\phi_i(s) = \frac{1}{(2\pi)^{1/2}\kappa} \exp(-\kappa|s - s_i|).$$

When κ gets large, the basis functions get sharper and the approximation gets worse. At its worst, the effective range of the kernel becomes shorter than grid of knots. Therefore, if κ becomes large relative to the grid spacing, both Kriging estimates and variance estimates will be badly infected by the poor choice of basis function. On the other hand, it can be shown that piecewise linear Galerkin finite element approximations are stable in the sense that the norm of the approximation can be bounded above by an appropriate norm of the true solution. This means that a piecewise linear basis, when used in this way, will *never* display this bad behaviour.

Another consideration when choosing sets of basis functions is the smoothness of the random field. Theorem 1 directly links the convergence of the approximate fields to how well the basis functions can approximate a certain class of functions. This will usually be the case. Typically, the smoother the field is, the more useful higher order “spectral” bases will be. For a smooth enough field, it may actually be cheaper to use a small, well chosen global basis than a large basis full of local functions.

3.7 Extending the models: anisotropy, drift, and space-time models on manifolds

One of the main advantages to the SPDE formulation is that it is easy to construct non-stationary variants. Non-stationary fields can now be defined by letting the parameters in (6), be space-dependent. For example, $\log \kappa$ can be expanded using a few weighted smooth basis functions

$$\log \kappa(\mathbf{s}) = \sum_i \beta_i b_i(\mathbf{s}) \quad (11)$$

and similar expansions can be used for τ . This extension requires only minimal changes to the method used in the stationary case. More interestingly, we can also incorporate models of spatially

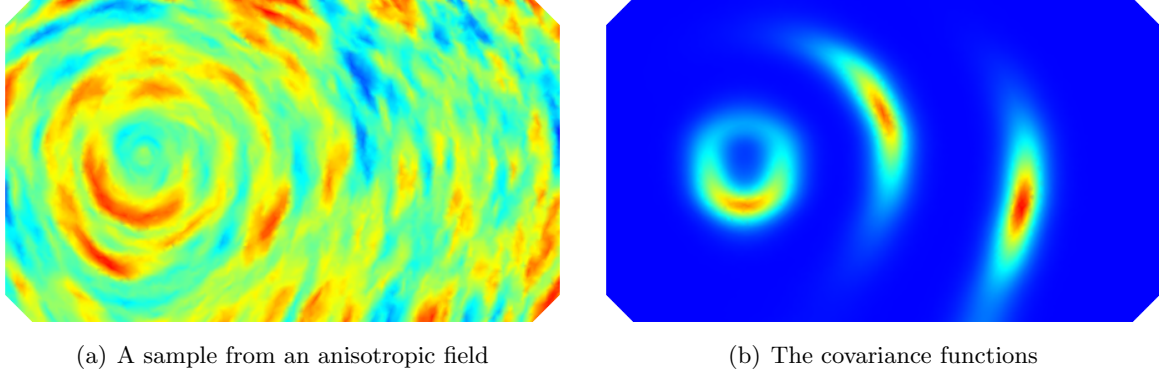


Figure 4: An anisotropic field defined through (6) when the Laplacian is replaced with the anisotropic variant $\nabla \cdot (\mathbf{D}(s)\nabla(\tau x(s)))$. The right hand figure shows the approximate covariance function evaluated at three points and the non-stationarity is evident.

varying anisotropy by replacing the Laplacian Δ in (6) with the more general operator

$$\nabla \cdot (\mathbf{D}(s)\nabla(\tau x(s))) \equiv \sum_{i,j=1}^2 \frac{\partial}{\partial s_i} \left(D_{ij}(s) \frac{\partial}{\partial s_j} (\tau(s)x(s)) \right).$$

These models can be extended even farther by incorporating a “drift” term, as well as temporal dependence, which leads to the general model

$$\frac{\partial}{\partial t}(\tau u(s, t)) + (\kappa^2(\tau x(s, t)) - \nabla \cdot (\mathbf{D}\nabla(\tau x(s, t)))) + \mathbf{b} \cdot \nabla(\tau x(s, t)) = W(s, t), \quad (12)$$

where t is the time variable and \mathbf{b} is a vector describing the direction of the drift and the dependence of κ , τ , \mathbf{D} and \mathbf{b} on s and t has been suppressed. The concepts behind the construction of Markovian approximations to the stationary SPDE transfer almost completely to this general setting, however more care needs to be taken with the exact form of the discretisation (Fuglstad, 2010).

A strong advantage of defining non-stationarity through (12) is that the underlying physics is well understood. For example the second order term $\nabla \cdot (\mathbf{D}\nabla(\tau x(s, t)))$ describes the diffusion, where the matrix $\mathbf{D}(s, t)$ describes the preferred direction of diffusion at location s and time t . A spatial random field with a non-constant diffusion matrix $\mathbf{D}(s)$ is shown in Figure 4. Similarly, the first order term $\mathbf{b} \cdot \nabla(\tau x(s, t))$ describes an advection effect and $\mathbf{b}(s, t)$ can be interpreted as the direction of the forcing field. We note that unlike other methods for modelling non-stationarity, the parametrisation in this model is *unconstrained*. This is in contrast to the deformation methods of Sampson and Guttorp (1992), in which the deformation must be constrained to be one-to-one. Initial work on parameterising and performing inference on these models has been undertaken by Fuglstad (2011).

An interesting consequence of defining our models through *local* stochastic partial differential equations, such as (12), is that the SPDEs still make sense when \mathbb{R}^d is replaced by a space that is only locally flat. We can, therefore, use (12) to *define* non-stationary Gaussian fields on manifolds, and still obtain a GMRF representation. Furthermore, the computations can be done in essentially the same way, the only change is that the Gaussian noise is now a Gaussian random measure, and that we need to take into account the local curvature when computing the integrals to obtain the solution. Figure 5 shows a realisation of a non-stationary Gaussian random field on a sphere, with a model similar to the one used in Figure 4. The solution is still explicit, so all elements of the precision matrix, for a fixed triangularisation, can be directly computed with no extra cost. The ability to construct computationally efficient representations of non-stationary random fields on a manifold is important, for example, when modelling environmental data on the sphere.

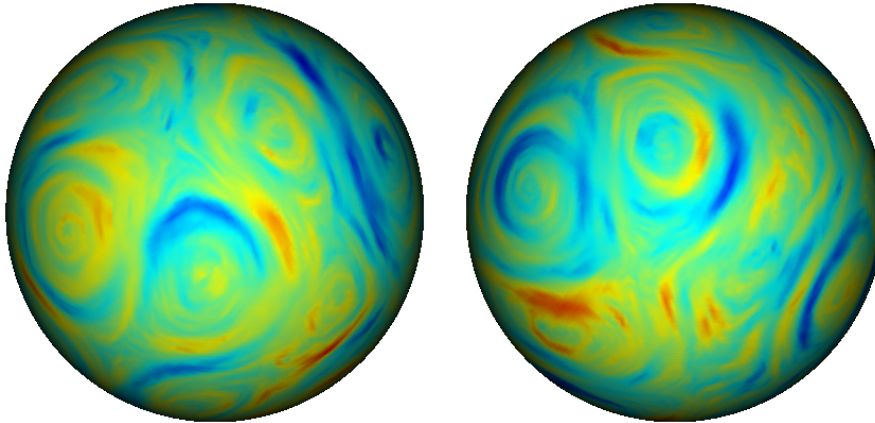


Figure 5: A sample from a non-stationary, anisotropic random field on the sphere.

3.8 Further extensions: areal data, multivariate random fields and non-integer α s

The combination of a continuously specified approximate random field and a Markovian structure allows for the straightforward assimilation of areal and point-referenced data. The trick is to utilise the matrix \mathbf{A} in the observation equation (5a). The link between the latent GMRF and the point-referenced data is as described above. For the areal data, the dependence is through integrals of the random field and these integrals can be computed exactly for $x_h(s)$ and the resulting sums go into \mathbf{A} . If a more complicated functional of the random field has been observed, this can be approximated using numerical quadrature. This has been recently applied by Simpson et al. (2011a) to inferring log-Gaussian point processes.

It is also possible to extend the SPDE framework considered above to a multivariate setting. Primarily, the idea is to replace a single SPDE with a system of SPDEs. It can be shown (work in progress) that these models overlap with, but are not equivalent to, the multivariate Matérn models constructed by Gneiting et al. (2010).

The only Markovian Matérn models are those where the smoothing parameter is $d/2$ less than an integer. It turns out, however, that we can construct good Markov random field approximations when α isn't an integer (see the authors' response in Lindgren et al., 2011). This is essentially a powerful extension of the GMRF approximations of Rue and Tjelmeland (2002). We are currently working on approximating non-Matérn random fields using this method.

4 Conclusion

In this paper we have surveyed the deep and fascinating link between the continuous Markov property and computationally efficient inference. The material in this paper barely scratches the surface of the questions, applications and oddities of these fields (for instance, Bolin (2011) uses a similar construction for non-Gaussian noise processes). However, we have demonstrated that by carefully considering the Markov property, we are able to greatly boost our modelling capabilities while at the same time ensuring that the models are fast to compute with. The combination of flexible modelling and fast computations ensures that investigations into Markovian random fields will continue to produce interesting and useful results well into the future.

References

R. J. Adler and J. Taylor. *Random Fields and Geometry*. Springer Monographs in Mathematics. Springer, 2007.

- R. Akerkar, S. Martino, and H. Rue. Approximate bayesian inference for survival models. *Scandinavian Journal of Statistics*, 28(3):514–528, 2010.
- E. Aune and D. P. Simpson. Hyper-parameter estimation involving high dimensional gaussian distributions. In *ISI 2011*, 2011.
- E. Aune, J. Eidsvik, and Y. Pokem. Iterative numerical methods for sampling from high dimensional gaussian distributions. Technical Report 4/2011, NTNU, 2011.
- M. Benzi, G. Golub, and J. Liesen. Numerical solutions of saddle point problems. *Acta Numerica*, 14:1–137, 2005.
- J. Besag. Spatial interaction and the statistical analysis of lattice systems (with discussion). *Journal of the Royal Statistical Society, Series B*, 36(2):192–225, 1974.
- D. Bolin. Spatial Matérn fields driven by non-Gaussian noise. Technical Report 2011:4, Lund University, 2011.
- D. Bolin and F. Lindgren. Wavelet Markov models as efficient alternatives to tapering and convolution fields. Technical Report 2009:13, Lund University, 2009.
- S. C. Brenner and R. Scott. *The Mathematical Theory of Finite Element Methods*. Springer, 3rd edition, 2007.
- M. Cameletti, F. Lindgren, D. P. Simpson, and H. Rue. Spatio-temporal modelling of particulate matter concentration through the SPDE approach. *AStA Adv Stat Anal*, Submitted, 2011.
- T. A. Davis. *Direct methods for sparse linear systems*. SIAM Book Series on the Fundamentals of Algorithms. SIAM, Philadelphia, 2006.
- Y. Fong, H. Rue, and J. Wakefield. Bayesian inference for generalized linear mixed models. *Biometrics*, 11(3):397–412, 2010.
- G.-A. Fuglstad. Approximating Solutions of Stochastic Differential Equations with Gaussian Markov Random Fields. Pre-Master’s thesis, Department of Mathematical Sciences, NTNU, 2010.
- G.-A. Fuglstad. Spatial Modelling and Inference with SPDE-based GMRFs. Master’s thesis, Department of Mathematical Sciences, NTNU, 2011.
- T. Gneiting, W. Kleiber, and M. Schlather. Matérn cross-covariance functions for multivariate random fields. *Journal of the American Statistical Association*, 105(491):1167–1177, 2010.
- G. H. Golub and C. F. van Loan. *Matrix Computations*. Johns Hopkins University Press, Baltimore, 3rd edition, 1996.
- J. B. Illian, S. H. Sørbye, and H. Rue. A toolbox for fitting complex spatial point process models using integrated Laplace transformation (INLA). *submitted*, 2011.
- F. Lindgren, H. Rue, and J. Lindström. An explicit link between Gaussian fields and Gaussian Markov random fields: The stochastic partial differential equation approach (with discussion). *Journal of the Royal Statistical Society. Series B. Statistical Methodology*, 73(4):423–498, September 2011.
- A. Riebler, L. Held, and H. Rue. Gender-specific differences and the impact of family integration on time trends in age-stratified swiss suicide rates. *Annals of Applied Statistics*, Accepted, 2011.
- J. A. Rozanov. Markov random fields and stochastic partial differential equations. *Math. USSR Sb.*, 32(4):515–534, 1977.

- H. Rue and L. Held. *Gaussian Markov Random Fields: Theory and Applications*, volume 104 of *Monographs on Statistics and Applied Probability*. Chapman & Hall, London, 2005.
- H. Rue and S. Martino. Approximate Bayesian inference for hierarchical Gaussian Markov random fields models. *Journal of Statistical Planning and Inference*, 137(10):3177–3192, 2007. Special Issue: Bayesian Inference for Stochastic Processes.
- H. Rue and H. Tjelmeland. Fitting Gaussian Markov random fields to Gaussian fields. *Scandinavian Journal of Statistics*, 29(1):31–50, 2002.
- H. Rue, S. Martino, and N. Chopin. Approximate Bayesian inference for latent Gaussian models using integrated nested Laplace approximations (with discussion). *Journal of the Royal Statistical Society, Series B*, 71(2):319–392, 2009.
- P. D. Sampson and P. Guttorp. Nonparametric estimation of nonstationary spatial covariance structure. *Journal of the American Statistical Association*, 87(417):108–119, 1992.
- B. Schrödle and L. Held. Spatio-temporal disease mapping using INLA. *Environmetrics*, 22(6):725–734, September 2011.
- D. P. Simpson. *Krylov subspace methods for approximating functions of symmetric positive definite matrices with applications to applied statistics and anomalous diffusion*. PhD thesis, Queensland University of Technology, 2009.
- D. P. Simpson, I. Turner, and A. Pettitt. Fast sampling from a Gaussian Markov random field using Krylov subspace approaches. Technical report, Queensland University of Technology, 2007.
- D. P. Simpson, I. W. Turner, and A. N. Pettitt. Sampling from Gaussian Markov random fields conditioned on linear constraints. *ANZIAM J (CTAC 2006)*, 48:C1041–C1053, July 2008. <http://anziamj.austms.org.au/ojs/index.php/ANZIAMJ/article/view/131> [July 17, 2008].
- D. P. Simpson, J. Illian, F. Lindgren, S. Sørbye, and H. Rue. Computationally efficient inference for log-Gaussian Cox processes. *Submitted*, 2011a.
- D. P. Simpson, F. Lindgren, and H. Rue. In order to make spatial statistics computationally feasible, we need to forget about the covariance function. *Environmetrics*, Accepted, 2011b.
- C. M. Strickland, D. P. Simpson, I. W. Turner, r. Denham, and K. L. Mengersen. Fast Bayesian analysis of spatial dynamic factor models for multitemporal remotely sensed imagery. *Journal of the Royal Statistical Society : Series C*, 60(1):109–124, January 2011.
- P. Whittle. On stationary processes in the plane. *Biometrika*, 41(3/4):434–449, 1954.
- P. Whittle. Stochastic processes in several dimensions. *Bull. Inst. Internat. Statist.*, 40:974–994, 1963.