

Approximating Gaussian random fields by Gaussian Markov random fields: A decade on

Daniel Simpson

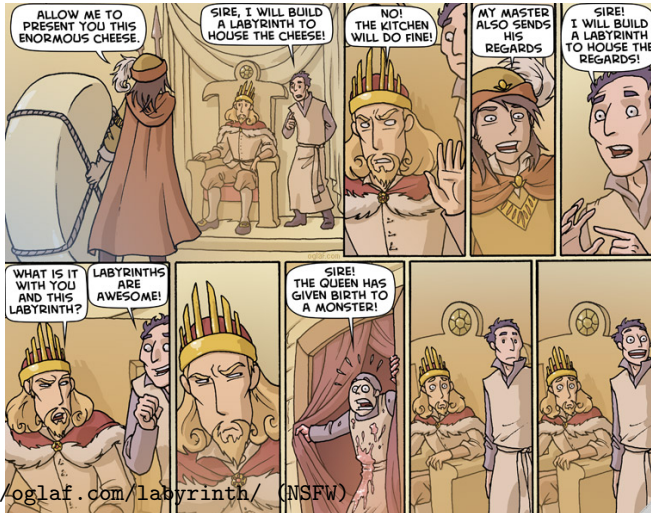
Joint work with Finn Lindgren and Håvard Rue (NTNU)
Ian Turner, Chris Strickland and Tony Pettitt (QUT)

Department of Biosciences, University of Helsinki

May, 2012

The Minotaur
Magic Dance
The Emperor's New Clothes
Real Love
Pass it on

"The minotaur justifies the labyrinth"—Jorge Luis Borges



<http://oglaf.com/labyrinth/> (NSFW)

Outline

- 1 The Minotaur
- 2 Magic Dance (aka The Labyrinth)
- 3 The emperor's new clothes: The fine art of preconditioning
- 4 Real Love
- 5 Pass it on

“Let’s start from the very beginning...”

Defn: Gaussian random fields

A random function $x(s)$ is a GRF iff there is a positive definite function $c(s, s')$ such that, for every finite collection of points $\{s_1, \dots, s_p\}$,

$$\mathbf{x} \equiv (x(s_1), \dots, x(s_p))^T \sim N(\mathbf{0}, \mathbf{\Sigma}),$$

where $\Sigma_{ij} = c(s_i, s_j)$.

- $\mathbf{\Sigma}$ will almost never be sparse.
- It is typically very hard to find families of parameterised positive definite functions.
- Even harder for multivariate (Xianping Hu), non-stationary (Geir-Arne Fuglstad) or spatiotemporal processes.

Second solution

Hi-Defn: Gaussian Random Fields

A Gaussian random field can be characterised by its law, which is a Gaussian measure specified through a mean functional (taken to be zero) and [suppressing technicalities!] its Cameron-Martin space $(H, \langle \cdot, \cdot \rangle_H)$, which is a Hilbert space with a *given* inner product.

- The inner product on H can be specified in terms of an operator Q (i.e. $\langle x, y \rangle_H = \langle x, Qx \rangle$). This is the precision operator and is the “inverse” of the covariance operator.
- Lots of things simplify, especially for hierarchical models!
- Furthermore Q is (pseudo-)differential operator, which we can approximate efficiently through standard wavelet (or Galerkin) methods.

“Being Joan Crawford at 17 was easy...”

Lovely definition, Daniel, but what if I want to actually do something?

- Parameter estimation: Requires the field at N data points, $\Sigma \in \mathbb{R}^{N \times N}$.
- Spatial prediction (Kriging): Requires the field at m points, probably densely through the domain, $\Sigma \in \mathbb{R}^{m \times m}$.
- Joint parameter and spatial estimation: Needs it at both, $\Sigma \in \mathbb{R}^{(N+m) \times (N+m)}$.

Clearly, if N or m is large, this will be *very* expensive. This is the classical problem of spatial statistics.

Det kommer bara leda till nåt ont..

Most of the methods aimed at reducing the “big N problem” in spatial statistics is based on some sort of low-dimensional approximation:

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

where \mathbf{w} is jointly Gaussian and $\phi_i(s)$ are deterministic functions.

There are *lots* of details!

The village green preservation society

With our model in hand, we now need to perform inference.

Direct methods

All Bayesian inference methods (and most Frequentist ones) require a factorisation of the covariance matrix $\Sigma = RR^T$ or the precision matrix $Q = \Sigma^{-1} = LL^T$.

- Making these factorisations computationally feasible is the main aim of modern spatial statistics.
- This was the (computational) state of the art 10 years ago and, with some minor blips, it is still the state of the art.

Outline

- 1 The Minotaur
- 2 Magic Dance (aka The Labyrinth)
- 3 The emperor's new clothes: The fine art of preconditioning
- 4 Real Love
- 5 Pass it on

“Oh those Russians!”—Indirect methods

For spatial problems in applied mathematics, physics, and engineering, direct methods are typically overlooked in favour of *iterative methods*.

- There are a huge variety of *Krylov subspace* methods, the most famous being the Conjugate Gradient method (also GMRES, BiCG-Stab, LSQR,...)
- These *do not* require the matrix, but rather access to matrix-vector products of the form Qx .
- Typically, these methods are exact if you run them for long enough, and they converge superlinearly in subspace size.

Ignition

Sampling from a Gaussian

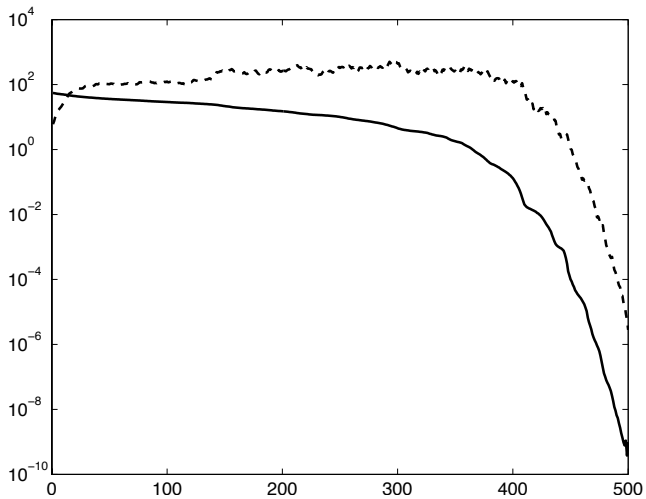
If we have a factorisation of the precision matrix $\mathbf{Q} = \mathbf{L}\mathbf{L}^T$, then it's easy to see that $\mathbf{x} = \mathbf{L}^{-T}\mathbf{z}$, \mathbf{z} i.i.d. standard normal, is a sample from $N(\mathbf{0}, \mathbf{Q}^{-1})$.

- We need a version of $\mathbf{L}^{-T}\mathbf{z}$ that we can compute using only matrix-vector products from \mathbf{Q} .
- It turns out, that we can use Krylov subspace methods to compute $\mathbf{Q}^{-1/2}\mathbf{z}$.
- Convergence is analogous to the conjugate gradient methods for computing $\mathbf{Q}^{-1}\mathbf{z}$

Outline

- 1 The Minotaur
- 2 Magic Dance (aka The Labyrinth)
- 3 The emperor's new clothes: The fine art of preconditioning**
- 4 Real Love
- 5 Pass it on

No easy way down



Stand by your manatee

What went wrong?

- The rate of convergence depends on the condition number of Q , which is $\mathcal{O}(h^4)$.
- This is a standard problem with Krylov subspace methods.
- When solving linear systems, the solution is to *precondition* the system, i.e. find $FF^T \approx Q$ and solve $F^{-1}QF^{-T}y = F^{-1}z$.
- A preconditioner is *optimal* if the condition number remains $\mathcal{O}(1)$ as $h \rightarrow 0$.

Can we precondition the sampling routine?

Ignition (Remix)

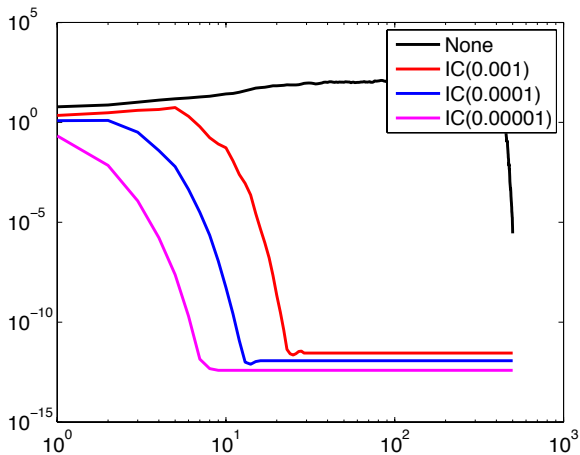
Rather than looking for a transformation that preserves the solution to the linear system, we look for one that gives the same *distribution*.

Preconditioned Sampling

Let Q and $M = FF^T$ be symmetric positive definite matrices. If $y \sim N(\mathbf{0}, (F^{-1}QF^{-T})^{-1})$, then the solution to $F^T x = y$ is a zero-mean Gaussian random vector with precision matrix Q .

- This replaces the problem of sampling from $N(\mathbf{0}, Q^{-1})$ with sampling from $\sim N(\mathbf{0}, (F^{-1}QF^{-T})^{-1})$, which should be better behaved.
- Generic choice of F is the incomplete Cholesky factorisation of Q .

Speed Lab



Breaking Glass

For stationary Gaussian random fields on a regular lattice (on a torus), the precision matrix (and the covariance matrix) is *circulant* and all of the calculations can be done in $\mathcal{O}(n \log n)$ operations using FFTs.

- Any operation involving *data* destroys the circulant structure, leading to precision matrices of the form $\mathbf{Q}_{post} = \mathbf{Q}_{prior} + \mathbf{\Lambda}$.
- This means that good MCMC methods that take into account the second order properties of the likelihood *cannot* be used!
- With Krylov Sampling, you can still do everything. Matrix-vector products with \mathbf{Q}_{post} still cost $\mathcal{O}(n \log n)$ and the number required to reach a prescribed accuracy is $\mathcal{O}(1)$ as $n \rightarrow \infty$.

All the best

Equivalent operator preconditioning

if $\|Q^{1/2}x\|_H$ and $\|M^{1/2}x\|_H$ are equivalent norms on $\mathcal{R}_H(Q^{-1/2})$ and Q_n and M_n are compatible discretisations of Q and M , then

$$c\langle x_n, M_n x_n \rangle \leq \langle x_n, Q_n x_n \rangle \leq C\langle x_n, M_n x_n \rangle$$

and the Krylov sampler required $\mathcal{O}(1)$ iterations to reach the prescribed accuracy ϵ .

- It is not all that difficult to show that if $Q_{prior}^{-1}\Lambda$ is nice (i.e. is Hilbert-Schmidt), then $M = Q_{prior}$ is an optimal preconditioner for Q_{post} .

Simon Smith and his amazing dancing bear

Q is generated using the exponential covariance function on a torus and the diagonals of Λ are $U[0, 10]$.

n	$L = Q^{1/2}$	$L = (Q + \bar{\Lambda})^{1/2}$
32	5	3
64	6	3
128	5	3
256	6	3
512	6	3
1024	6	3

Outline

- 1 The Minotaur
- 2 Magic Dance (aka The Labyrinth)
- 3 The emperor's new clothes: The fine art of preconditioning
- 4 Real Love**
- 5 Pass it on

Less talk, more rock

A classic latent Gaussian process in which the latent field is almost always modelled as block circulant (or block Toeplitz) is the log-Gaussian Cox process model for point pattern data.

$$y_i | \boldsymbol{\eta} \sim \text{Po}(h^2 e^{\eta_i})$$

$$\boldsymbol{\eta} | \boldsymbol{\theta} \sim N(\mathbf{0}, \mathbf{Q}(\boldsymbol{\theta})^{-1})$$

$$\boldsymbol{\theta} \sim \pi(\boldsymbol{\theta}).$$

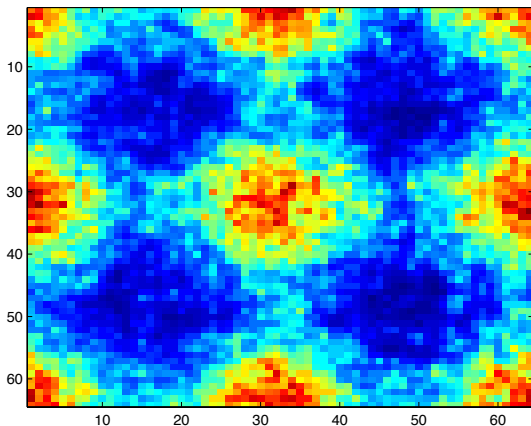
- Here \mathbf{Q} is a circulant matrix that has possibly been extended to include fixed effects.
- The preconditioner for in the case of fixed effects is the same! (block diagonal with the circulant preconditioner and maybe a scaling for the fixed effects components).

“Further Complications.”

There is one problem: computing the acceptance ratio (or the Hamiltonian) requires the computation of the determinant.

- This is trivial if you're using a direct method!
- For indirect methods, $\log\det(\mathbf{A}) = \mathbb{E}(\mathbf{u}^T \log(\mathbf{A})\mathbf{u})$ can be computed (and preconditioned) in the same way as the samples.
- NB: $\log\det(\mathbf{Q}) = \log\det(\mathbf{F}^{-1}\mathbf{Q}\mathbf{L}^{-T}) - 2\log\det(\mathbf{F})$
- Unbiased MC estimate for the log-determinant can be converted to an unbiased estimate for the determinant using results from Physics. The resulting MCMC routine still asymptotically correct, but the cost is a higher variance.

Are you the one I've been waiting for?



Outline

- 1 The Minotaur
- 2 Magic Dance (aka The Labyrinth)
- 3 The emperor's new clothes: The fine art of preconditioning
- 4 Real Love
- 5 Pass it on

I could never take the place of your man

Hopefully, I have convinced you that there are a suite of iterative methods that can be used as efficient replacements for traditional methods.

That being said, these are still methods for *HARD* problems—if existing methods are satisfactory, there is no reason to change!

Now—Later—Soon

- Seeing as we can do everything as long as we have matrix-vector products, we should focus expanding our work to these models! When life is non-stationary, or unequally spaced, we should really spend some quality time with wavelet.
- We probably need more work on MCMC theory. There are results for IEEE floating point arithmetic, but in this case the cut-off function is not the same for each iteration. Is the chain still ergodic? And what does it converge to? How far are we from correct?
- We can actually do all of the things required for INLA. With some programming effort, we may get some interesting results for *HUGE* models!

Playlist: <http://spoti.fi/KubkvQ>