

Computation fundamentals of discrete GMRF representations of continuous domain spatial models

Finn Lindgren

September 23, 2015
v0.2.2

Abstract

The fundamental formulas and algorithms for Bayesian spatial statistics using continuous domain GMRF models are presented. The construction of the models as approximate SPDE solutions is not covered. The results are applicable to non-spatial GMRF models as well, but the important property of observations being *local* and *non-local* is only discussed in the spatial context.

Contents

1	Basic spatial model in GMRF representation	2
1.1	Bayesian hierarchical formulation	2
1.2	Joint distribution	2
1.3	Conditional/posterior distribution (or Kriging)	3
2	Expectations, covariances, and samples	3
2.1	Solving with the posterior precision	3
2.2	Prior and posterior variances	4
2.3	Leave-one-out cross-validation	4
2.4	Sampling from the prior	5
2.5	Sampling from the posterior	5
2.5.1	Direct method	5
2.5.2	Least squares	5
2.5.3	Sampling with conditioning by Kriging	5
3	Likelihood evaluation	6
3.1	Conditional marginal data likelihood	6
3.2	Posterior parameter likelihood and total marginal data likelihood	7
3.3	A note on non-Gaussian data and INLA	7
4	Linear constraints	7
4.1	Non-interacting hard constraints	8
4.2	Conditioning by Kriging	8
4.3	Conditional covariances	9
4.4	Constrained likelihood evaluation	10

1 Basic spatial model in GMRF representation

Basic linear spatial model:

- Spatial basis expansion model $u(\mathbf{s}) = \sum_{k=1}^m \psi_k(\mathbf{s})u_k$, with m basis functions $\psi_k(\mathbf{s})$ with local (compact) support, and GMRF coefficients $(u_1, \dots, u_m) = \mathbf{u} \sim \mathcal{N}(\boldsymbol{\mu}_u, \mathbf{Q}_u^{-1})$, where \mathbf{Q}_u is derived from an SPDE construction with parameters θ_u and $\boldsymbol{\mu}_u$ is usually $\mathbf{0}$.
- Covariate matrix \mathbf{B} , p coefficients $\mathbf{b} \sim \mathcal{N}(\boldsymbol{\mu}_b, \mathbf{Q}_b^{-1})$. Typically, $\boldsymbol{\mu}_b = \mathbf{0}$, and $\mathbf{Q}_b = \tau_b \mathbf{I}_p$ with $\tau_b \approx 0$.
- n observations $\mathbf{y} = (y_1, \dots, y_n)$ at locations \mathbf{s}_i , and additive zero mean Gaussian noise e_i with $(e_1, \dots, e_n) = \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_\epsilon^{-1})$. Typically, $\mathbf{Q}_\epsilon = \tau_\epsilon \mathbf{I}$.

Define \mathbf{A} so that $A_{i,j} = \psi_j(\mathbf{s}_i)$. The observation model is then

$$\mathbf{y} = \mathbf{B}\mathbf{b} + \mathbf{A}\mathbf{u} + \boldsymbol{\epsilon}.$$

Since the covariate effects and the spatial field are jointly Gaussian, it's often practical to join them into a single vector $\mathbf{x} = (\mathbf{u}, \mathbf{b})$, and

$$\begin{aligned} \boldsymbol{\mu}_x &= \begin{bmatrix} \boldsymbol{\mu}_u \\ \boldsymbol{\mu}_b \end{bmatrix}, \\ \mathbf{Q}_x &= \begin{bmatrix} \mathbf{Q}_u & \mathbf{0} \\ \mathbf{0} & \mathbf{Q}_b \end{bmatrix}, \\ \mathbf{A}_x &= [\mathbf{A} \quad \mathbf{B}], \\ \mathbf{y} &= \mathbf{A}_x \mathbf{x} + \boldsymbol{\epsilon}. \end{aligned}$$

1.1 Bayesian hierarchical formulation

Written as a full Bayesian hierarchical model, the model is

$$\begin{aligned} \theta &= (\theta_u, \tau_b, \tau_\epsilon) \sim \pi(\theta), \\ (\mathbf{x} \mid \theta) &\sim \mathcal{N}(\boldsymbol{\mu}_x, \mathbf{Q}_x^{-1}), \\ (\mathbf{y} \mid \theta, \mathbf{x}) &\sim \mathcal{N}(\mathbf{A}_x \mathbf{x}, \mathbf{Q}_\epsilon^{-1}). \end{aligned}$$

It will be important to note that for local basis functions (and at most Markov dependent noise $\boldsymbol{\epsilon}$), $\mathbf{A}^\top \mathbf{Q}_\epsilon \mathbf{A}$ is sparse, and that

$$\mathbf{A}_x^\top \mathbf{Q}_\epsilon \mathbf{A}_x = \begin{bmatrix} \mathbf{A}^\top \mathbf{Q}_\epsilon \mathbf{A} & \mathbf{A}^\top \mathbf{Q}_\epsilon \mathbf{B} \\ \mathbf{B}^\top \mathbf{Q}_\epsilon \mathbf{A} & \mathbf{B}^\top \mathbf{Q}_\epsilon \mathbf{B} \end{bmatrix}$$

is therefore also sparse if p is small.

1.2 Joint distribution

The joint distribution for the observations and the latent variables, (\mathbf{y}, \mathbf{x}) is given by

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{x} \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mathbf{A}_x \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_x \end{bmatrix}, \begin{bmatrix} \mathbf{Q}_\epsilon & -\mathbf{Q}_\epsilon \mathbf{A}_x \\ -\mathbf{A}_x^\top \mathbf{Q}_\epsilon & \mathbf{Q}_x + \mathbf{A}_x^\top \mathbf{Q}_\epsilon \mathbf{A}_x \end{bmatrix}^{-1} \right)$$

1.3 Conditional/posterior distribution (or Kriging)

The conditional distribution for \mathbf{x} given \mathbf{y} and θ is given by

$$\begin{aligned}(\mathbf{x} \mid \theta, \mathbf{y}) &\sim \mathcal{N}(\boldsymbol{\mu}_{x|y}, \mathbf{Q}_{x|y}^{-1}), \\ \mathbf{Q}_{x|y} &= \mathbf{Q}_x + \mathbf{A}_x^\top \mathbf{Q}_\epsilon \mathbf{A}_x, \\ \boldsymbol{\mu}_{x|y} &= \boldsymbol{\mu}_x + \mathbf{Q}_{x|y}^{-1} \mathbf{A}_x^\top \mathbf{Q}_\epsilon (\mathbf{y} - \mathbf{A}_x \boldsymbol{\mu}_x).\end{aligned}$$

The proof is by completing the square in the exponent of the Gaussian conditional density expression

$$\pi(\mathbf{x} \mid \boldsymbol{\theta}, \mathbf{y}) = \frac{\pi(\mathbf{x}, \mathbf{y} \mid \boldsymbol{\theta})}{\pi(\mathbf{y} \mid \boldsymbol{\theta})} \propto \pi(\mathbf{x} \mid \boldsymbol{\theta}) \pi(\mathbf{y} \mid \boldsymbol{\theta}, \mathbf{x}).$$

Note that the elements of $\boldsymbol{\mu}_{x|y}$ are the basis function coefficients and covariate effect estimates in the universal Kriging predictor, and that we've arrived at it without going "the long way around" via the traditional covariance based Kriging equations and matrix inversion lemmas.

2 Expectations, covariances, and samples

Direct calculations and simulation (for models smaller than approximately $m = 10^6$) is typically done with sparse Cholesky decomposition. Reordering methods are required to achieve sparsity (CAMD, nested dissection, etc.). Here, assume for ease of notation that the nodes are already in a suitable order, and that both the prior and posterior precision decompositions are sparse:

$$\begin{aligned}\mathbf{Q}_x &= \mathbf{L}_x \mathbf{L}_x^\top, \\ \mathbf{Q}_{x|y} &= \mathbf{L}_{x|y} \mathbf{L}_{x|y}^\top,\end{aligned}$$

where all \mathbf{L}_\cdot are lower triangular sparse matrices.

See Sections 2.5.2 and 4 for examples of how to handle models where only the prior precision has a sparse Cholesky decomposition. In several instances, the needed adjustments for the use of iterative equation solvers are also mentioned.

2.1 Solving with the posterior precision

Kriging requires solving a linear system

$$\begin{aligned}\mathbf{Q}_{x|y} \mathbf{z} &= \mathbf{w}, \\ \mathbf{z} &= \mathbf{L}_{x|y}^{-\top} \left(\mathbf{L}_{x|y}^{-1} \mathbf{w} \right),\end{aligned}$$

requiring one forward solve and one backward solve. Several systems with the same posterior precision can be solved simultaneously by having a multi-column \mathbf{w} . Often, practical implementations avoid actually transposing the large sparse matrices, effectively doing

$$\mathbf{z} = \left\{ \left(\mathbf{L}_{x|y}^{-1} \mathbf{w} \right)^\top \mathbf{L}_{x|y}^{-1} \right\}^\top$$

instead.

The Kriging predictor (the posterior expectation) is obtained for $\mathbf{w} = \mathbf{A}_x^\top \mathbf{Q}_\epsilon (\mathbf{y} - \mathbf{A}_x \boldsymbol{\mu}_x)$, with $\boldsymbol{\mu}_{x|y} = \boldsymbol{\mu}_x + \mathbf{z}$.

2.2 Prior and posterior variances

Even though the full covariance matrix is too large to be stored, the variances and covariances of neighbouring Markov nodes can be calculated, using an algorithm related to the Cholesky decomposition. In R-INLA, this is available as $\mathbf{S} = \text{inla.qinv}(Q)$, where

$$S_{ij} = \begin{cases} (Q^{-1})_{ij}, & \text{when } Q_{ij} \neq 0 \text{ or } (i, j) \text{ is in the Cholesky infill, and} \\ 0 & \text{otherwise.} \end{cases} \quad (1)$$

This means that the posterior neighbour covariances $\text{Cov}(x_i, x_j \mid \boldsymbol{\theta}, \mathbf{y})$ for neighbours (i, j) are obtained by calling $\text{inla.qinv}()$ with the posterior precision $\mathbf{Q}_{x|y}$.

Provided that the desired posterior predictions $\mathbf{A}_x^{\text{pred}} \mathbf{x}$ have the same sparse connectivity structure as the prior model itself, e.g. when $\mathbf{A}_x^{\text{pred}} = \mathbf{A}_x$, the sparse symmetric \mathbf{S} matrix is sufficient to evaluate marginal posterior predictive variances, as the diagonal of $\mathbf{A}_x^{\text{pred}} \mathbf{S} (\mathbf{A}_x^{\text{pred}})^\top$:

$$\text{Var} \left((\mathbf{A}_x^{\text{pred}} \mathbf{x})_i \right) = \sum_{j=1}^{m+p} (\mathbf{A}_x^{\text{pred}})_{ij} \sum_{k=1}^{m+p} S_{jk} (\mathbf{A}_x^{\text{pred}})_{ik},$$

which in R would typically be implemented as `rowSums(A * (A %*% S))`.

When the desired predictions are dense relative to the model (e.g. when predicting the overall spatial average) solves with $\mathbf{L}_{x|y}$ are required. A small number of predictions (i.e. the number of rows of $\mathbf{A}_x^{\text{pred}}$ is small) can be evaluated simultaneously with

```
LAt <- solve(L, t(A))
colSums(LAt * LAt)
```

Data predictive marginal variances, for $\mathbf{A}_x^{\text{pred}} \mathbf{x} + \epsilon^{\text{pred}}$, are obtained by adding τ_ϵ^{-1} to the variances for $\mathbf{A}_x^{\text{pred}} \mathbf{x}$.

2.3 Leave-one-out cross-validation

In the case of conditionally independent observations, expressions for expectations and variances of leave-one-out cross-validation predictive distributions can be obtained via rank one modification of the full posterior precision matrix, and are obtained “for free” by manipulating the above expressions.

Let $\mathbf{S}_{x|y}$ be the neighbour covariance matrix obtained from (1) applied to $\mathbf{Q}_{x|y}$, and define

$$\begin{aligned} E_i &= (\mathbf{A}_x \boldsymbol{\mu}_{x|y})_i, \\ V_i &= (\mathbf{A}_x)_i \cdot \mathbf{S}_{x|y} (\mathbf{A}_x)_i^\top. \end{aligned}$$

Then the predictive distribution for x_i based on the complete data is

$$((\mathbf{A}_x \mathbf{x})_i \mid \boldsymbol{\theta}, \mathbf{y}) \sim N(E_i, V_i),$$

and the leave-one-out predictive distribution is

$$((\mathbf{A}_x \mathbf{x})_i \mid \boldsymbol{\theta}, \mathbf{y}_j, j \neq i) \sim N(E_{(i)}, V_{(i)}),$$

where

$$\begin{aligned} E_{(i)} &= y_i - \frac{y_i - E_i}{1 - q_i V_i}, \\ V_{(i)} &= \frac{V_i}{1 - q_i V_i}. \end{aligned}$$

where q_i is the conditional precision of observation i , from $\mathbf{Q}_\epsilon = \text{diag}(q_1, \dots, q_n)$.

2.4 Sampling from the prior

Let $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{m+p})$, then

$$\mathbf{x} = \boldsymbol{\mu}_x + \mathbf{L}_x^{-\top} \mathbf{w}$$

is a sample from the prior model for \mathbf{x} . Multiple samples are obtained by repeated forward solves.

2.5 Sampling from the posterior

2.5.1 Direct method

First compute $\boldsymbol{\mu}_{x|y}$ as in Section 2.1, and then let $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{m+p})$. Then

$$\mathbf{x} = \boldsymbol{\mu}_{x|y} + \mathbf{L}_{x|y}^{-\top} \mathbf{w}$$

is a sample from the posterior for \mathbf{x} given \mathbf{y} . Multiple samples are obtained by repeated forward solves.

2.5.2 Least squares

In some models, the posterior precision matrix does not have a sparse Cholesky decomposition, but iterative linear algebra methods may still be used to solve systems with $\mathbf{Q}_{x|y}$. If the prior precision can be decomposed as $\mathbf{Q}_x = \mathbf{H}_x^\top \mathbf{H}_x$ for some sparse matrix \mathbf{H}_x (not necessarily a Cholesky factor), then one approach to sampling is least squares. Given that the posterior expectation has already been calculated using an iterative solution to the equations in Section 1.3, a sample of \mathbf{x} conditionally on \mathbf{y} is given by the solution to

$$\mathbf{Q}_{x|y}(\mathbf{x} - \boldsymbol{\mu}_{x|y}) = [\mathbf{H}_x \quad \mathbf{A}_x^\top \mathbf{L}_\epsilon] \mathbf{w},$$

where $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{m+p+n})$.

2.5.3 Sampling with conditioning by Kriging

The *conditioning by Kriging* approach to sampling is to construct samples from the prior and then *correct* them to become samples from the posterior. For GMRF models, this is less efficient than the direct approach in Section 2.5.1, except when some observations are non-local, which leads to a dense posterior precision.

In its traditional form, the samples are constructed as follows:

$$\begin{aligned} \mathbf{x}^* &= \boldsymbol{\mu}_x + \mathbf{L}_x^{-\top} \mathbf{w}_x, & \mathbf{w}_x &\sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{m+p}), \\ \mathbf{y}^* &= \mathbf{y} + \mathbf{L}_\epsilon^{-\top} \mathbf{w}_y, & \mathbf{w}_y &\sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n), \\ \mathbf{x} &= \mathbf{x}^* - \mathbf{Q}_x^{-1} \mathbf{A}_x^\top \left(\mathbf{A}_x \mathbf{Q}_x^{-1} \mathbf{A}_x^\top + \mathbf{Q}_\epsilon^{-1} \right)^{-1} (\mathbf{A}_x \mathbf{x}^* - \mathbf{y}^*) \end{aligned}$$

where the final expression is the covariance based Kriging equation. For anything but very small n , the inner matrix is dense. Using the Woodbury identity, the expression can be rewritten to recover the precision based equation from Section 1.3,

$$\begin{aligned} \mathbf{x} &= \mathbf{x}^* + \left(\mathbf{Q}_x + \mathbf{A}_x^\top \mathbf{Q}_\epsilon \mathbf{A}_x \right)^{-1} \mathbf{A}_x^\top \mathbf{Q}_\epsilon (\mathbf{y}^* - \mathbf{A}_x \mathbf{x}^*) \\ &= \mathbf{x}^* + \mathbf{Q}_{x|y}^{-1} \mathbf{A}_x^\top \mathbf{Q}_\epsilon (\mathbf{y}^* - \mathbf{A}_x \mathbf{x}^*), \end{aligned}$$

which is calculated as before, with forward and backward solves on $\mathbf{L}_{x|y}$.

If N samples, as well as the posterior mean, are desired, conditioning by Kriging requires N solves with \mathbf{L}_x and \mathbf{L}_ϵ (to obtain \mathbf{x}_0 and \mathbf{y}_0), and $2 + 2N$ solves with $\mathbf{L}_{x|y}$. The direct method only requires $2 + N$ solves with $\mathbf{L}_{x|y}$, and no solves with \mathbf{L}_x or \mathbf{L}_ϵ , which means that it is nearly always preferable to conditioning by Kriging. The exception is non-local observations, which can be handled separately as soft constraints, see Section 4.

3 Likelihood evaluation

The expensive part of evaluating likelihoods is calculating precision determinants, but they come at no extra cost given the Cholesky factors:

$$\begin{aligned}\log \pi(\mathbf{x} | \boldsymbol{\theta}) &= -\frac{m+p}{2} \log(2\pi) + \frac{1}{2} \log \det(\mathbf{Q}_x) - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_x)^\top \mathbf{Q}_x (\mathbf{x} - \boldsymbol{\mu}_x), \\ \log \det(\mathbf{Q}_x) &= 2 \sum_{j=1}^{m+p} \log(\mathbf{L}_x)_{jj}.\end{aligned}$$

Corresponding expressions hold for $\pi(\mathbf{y} | \boldsymbol{\theta}, \mathbf{x})$ and $\pi(\mathbf{x} | \boldsymbol{\theta}, \mathbf{y})$, using \mathbf{L}_ϵ and $\mathbf{L}_{x|y}$.

Marginal data likelihoods require more effort.

3.1 Conditional marginal data likelihood

For any arbitrary \mathbf{x}^* ,

$$\begin{aligned}\pi(\mathbf{y} | \boldsymbol{\theta}) &= \frac{\pi(\boldsymbol{\theta}, \mathbf{y})}{\pi(\boldsymbol{\theta})} = \frac{\pi(\boldsymbol{\theta}, \mathbf{y})\pi(\mathbf{x} | \boldsymbol{\theta}, \mathbf{y})}{\pi(\boldsymbol{\theta})\pi(\mathbf{x} | \boldsymbol{\theta}, \mathbf{y})} \Big|_{\mathbf{x}=\mathbf{x}^*} = \frac{\pi(\boldsymbol{\theta}, \mathbf{y}, \mathbf{x})}{\pi(\boldsymbol{\theta})\pi(\mathbf{x} | \boldsymbol{\theta}, \mathbf{y})} \Big|_{\mathbf{x}=\mathbf{x}^*} \\ &= \frac{\pi(\mathbf{x} | \boldsymbol{\theta})\pi(\mathbf{y} | \boldsymbol{\theta}, \mathbf{x})}{\pi(\mathbf{x} | \boldsymbol{\theta}, \mathbf{y})} \Big|_{\mathbf{x}=\mathbf{x}^*}.\end{aligned}$$

In practice, this is evaluated for $\mathbf{x}^* = \boldsymbol{\mu}_{x|y}$, which gives better numerical stability than $\mathbf{x}^* = \mathbf{0}$. Written explicitly,

$$\begin{aligned}\log \pi(\mathbf{y} | \boldsymbol{\theta}) &= -\frac{n}{2} \log(2\pi) + \frac{1}{2} \log \det(\mathbf{Q}_x) + \frac{1}{2} \log \det(\mathbf{Q}_\epsilon) - \frac{1}{2} \log \det(\mathbf{Q}_{x|y}) \\ &\quad - \frac{1}{2} (\boldsymbol{\mu}_{x|y} - \boldsymbol{\mu}_x)^\top \mathbf{Q}_x (\boldsymbol{\mu}_{x|y} - \boldsymbol{\mu}_x) - \frac{1}{2} (\mathbf{y} - \mathbf{A}_x \boldsymbol{\mu}_{x|y})^\top \mathbf{Q}_\epsilon (\mathbf{y} - \mathbf{A}_x \boldsymbol{\mu}_{x|y}),\end{aligned}$$

where an anticipated third quadratic form, from $\pi(\mathbf{x} | \boldsymbol{\theta}, \mathbf{y})$, has been eliminated due to the choice of \mathbf{x}^* . Comparing with a standard marginal Gaussian log-likelihood expression, the log-determinant and quadratic form can be identified,

$$\begin{aligned}\log \det(\mathbf{Q}_y) &= \log \det(\mathbf{Q}_x) + \log \det(\mathbf{Q}_\epsilon) - \log \det(\mathbf{Q}_{x|y}) \\ (\mathbf{y} - \mathbf{A}_x \boldsymbol{\mu}_x)^\top \mathbf{Q}_y (\mathbf{y} - \mathbf{A}_x \boldsymbol{\mu}_x) &= (\boldsymbol{\mu}_{x|y} - \boldsymbol{\mu}_x)^\top \mathbf{Q}_x (\boldsymbol{\mu}_{x|y} - \boldsymbol{\mu}_x) + (\mathbf{y} - \mathbf{A}_x \boldsymbol{\mu}_{x|y})^\top \mathbf{Q}_\epsilon (\mathbf{y} - \mathbf{A}_x \boldsymbol{\mu}_{x|y}),\end{aligned}$$

even though the marginal precision matrix itself, \mathbf{Q}_y , is dense and cannot be stored.

3.2 Posterior parameter likelihood and total marginal data likelihood

Using the result for the conditional marginal data likelihood, we can express the posterior parameter likelihood as

$$\pi(\boldsymbol{\theta} \mid \mathbf{y}) = \frac{\pi(\boldsymbol{\theta}, \mathbf{y})}{\pi(\mathbf{y})} = \frac{\pi(\boldsymbol{\theta})\pi(\mathbf{y} \mid \boldsymbol{\theta})}{\pi(\mathbf{y})},$$

where the only unknown is the total marginal data likelihood $\pi(\mathbf{y})$, which can be approximated with numerical integration over $\boldsymbol{\theta}$,

$$\pi(\mathbf{y}) = \int \pi(\boldsymbol{\theta})\pi(\mathbf{y} \mid \boldsymbol{\theta}) \, d\boldsymbol{\theta}.$$

3.3 A note on non-Gaussian data and INLA

For non-Gaussian data, Laplace approximation can be used, which in essence involves replacing the non-Gaussian conditional posterior $\pi(\mathbf{x} \mid \boldsymbol{\theta}, \mathbf{y})$ with a Gaussian approximation in the expression for $\pi(\mathbf{y} \mid \boldsymbol{\theta})$, taking \mathbf{x}^* as the mode of $\pi(\mathbf{x} \mid \boldsymbol{\theta}, \mathbf{y})$.

The INLA implementation finds the mode of $\pi(\boldsymbol{\theta} \mid \mathbf{y})$ using numerical optimisation, and further approximates the marginal posteriors for the latent variables x_k , $\pi(x_k \mid \mathbf{y})$, with numerical integration of

$$\pi(x_k \mid \mathbf{y}) = \int \pi(x_k \mid \boldsymbol{\theta}, \mathbf{y})\pi(\boldsymbol{\theta} \mid \mathbf{y}) \, d\boldsymbol{\theta},$$

where a further laplace approximation of $\pi(x_k \mid \boldsymbol{\theta}, \mathbf{y})$ is made (as well as skewness corrections not discussed here),

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

This combination of Laplace approximations is the ‘‘nested Laplace’’ part of INLA. In the case of Gaussian data, the only approximation in the method is the numerical integration itself.

4 Linear constraints

There are two main classes of linear constraints:

1. *Hard constraints*, which fall into two subclasses,
 - (a) *Interacting hard constraints*, where linear combinations of nodes are constrained,
 - (b) *Non-interacting hard constraints*, where each constraint acts only on a single node, in effect specifying those nodes explicitly,
2. *Soft constraints*, where linear combinations are specified with some uncertainty, equivalent to posterior conditioning on observations of these linear combinations. Observations are treated under this framework when they break the Markov structure in the posterior distribution, e.g. by affecting all the nodes in the field directly.

All these cases can be written on matrix form as

$$\mathbf{A}_c \mathbf{x} + \boldsymbol{\epsilon}_c = \mathbf{e}_c,$$

where $\epsilon_c = \mathbf{0}$ in the hard constraint cases, and

$$\epsilon_c \sim N(\mathbf{0}, \mathbf{Q}_c^{-1}),$$

in the soft constraint case. The aim is to sample from and evaluate properties of the conditional distribution for $(\mathbf{x} \mid \mathbf{e}_c)$, when

$$\begin{aligned} \mathbf{x} &\sim N(\boldsymbol{\mu}, \mathbf{Q}^{-1}), \\ (\mathbf{e}_c \mid \mathbf{x}) &\sim N(\mathbf{A}_c \mathbf{x}, \mathbf{Q}_c^{-1}), \end{aligned}$$

where \mathbf{x} comes from a known GMRF model, typically either a prior or a posterior distribution.

Let r denote the number of constraints, i.e. the number of rows in \mathbf{A}_c and the number of elements in \mathbf{e}_c .

4.1 Non-interacting hard constraints

Non-interacting hard constraints are easiest to handle, as the constrained nodes can be completely removed from all calculations in a preprocessing step, as they are in effect known deterministic values, and the conditional precision for the remaining nodes can be trivially extracted as a block from the joint precision matrix. Without loss of practical generality, we can assume that all scaling is done in \mathbf{e}_c , so that each row of \mathbf{A}_c has a single non-zero entry that is equal to 1.

For ease of notation, assume that the nodes are ordered so that the unconstrained nodes U are followed by the constrained, C , so that the joint distribution is

$$\begin{bmatrix} \mathbf{x}_U \\ \mathbf{x}_C \end{bmatrix} \sim N \left(\begin{bmatrix} \boldsymbol{\mu}_U \\ \boldsymbol{\mu}_C \end{bmatrix}, \begin{bmatrix} \mathbf{Q}_{UU} & \mathbf{Q}_{UC} \\ \mathbf{Q}_{CU} & \mathbf{Q}_{CC} \end{bmatrix}^{-1} \right).$$

The conditional distribution for the unconstrained nodes is given by

$$\begin{aligned} (\mathbf{x}_U \mid \mathbf{x}_C = \mathbf{e}_c) &\sim N(\boldsymbol{\mu}_{U|C}, \mathbf{Q}_{UU}^{-1}), \\ \boldsymbol{\mu}_{U|C} &= \boldsymbol{\mu}_U - \mathbf{Q}_{UU}^{-1} \mathbf{Q}_{UC} (\mathbf{e}_c - \boldsymbol{\mu}_C). \end{aligned}$$

The constrained mean and constrained samples are obtained with the same forward and backward solving techniques as in Section 1.3 and Section 2.5.1, using the Cholesky decomposition of \mathbf{Q}_{UU} . The joint constrained distribution can be written as a degenerate Normal distribution,

$$\left(\begin{bmatrix} \mathbf{x}_U \\ \mathbf{x}_C \end{bmatrix} \mid \mathbf{x}_C = \mathbf{e}_c \right) \sim N \left(\begin{bmatrix} \boldsymbol{\mu}_{U|C} \\ \mathbf{e}_c \end{bmatrix}, \begin{bmatrix} \mathbf{Q}_{UU}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \right).$$

4.2 Conditioning by Kriging

Conditioning on local observations preserves the Markov structure, and Markov breaking conditioning is therefore done as a final step, when computing posterior means, and sampling from the posterior as well as from the prior. Using the conditioning by Kriging introduced in Section 2.5.3, conditioning on linear constraints or on non-local observations can be written in a common framework. Let $\mathbf{Q}_c = \mathbf{L}_c \mathbf{L}_c^\top$ be the Cholesky decomposition of \mathbf{Q}_c in the soft constraint case, and for notational convenience, define $\mathbf{Q}_c^{-1} = \mathbf{0}$ and $\mathbf{L}_c^{-1} = \mathbf{0}$ for hard constraints.

The method proceeds by constructing a sample from the unconstrained model, and then correcting for the constraints,

$$\begin{aligned} \mathbf{x}^* &= \boldsymbol{\mu} + \mathbf{L}^{-\top} \mathbf{w}_x, & \mathbf{w}_x &\sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{m+p}), \\ \mathbf{e}_c^* &= \mathbf{e}_c + \mathbf{L}_c^{-\top} \mathbf{w}_c, & \mathbf{w}_c &\sim \mathcal{N}(\mathbf{0}, \mathbf{I}_r), \\ \mathbf{x} &= \mathbf{x}^* - \mathbf{Q}^{-1} \mathbf{A}_c^\top \left(\mathbf{A}_c \mathbf{Q}^{-1} \mathbf{A}_c^\top + \mathbf{Q}_c^{-1} \right)^{-1} (\mathbf{A}_c \mathbf{x}^* - \mathbf{e}_c^*). \end{aligned}$$

In order to obtain the conditional mean instead of a sample, set $\mathbf{x}^* = \boldsymbol{\mu}$ and $\mathbf{e}_c^* = \mathbf{e}_c$. The dense but small inner $r \times r$ matrix can be evaluated using a single forward solve with \mathbf{L} ,

$$\mathbf{A}_c \mathbf{Q}^{-1} \mathbf{A}_c^\top + \mathbf{Q}_c^{-1} = \left(\mathbf{L}^{-1} \mathbf{A}_c^\top \right)^\top \mathbf{L}^{-1} \mathbf{A}_c^\top + \mathbf{Q}_c^{-1}.$$

Defining $\tilde{\mathbf{A}}_c^\top = \mathbf{L}^{-1} \mathbf{A}_c^\top$, the full expression becomes

$$\mathbf{x} = \mathbf{x}^* - \mathbf{L}^{-\top} \tilde{\mathbf{A}}_c^\top \left(\tilde{\mathbf{A}}_c \tilde{\mathbf{A}}_c^\top + \mathbf{Q}_c^{-1} \right)^{-1} (\mathbf{A}_c \mathbf{x}^* - \mathbf{e}_c^*),$$

where the remaining inner $r \times r$ solve is evaluated using, e.g., Cholesky decomposition, and the result fed into a final backward solve with \mathbf{L}^\top .

In order to avoid unnecessary recalculations for multiple samples, further precomputation can be used. A fully Cholesky based approach that also simplifies later calculations is

$$\begin{aligned} \tilde{\mathbf{L}}_c \tilde{\mathbf{L}}_c^\top &= \tilde{\mathbf{A}}_c \tilde{\mathbf{A}}_c^\top + \mathbf{L}_c^\top \mathbf{L}_c^{-1}, \\ \mathbf{B} &= \mathbf{L}^{-\top} \tilde{\mathbf{A}}_c^\top \tilde{\mathbf{L}}_c^{-\top}, \\ \mathbf{x} &= \mathbf{x}^* - \mathbf{B} \tilde{\mathbf{L}}_c^{-1} (\mathbf{A}_c \mathbf{x}^* - \mathbf{e}_c^*), \end{aligned}$$

where $\tilde{\mathbf{L}}_c$ is a dense $r \times r$ Cholesky factor, and \mathbf{B} is dense, $(m+p) \times r$. For models that require iterative solves for \mathbf{Q} instead of Cholesky decomposition, the pre-calculations are modified slightly,

$$\begin{aligned} \tilde{\tilde{\mathbf{A}}}_c^\top &= \mathbf{Q}^{-1} \mathbf{A}_c^\top \\ \tilde{\mathbf{L}}_c \tilde{\mathbf{L}}_c^\top &= \tilde{\tilde{\mathbf{A}}}_c \tilde{\tilde{\mathbf{A}}}_c^\top + \mathbf{L}_c^\top \mathbf{L}_c^{-1}, \\ \mathbf{B} &= \tilde{\tilde{\mathbf{A}}}_c^\top \tilde{\mathbf{L}}_c^{-\top}, \end{aligned}$$

but the final step for constructing \mathbf{x} remains unchanged.

4.3 Conditional covariances

From the conditioning by Kriging algorithm, we obtain

$$\text{Cov}(\mathbf{x}, \mathbf{x} \mid \mathbf{e}_c) = \mathbf{Q}^{-1} - \mathbf{B} \mathbf{B}^\top,$$

which can be combined with the neighbour covariance algorithm from Section 2.2 to yield the conditional covariances between nodes that are neighbours in the unconstrained model,

$$(\mathbf{S}_c)_{ij} = \begin{cases} \mathbf{S}_{ij} - \sum_k \mathbf{B}_{ik} \mathbf{B}_{jk}, & \text{when } Q_{ij} \neq 0 \text{ or } (i, j) \text{ is in the Cholesky infill, and} \\ 0 & \text{otherwise.} \end{cases}$$

4.4 Constrained likelihood evaluation

The constrained likelihood is given by

$$\pi(\mathbf{x} \mid \mathbf{e}_c) = \frac{\pi(\mathbf{x})\pi(\mathbf{e}_c \mid \mathbf{x})}{\pi(\mathbf{e}_c)},$$

where $\pi(\mathbf{e}_c \mid \mathbf{x})$ is a degenerate density in the hard constraint case.

The marginal distributions for \mathbf{x} and \mathbf{e}_c are

$$\begin{aligned}\mathbf{x} &\sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{Q}^{-1}), \\ \mathbf{e}_c &\sim \mathcal{N}(\mathbf{A}_c\boldsymbol{\mu}, \tilde{\mathbf{L}}_c\tilde{\mathbf{L}}_c^\top),\end{aligned}$$

and the densities can be evaluated with the usual approach, with

$$\begin{aligned}\log \pi(\mathbf{x}) &= -\frac{m+p}{2} \log(2\pi) + \log \det(\mathbf{L}) - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \mathbf{Q}(\mathbf{x} - \boldsymbol{\mu}), \\ \log \pi(\mathbf{e}_c) &= -\frac{r}{2} \log(2\pi) - \log \det(\tilde{\mathbf{L}}_c) - \frac{1}{2}(\mathbf{e}_c - \mathbf{A}_c\boldsymbol{\mu})^\top \tilde{\mathbf{L}}_c^{-\top} \tilde{\mathbf{L}}_c^{-1}(\mathbf{e}_c - \mathbf{A}_c\boldsymbol{\mu}).\end{aligned}$$

In the hard constraint case, Section 2.3.3 of (Rue and Held, 2005) provides the degenerate conditional density through

$$\log \pi(\mathbf{e}_c \mid \mathbf{x}) = \log \pi(\mathbf{A}_c\mathbf{x} \mid \mathbf{x}) = -\frac{1}{2} \log \det(\mathbf{A}_c\mathbf{A}_c^\top)$$

which can be evaluated using the Cholesky decomposition of the $r \times r$ matrix $\mathbf{A}_c\mathbf{A}_c^\top$. In the soft constraint case,

$$\log \pi(\mathbf{e}_c \mid \mathbf{x}) = -\frac{r}{2} \log(2\pi) + \log \det(\mathbf{L}_c) - \frac{1}{2}(\mathbf{e}_c - \mathbf{A}_c\mathbf{x})^\top \mathbf{Q}_c(\mathbf{e}_c - \mathbf{A}_c\mathbf{x}).$$

References

Rue, H. and L. Held (2005). *Gaussian Markov Random Fields; Theory and Applications*. Chapman & Hall/CRC.