

# Inference for Gaussian Markov random fields and the INLA approach

Johan Lindström<sup>1</sup>

<sup>1</sup>Centre for Mathematical Sciences  
Lund University

Lund 2016-02-11



## Whom I “borrowed” slides from

David Bolin’s PhD course on GMRFs from 2015.  
[www.math.chalmers.se/~bodavid/GMRF2015/](http://www.math.chalmers.se/~bodavid/GMRF2015/)

Various presentations from the INLA group.  
[www.r-inla.org/](http://www.r-inla.org/)

Daniel Simpson’s Helsinki presentation.  
[www.math.ntnu.no/~danesi/slides.pdf](http://www.math.ntnu.no/~danesi/slides.pdf)

### Spatial interpolation

Given observations at some locations,  $Y(\mathbf{s}_i)$ ,  $i = 1 \dots n$   
we want to make statements about the value at unobserved location(s),  $X(\mathbf{s})$ .

The general model formulation is

$$Y(\mathbf{s}_i) | \mathbf{X}, \vartheta \sim F(X(\mathbf{s}_i), \vartheta)$$

$$\mathbf{X} | \vartheta \sim \text{Spatial field}(\vartheta)$$

often we assume  $\mathbf{X}$  to be a Gaussian field.

### Bayesian hierarchical modelling using GMRF

Data model,  $p(\mathbf{y}|\boldsymbol{\eta}, \boldsymbol{\vartheta})$ : Describing how **observations** arise assuming a **known latent field  $\boldsymbol{\eta}$** .

Latent model,  $p(\boldsymbol{\eta}|\boldsymbol{\vartheta})$ : Describing how the **latent field** behaves.

$$\boldsymbol{\eta} = \mathbf{A}\mathbf{x} + \mathbf{B}\boldsymbol{\beta} \quad \mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}^{-1}(\boldsymbol{\vartheta}))$$

Parameters,  $p(\boldsymbol{\vartheta})$ : **Prior knowledge** of the parameters.

### Bayesian hierarchical modelling using GMRF

In general these models can be seen as GLMs with

Data model The data is assumed to be **conditionally independent**

$$p(\mathbf{y}|\boldsymbol{\eta}, \boldsymbol{\vartheta}) = \prod_i p(y_i|\eta_i, \boldsymbol{\vartheta})$$

Latent model The latent model is often constructed as a linear combination of **spatial effects** and **regression mean**

$$\boldsymbol{\eta} = \mathbf{A}\mathbf{x} + \mathbf{B}\boldsymbol{\beta} = \begin{bmatrix} \mathbf{A} & \mathbf{B} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \boldsymbol{\beta} \end{bmatrix}$$

$$\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}^{-1}(\boldsymbol{\vartheta})) \quad \boldsymbol{\beta} \sim \mathcal{N}(\mathbf{0}, q_\beta^{-1} \cdot \mathbf{I})$$

$$\begin{bmatrix} \mathbf{x} \\ \boldsymbol{\beta} \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \mathbf{Q}(\boldsymbol{\vartheta}) & \mathbf{0} \\ \mathbf{0} & q_\beta \cdot \mathbf{I} \end{bmatrix}^{-1}\right)$$

Here  $q_\beta \cdot \mathbf{I}$  represents a vague prior on  $\boldsymbol{\beta}$ .

### Inference

Given a Bayesian hierarchical model we are interested in

Posteriors for the parameters

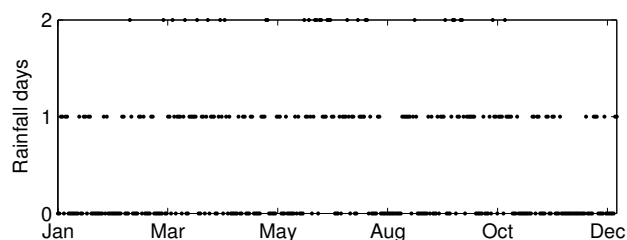
$$p(\boldsymbol{\vartheta}|\mathbf{y}) \propto p(\mathbf{y}|\boldsymbol{\vartheta})p(\boldsymbol{\vartheta})$$

Posteriors for the latent field

$$p(\boldsymbol{\eta}|\mathbf{y}) \propto \int p(\boldsymbol{\eta}|\mathbf{y}, \boldsymbol{\vartheta})p(\boldsymbol{\vartheta}|\mathbf{y}) d\boldsymbol{\vartheta}$$

## Example: Tokyo rainfall (Kitagawa, 1987)

- ▶ During 1983 and 1984 each day with with more than 1 mm rainfall in Tokyo was recorded.
- ▶ We want to estimate the probability of rainfall,  $p_t$ , for a given calendar day,  $t = 1, \dots, 366$ .



## Model for the Tokyo rainfall

The Tokyo data can be modelled using binomial-observations of a latent Gaussian process:

$$\begin{aligned}
 y_t | x_t &\sim \text{Bin}(N_t, p_t) & p_t &= \frac{e^{x_t}}{1 + e^{x_t}} \\
 \mathbf{x} | \tau &\sim \mathbf{N}(\mathbf{0}, \tau^{-1} \mathbf{Q}^{-1}) \\
 \tau &\sim \Gamma(a, b) & a = 1 & \quad b = 10^{-5}
 \end{aligned}$$

To account for the seasonality, the latent process is modelled as a **cyclic random walk**.

For this model,  $\eta = \mathbf{x}$  and we use  $\mathbf{x}$  for the rest of the example.

## Latent field

The latent field is taken as the solution, on  $t = 1, 2, \dots, 366$ , to the SPDE

$$(\kappa^2 - \Delta)^{\alpha/2} \mathbf{x}(t) = \frac{1}{\sqrt{\tau}} \mathcal{W}(t).$$

with cyclical boundaries,  $x(1) \approx x(T)$ , and  $\kappa^2 = 0$ .

The resulting  $\mathbf{G}$  and  $\mathbf{C}$  matrices are

$$\mathbf{G} = \begin{bmatrix} 2 & -1 & 0 & \dots & 0 & -1 \\ -1 & 2 & -1 & 0 & \dots & \\ 0 & -1 & 2 & -1 & 0 & \dots \\ & & & \ddots & & \\ -1 & 0 & \dots & 0 & -1 & 2 \end{bmatrix} \quad \mathbf{C} = \mathbf{I}$$

giving  $\mathbf{Q}_{\alpha=1} = \mathbf{G}$  (crw1) and  $\mathbf{Q}_{\alpha=2} = \mathbf{G}\mathbf{G}$  (crw2).

## Naïve Inference

Due to the markov structure of the model we have conditional posteriors

$$p(\mathbf{x}_t | \mathbf{x}_{-t}, \mathbf{y}_t) \propto (1 + e^{\mathbf{x}_t})^{-N_t} \exp\left(\mathbf{y}_t \mathbf{x}_t - \frac{\tau Q_{tt}}{2} \left(\mathbf{x}_t + \frac{\sum_{s \neq t} Q_{ts} \mathbf{x}_s}{Q_{tt}}\right)^2\right)$$

and

$$\tau | \mathbf{x} \sim \Gamma\left(\frac{N_x - 1}{2} + a, \frac{\mathbf{x}^\top \mathbf{Q} \mathbf{x}}{2} + b\right)$$

The posterior  $p(\mathbf{x}_t | \mathbf{x}_{-t}, \mathbf{y}_t)$  is **log-concave** and can be sampled exactly using **adaptive rejection sampling** (Gilks and Wild, 1992).

Given these conditionals, the posterior  $p(\mathbf{x}, \tau | \mathbf{y})$  can be estimated using Gibbs-sampling.

## Naïve Inference — Auxiliary variables

As an alternative to the **adaptive rejection sampling** the MCMC estimation can be simplified through the use of **auxiliary variables**

For **student's t** we have:

If  $X \sim N(0, 1)$  and  $Z \sim \Gamma(\nu/2, \nu/2)$  then

$$X\sqrt{Z} \sim t(\nu)$$

## Auxiliary variables for Binomials

For binary models with probit link (Albert and Chib, 1993) the model

$$y_t \sim \text{Bin}(1, p_t) \quad p_t = \Phi(\mathbf{x}_t) \quad \text{probit link}$$

is equivalent to

$$\begin{aligned} w_t &= \mathbf{x}_t + \varepsilon_t, & \varepsilon_t &\sim N(0, 1) \\ y_t &= \begin{cases} 1 & \text{if } w_t > 0 \\ 0 & \text{otherwise} \end{cases} \end{aligned}$$

and the problem reduces to sampling from a normal,  $\mathbf{x} | \mathbf{w} \sim N$  and a truncated normal  $w_t | x_t, y_t \sim N_{\text{trunc}}$ .

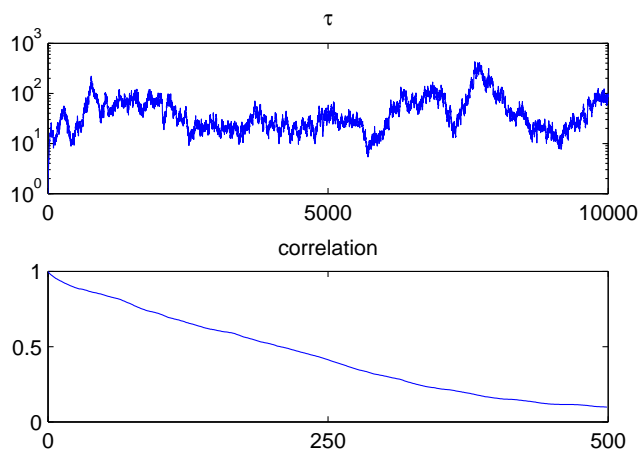
Auxiliary variables for logit links also exist (Held and Holmes, 2006).

## Results — Single site Gibbs

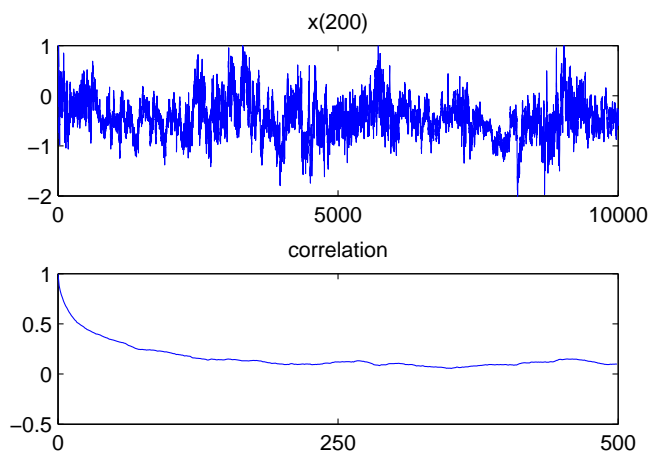
The single site Gibbs-sampler is easy to implement but has several drawbacks:

- ▶ **Very long run time** (hard to vectorise code)
- ▶ **Very slow mixing** of the chains.
- ▶ Mixing is poor between  $\tau$  and the latent field  $\mathbf{x}$
- ▶ And mixing is poor within the latent field  $\mathbf{x}$

## $\tau$ — Single site Gibbs



## $x_{200}$ — Single site Gibbs



## Comments

For inference in latent fields it is well known (Knorr-Held and Rue, 2002) that **blocking improves mixing**. Using the auxiliary variables allows us to block the updates as

1.  $\mathbf{x}|\mathbf{w} \sim \mathbf{N}$  can be sampled jointly
2.  $\mathbf{w}_t|\mathbf{x}_t, \mathbf{y}_t \sim \mathbf{N}_{\text{trunc}}$  are conditionally independent and can be sampled jointly.
3.  $\tau|\mathbf{x} \sim \Gamma\left(\frac{M_{\mathbf{x}}-1}{2} + a, \frac{\mathbf{x}^\top \mathbf{Q} \mathbf{x}}{2} + b\right)$

This will improve mixing but still leaves poor mixing between  $\tau$  and  $\mathbf{x}$ .

MCMC is “right” at the limit (infinite run time). In practice, **fast approximate solutions** are better than exact slow solutions.

## Inference for parameters and latent field

Given a model with **conditionally independent** observations

$$p(\mathbf{y}|\mathbf{x}, \boldsymbol{\vartheta}) = \prod_i p(y_i|x_i, \boldsymbol{\vartheta})$$

from a latent GMRF  $\mathbf{x} \sim \mathbf{N}(\mathbf{0}, \mathbf{Q}^{-1}\boldsymbol{\vartheta})$  we want:

Posteriors for the parameters

$$p(\boldsymbol{\vartheta}|\mathbf{y}) \propto p(\mathbf{y}|\boldsymbol{\vartheta}) p(\boldsymbol{\vartheta}) = \int p(\mathbf{y}|\mathbf{x}, \boldsymbol{\vartheta}) p(\mathbf{x}|\boldsymbol{\vartheta}) d\mathbf{x} p(\boldsymbol{\vartheta})$$

Maximum posterior estimates

$$\arg \max_{\boldsymbol{\vartheta}} p(\boldsymbol{\vartheta}|\mathbf{y})$$

Posteriors for the latent field

$$p(\mathbf{x}|\mathbf{y}) \propto \int p(\mathbf{x}|\mathbf{y}, \boldsymbol{\vartheta}) p(\boldsymbol{\vartheta}|\mathbf{y}) d\boldsymbol{\vartheta}$$

## Computing the posterior — A “trick”

To avoid explicitly computing the integral

$$\int p(\mathbf{y}, \mathbf{x}|\boldsymbol{\vartheta}) d\mathbf{x} = \int p(\mathbf{y}|\mathbf{x}, \boldsymbol{\vartheta}) p(\mathbf{x}|\boldsymbol{\vartheta}) d\mathbf{x}$$

we note that conditional distributions provide the equality

$$p(\mathbf{y}|\mathbf{x}, \boldsymbol{\vartheta}) p(\mathbf{x}|\boldsymbol{\vartheta}) = p(\mathbf{y}, \mathbf{x}|\boldsymbol{\vartheta}) = p(\mathbf{x}|\mathbf{y}, \boldsymbol{\vartheta}) p(\mathbf{y}|\boldsymbol{\vartheta})$$

This gives

$$p(\boldsymbol{\vartheta}|\mathbf{y}) \propto \underbrace{\frac{p(\mathbf{y}|\mathbf{x}, \boldsymbol{\vartheta}) p(\mathbf{x}|\boldsymbol{\vartheta})}{p(\mathbf{x}|\mathbf{y}, \boldsymbol{\vartheta})}}_{p(\mathbf{y}|\boldsymbol{\vartheta})} \cdot p(\boldsymbol{\vartheta}) \quad \text{for any } \mathbf{x}.$$

With  $p(\mathbf{x}|\mathbf{y}, \boldsymbol{\vartheta})$  being the posterior of  $\mathbf{x}$  given  $\mathbf{y}$ .

## Can we compute $p(\mathbf{x}|\mathbf{y}, \vartheta)$ ?

For Gaussian observations we have

$$\begin{aligned}\mathbf{y}|\mathbf{A}\mathbf{x}, \vartheta &\sim N(\mathbf{A}\mathbf{x}, \mathbf{Q}_\varepsilon^{-1}(\vartheta)), \\ \mathbf{x}|\vartheta &\sim N(\mathbf{0}, \mathbf{Q}^{-1}(\vartheta)),\end{aligned}$$

with posteriors

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

where

$$\begin{aligned}\boldsymbol{\mu}_{\mathbf{x}|\mathbf{y}}(\vartheta) &= \mathbf{Q}_{\mathbf{x}|\mathbf{y}}^{-1}(\vartheta) \mathbf{A}^\top \mathbf{Q}_\varepsilon(\vartheta) \mathbf{y}, \\ \mathbf{Q}_{\mathbf{x}|\mathbf{y}}(\vartheta) &= \mathbf{Q}(\vartheta) + \mathbf{A}^\top \mathbf{Q}_\varepsilon(\vartheta) \mathbf{A}.\end{aligned}$$

## Parameter estimation — Gaussian case

We now have:

$$p(\vartheta|\mathbf{y}) \propto p(\mathbf{y}|\vartheta)p(\vartheta) = \frac{p(\mathbf{y}|\mathbf{x}, \vartheta)p(\mathbf{x}|\vartheta)}{p(\mathbf{x}|\mathbf{y}, \vartheta)} \cdot p(\vartheta) \quad \text{for any } \mathbf{x}.$$

Pick  $\mathbf{x} = \mathbf{E}(\mathbf{x}|\mathbf{y}, \vartheta) = \boldsymbol{\mu}_{\mathbf{x}|\mathbf{y}}$

$$\begin{aligned}p(\mathbf{y}|\vartheta) \propto \frac{|\mathbf{Q}|^{1/2} |\mathbf{Q}_\varepsilon|^{1/2}}{|\mathbf{Q}_{\mathbf{x}|\mathbf{y}}|^{1/2}} \exp\left(-\frac{1}{2} \left[ \boldsymbol{\mu}_{\mathbf{x}|\mathbf{y}}^\top \mathbf{Q} \boldsymbol{\mu}_{\mathbf{x}|\mathbf{y}} \right. \right. \\ \left. \left. + (\mathbf{y} - \mathbf{A}\boldsymbol{\mu}_{\mathbf{x}|\mathbf{y}})^\top \mathbf{Q}_\varepsilon (\mathbf{y} - \mathbf{A}\boldsymbol{\mu}_{\mathbf{x}|\mathbf{y}}) \right] \right).\end{aligned}$$

Estimate of  $\vartheta$

$$\hat{\vartheta} = \arg \max_{\vartheta} p(\vartheta|\mathbf{y}) = \arg \max_{\vartheta} p(\mathbf{y}|\vartheta)p(\vartheta)$$

## Non-Gaussian Observations

If the observations  $p(y_i|x_i, \vartheta)$  are **non-gaussian** the posterior

$$p(\mathbf{x}|\mathbf{y}, \vartheta)$$

**does not have a closed form** expression.

Note that the ordinary Bayesian solution of taking

$$p(\mathbf{x}|\mathbf{y}, \vartheta) \propto p(\mathbf{y}|\mathbf{x}, \vartheta)p(\mathbf{x}|\vartheta)$$

is **insufficient** since we need the **explicit dependence** on  $\vartheta$ !

We need an **approximation** of  $p(\mathbf{x}|\mathbf{y}, \vartheta)$ .

## Taylor Series

Given a univariate function  $f(x)$  the **Taylor series** around  $x^{(0)}$  is given by

$$f(x) \approx f(x^{(0)}) + f'(x^{(0)}) (x - x^{(0)}) + \frac{f''(x^{(0)})}{2!} (x - x^{(0)})^2 + \dots$$

## Taylor Series — Multivariate

Given a function  $f(\mathbf{x})$  where  $\mathbf{x}$  is a column vector, then

$$f(\mathbf{x}) \approx f(\mathbf{x}^{(0)}) + (\mathbf{x} - \mathbf{x}^{(0)})^\top \nabla f^{(0)} + \frac{1}{2!} (\mathbf{x} - \mathbf{x}^{(0)})^\top \mathbf{H}^{(0)} (\mathbf{x} - \mathbf{x}^{(0)}) + \dots$$

Here  $\nabla f^{(0)}$  is the **gradient** of  $f(\mathbf{x})$  evaluated at  $\mathbf{x}^{(0)}$  and  $\mathbf{H}^{(0)}$  is the **Hessian (second derivative)** matrix of  $f(\mathbf{x})$ .

## The classical Laplace approximation

To compute and approximate integrals of the form

$$\int \exp(ng(x)) dx$$

Pierre-Simon Laplace suggested a Taylor expansion of  $g(x)$  around it's maximum,  $x_0$ , (i.e.  $g'(x_0) = 0$  and  $g''(x_0) < 0$ )

$$g(x) \approx g(x_0) + \frac{1}{2} g''(x_0) (x - x_0)^2 + \dots$$

Giving

$$\begin{aligned} \int \exp(ng(x)) dx &\approx \exp(ng(x_0)) \int \exp\left(-\frac{n|g''(x_0)|}{2} (x - x_0)^2\right) dx \\ &= \exp(ng(x_0)) \sqrt{\frac{2\pi}{n|g''(x_0)|}} \end{aligned}$$

## Laplace approximation — Use in Statistics

Instead of approximating the integral  $\int \exp(ng(x)) dx$  we could use similar ideas to approximate densities of the form

$$p(\mathbf{x}) \propto \exp(g(\mathbf{x}))$$

Given the Taylor expansion of  $g(\mathbf{x})$

$$\begin{aligned} g(\mathbf{x}) &\approx g(x_0) + g'(x_0)(x - x_0) + \frac{1}{2} g''(x_0)(x - x_0)^2 \\ &= (g'(x_0) - x_0 g''(x_0))x + \frac{g''(x_0)}{2} x^2 + \text{const.} \end{aligned}$$

and assuming  $g''(x_0) < 0$  we have

$$X \overset{\text{approx.}}{\sim} N\left(\frac{g'(x_0) - x_0 g''(x_0)}{-g''(x_0)}, \frac{-1}{g''(x_0)}\right)$$



## Example — Gamma observations

Assuming a simple model with Gamma-distributed observations of one pixel

$$x \sim N(0, q^{-1}) \quad y|x \sim \Gamma\left(b, \frac{e^x}{b}\right)$$

with density functions

$$p(x) = \sqrt{\frac{q}{2\pi}} e^{-\frac{q}{2}x^2} \quad p(y|x) = \frac{y^{b-1}}{\Gamma(b)} \left(\frac{b}{e^x}\right)^b e^{-\frac{b \cdot y}{e^x}}$$

## Example — Posterior

For the posterior we now have

$$\begin{aligned} \log p(x|y) &= \log p(y|x) + \log p(x) - \log p(y) \\ &= -b \cdot x - b \cdot y \cdot e^{-x} - \frac{q \cdot x^2}{2} + \text{const.} \end{aligned}$$

with derivatives

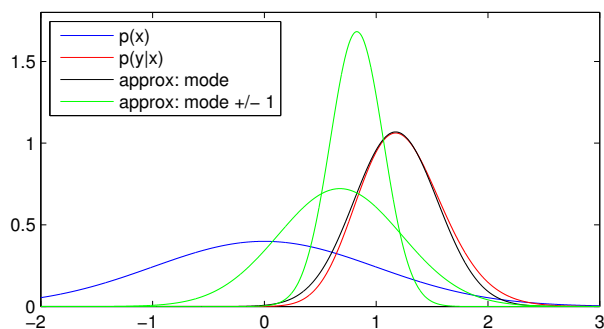
$$\begin{aligned} \frac{\partial \log p(x|y)}{\partial x} &= -b + b y e^{-x} - q x \\ \frac{\partial^2 \log p(x|y)}{\partial x^2} &= -b y e^{-x} - q \end{aligned}$$

giving

$$x|y \stackrel{\text{approx.}}{\sim} N\left(\frac{b y e^{-x_0} (1 + x_0) - b}{q + b y e^{-x_0}}, \frac{1}{q + b y e^{-x_0}}\right)$$

## Example — Effect of $x^{(0)}$

Posterior given  $q = 1$ ,  $b = 5$  and  $y = 4$  (or  $\log y = 1.39$ ).



Comparing the effect of different choices for  $x^{(0)}$ ; the mode of  $\log p(y|x)$  provides the best approximation.

## Parameter estimation using the Laplace approximation

Maximum posterior estimates of the parameters are given by maximising

$$p(\vartheta|\mathbf{y}) \propto \frac{p(\mathbf{y}|\mathbf{x}, \vartheta) p(\mathbf{x}|\vartheta)}{p(\mathbf{x}|\mathbf{y}, \vartheta)} \cdot p(\vartheta) \quad \text{for any } \mathbf{x}.$$

Replacing  $p(\mathbf{x}|\mathbf{y}, \vartheta)$  with the Laplace approximation  $p_G(\mathbf{x}|\mathbf{y}, \vartheta)$  we obtain the approximate posterior

$$\tilde{p}(\vartheta|\mathbf{y}) \propto \frac{p(\mathbf{y}|\mathbf{x}, \vartheta) p(\mathbf{x}|\vartheta)}{p_G(\mathbf{x}|\mathbf{y}, \vartheta)} \cdot p(\vartheta)$$

## Approximating the posterior

For (non-Gaussian) observations we have that

$$\begin{aligned} \log p(\mathbf{x}|\mathbf{y}, \vartheta) &= \log p(\mathbf{y}|\mathbf{x}, \vartheta) + \log p(\mathbf{x}|\vartheta) + \text{const.} \\ &= \sum_i \log p(y_i|x_i, \vartheta) + \log p(\mathbf{x}|\vartheta) + \text{const.} \\ &= \sum_i \log p(y_i|x_i, \vartheta) - \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} + \text{const.} \end{aligned}$$

Our plan is to construct a Gaussian approximation of the posterior  $p(\mathbf{x}|\mathbf{y}, \vartheta)$  by computation of a **second order Laplace (Taylor) approximation** of

$$f_i(x_i) = \log p(y_i|x_i, \vartheta) \quad \text{around } \mathbf{x}^{(0)}$$

## Laplace approximation

Given the log-posterior

$$\begin{aligned} \log p(\mathbf{x}|\mathbf{y}, \vartheta) &= \sum_i f_i(x_i) - \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} + \text{const.} \\ f_i(x_i) &= \log p(y_i|x_i, \vartheta) \end{aligned}$$

a Taylor expansion of  $f_i(x_i)$  around  $\mathbf{x}^{(0)}$  gives

$$f_i(x_i) \approx f_i(x_i^{(0)}) + f_i'(x_i^{(0)}) (x_i - x_i^{(0)}) + \frac{1}{2} f_i''(x_i^{(0)}) (x_i - x_i^{(0)})^2$$

Or in vector form

$$\sum_i f_i(x_i) \approx \sum_i f_i(x_i^{(0)}) + (\mathbf{x} - \mathbf{x}^{(0)})^\top \nabla f^{(0)} + \frac{1}{2} (\mathbf{x} - \mathbf{x}^{(0)})^\top \mathbf{H}_f^{(0)} (\mathbf{x} - \mathbf{x}^{(0)})$$

where  $\nabla f^{(0)} = [f_i'(x_i^{(0)})]_i$  and  $\mathbf{H}_f^{(0)} = \text{diag}(f_i''(x_i^{(0)}))$ .

## Laplace approximation

Collecting terms and accounting for the Gaussian-prior on  $\mathbf{x}$  we get

$$\begin{aligned}\log p(\mathbf{x}|\mathbf{y}, \vartheta) &= \sum_i f_i(x_i) - \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} + \text{const.} \\ &\approx \mathbf{x}^\top (\nabla f^{(0)} - \mathbf{H}_f^{(0)} \mathbf{x}^{(0)}) - \frac{1}{2} \mathbf{x}^\top (\mathbf{Q} - \mathbf{H}_f^{(0)}) \mathbf{x} + \text{const.}\end{aligned}$$

A Gaussian approximation  $p_G(\mathbf{x}|\mathbf{y}, \vartheta)$  is obtained with

$$\begin{aligned}\mathbf{E}_{\mathbf{x}^{(0)}}(\mathbf{x}|\mathbf{y}, \vartheta) &= \mathbf{Q}_{\mathbf{x}|\mathbf{y}}^{-1} (\nabla f^{(0)} - \mathbf{H}_f^{(0)} \mathbf{x}^{(0)}) \\ \mathbf{V}_{\mathbf{x}^{(0)}}(\mathbf{x}|\mathbf{y}, \vartheta) &= \mathbf{Q}_{\mathbf{x}|\mathbf{y}}^{-1} = (\mathbf{Q} - \mathbf{H}_f^{(0)})^{-1}\end{aligned}$$

Gaussian approximation — Suitable  $\mathbf{x}^{(0)}$ 

With the Gaussian approximation of the denominator we have

$$\tilde{p}(\vartheta|\mathbf{y}) \propto \frac{p(\mathbf{y}|\mathbf{x}, \vartheta) p(\mathbf{x}|\vartheta)}{p_G(\mathbf{x}|\mathbf{y}, \vartheta)} \cdot p(\vartheta) \quad \text{for any } \mathbf{x}.$$

We want a good approximation for the **most likely values** of the posterior distribution; expand around the **mode**

$$\begin{aligned}\mathbf{x}^{(0)} &= \arg \max_{\mathbf{x}} p(\mathbf{x}|\mathbf{y}, \vartheta) \\ &= \arg \max_{\mathbf{x}} \underbrace{(\log p(\mathbf{y}|\mathbf{x}, \vartheta) + \log p(\mathbf{x}|\vartheta))}_{f(\mathbf{x})} \\ &= \arg \max_{\mathbf{x}} \left( \sum_i f_i(x_i) - \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} \right).\end{aligned}$$

Gaussian approximation — Suitable  $\mathbf{x}^{(0)}$ 

Since the derivative of

$$\sum_i f_i(x_i) - \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x}$$

has to be zero at the mode,  $\mathbf{x}^{(0)}$ , we have

$$\nabla f^{(0)} - \mathbf{Q} \mathbf{x}^{(0)} = 0 \quad \iff \quad \nabla f^{(0)} = \mathbf{Q} \mathbf{x}^{(0)}$$

Inserting this into the conditional expectation for the Gaussian approximation we have

$$\begin{aligned}\mathbf{E}_{\mathbf{x}^{(0)}}(\mathbf{x}|\mathbf{y}, \vartheta) &= \mathbf{Q}_{\mathbf{x}|\mathbf{y}}^{-1} (\nabla f^{(0)} - \mathbf{H}_f^{(0)} \mathbf{x}^{(0)}) \\ &= \mathbf{Q}_{\mathbf{x}|\mathbf{y}}^{-1} (\mathbf{Q} \mathbf{x}^{(0)} - \mathbf{H}_f^{(0)} \mathbf{x}^{(0)}) \\ &= \mathbf{Q}_{\mathbf{x}|\mathbf{y}}^{-1} (\mathbf{Q} - \mathbf{H}_f^{(0)}) \mathbf{x}^{(0)} = \mathbf{x}^{(0)}\end{aligned}$$

## Parameter Estimation — Approximate posterior

To evaluate  $p(\vartheta|\mathbf{y})$  using the Taylor (Laplace) approximation do:

1. For a given  $\vartheta$  find the mode

$$\mathbf{x}^{(0)} = \arg \max_{\mathbf{x}} p(\mathbf{x}|\mathbf{y}, \vartheta)$$

2. Taylor expansion of  $f(\mathbf{x})$  around  $\mathbf{x}^{(0)}$
3. The approximation of  $p(\vartheta|\mathbf{y})$  is

$$\tilde{p}(\vartheta|\mathbf{y}) \propto \frac{p(\mathbf{y}|\mathbf{x}^{(0)}, \vartheta) p(\mathbf{x}^{(0)}|\vartheta)}{p_G(\mathbf{x}^{(0)}|\mathbf{y}, \vartheta)} \cdot p(\vartheta)$$

The (approximate) estimate of  $\vartheta$  is

$$\vartheta_{\text{MAP}} \approx \arg \max_{\vartheta} \tilde{p}(\vartheta|\mathbf{y})$$

## Parameter Estimation — Details

- Technically we have

$$\tilde{p}(\vartheta|\mathbf{y}) \propto \frac{p(\mathbf{y}|\mathbf{x}, \vartheta) p(\mathbf{x}|\vartheta)}{p_G(\mathbf{x}|\mathbf{y}, \vartheta)} \cdot p(\vartheta) \quad \text{for any } \mathbf{x}.$$

so the  $\mathbf{x}$  at which we evaluate and the  $\mathbf{x}^{(0)}$  for the Taylor expansion can be different.

- The approximation  $p_G(\mathbf{x}|\mathbf{y}, \vartheta)$  is best at  $\mathbf{x}^{(0)}$  making this a good choice to evaluate.
- Since the expectation and mode coincides the approximate posterior then simplifies to

$$\tilde{p}(\vartheta|\mathbf{y}) \propto \frac{p(\mathbf{y}|\mathbf{x}^{(0)}, \vartheta) p(\mathbf{x}^{(0)}|\vartheta)}{|\mathbf{Q}_{\mathbf{x}|\mathbf{y}}|^{1/2}} \cdot p(\vartheta)$$

## INLA (Rue et al., 2009)

The above estimate of  $\vartheta_{\text{MAP}}$  describes the first part of INLA (integrated nested Laplace approximations). Recall that we also want

- $p(\vartheta|\mathbf{y})$
- $p(\mathbf{x}|\mathbf{y})$

These can not be computed, however INLA provides a way of computing the marginals

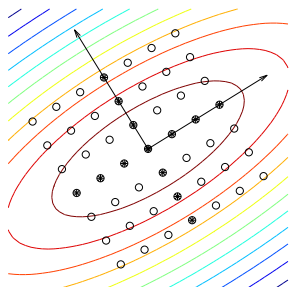
- $p(\vartheta_i|\mathbf{y})$
- $p(\mathbf{x}_i|\mathbf{y})$

which is often sufficient for inference.

## Posteriors for $\vartheta_j$

Given the mode  $\vartheta_{\text{MAP}}$  we compute the posterior  $p(\vartheta_j|\mathbf{y})$  by

1. Compute the Hessian of  $\tilde{p}(\vartheta|\mathbf{y})$  at the mode.
2. Use the Hessian to construct an integration grid.
3. Do numerical integration.



## Posteriors for $x_j$

We can write the posteriors of  $p(x_j|\mathbf{y})$  as

$$p(x_j|\mathbf{y}) = \int p(x_j|\mathbf{y}, \vartheta)p(\vartheta|\mathbf{y}) d\vartheta \\ \approx \int p_G(x_j|\mathbf{y}, \vartheta)\tilde{p}(\vartheta|\mathbf{y}) d\vartheta \approx \sum_k p_G(x_j|\mathbf{y}, \vartheta_k)\tilde{p}(\vartheta_k|\mathbf{y})\Delta_k$$

$\Delta_k$  are integration weights due to the point configuration.

- ▶ The part  $\tilde{p}(\vartheta_k|\mathbf{y})$  can be computed by evaluating the approximate posterior for each configuration  $\vartheta_k$ .
- ▶ The Gaussian approximation  $p_G(x_j|\mathbf{y}, \vartheta)$  can be obtained from the approximation  $p_G(\mathbf{x}|\mathbf{y}, \vartheta)$  as

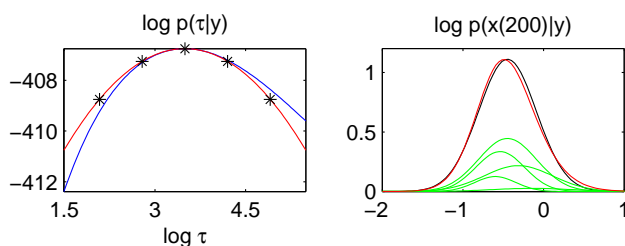
$$x_j|\mathbf{y}, \vartheta \stackrel{\text{approx.}}{\sim} \mathbf{N}(x_j^{(0)}, (Q_{x_j|\mathbf{y}}^{-1})_{jj})$$

- ▶ The multifrontal approach (Campbell and Davis, 1995) can be used to compute relevant elements of  $Q_{x_j|\mathbf{y}}^{-1}$ .

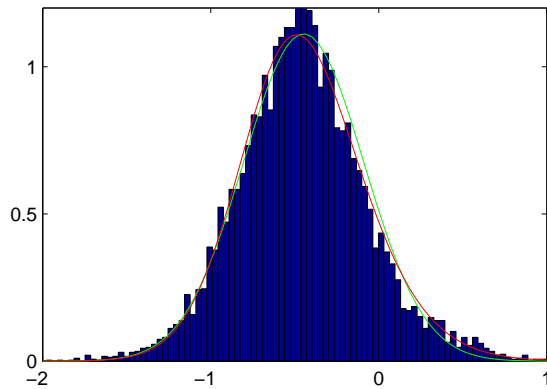
## Example: Tokyo rainfall

For the Tokyo rainfall the approximate posterior,  $\tilde{p}(\tau|\mathbf{y})$ , is evaluated and an integration grid using 5 points taken at

$$\tau_k = \tau_{\text{MAP}} \pm \frac{k}{\sqrt{-\tilde{p}''(\tau|\mathbf{y})}} \quad k = -2, -1, 0, 1, 2$$



## Tokyo rainfall — Comparing to MCMC



green: Approximation using one Gaussian.  
 red: Approximation using five Gaussians.

## INLA Details

For the integration two options could be used:

**Grid** More accurate, but time-consuming.

**Central Composite Design (CCD)** Faster than grid, often sufficient for computations of  $p(x_i|y)$  (less accurate for  $p(\vartheta_i|y)$ ).

When constructing the integration grid, the Hessian can be adjusted to account for different variances on each side of the mode.

Computations of  $p(x_i|y)$  can be improved by using a **skew-Normal** approximation instead of  $p_G(x_i|y, \vartheta)$ .

Implemented in the **R-INLA** package.

## INLA Assumptions

**Latent GMRF:**

Sparse of  $Q$  (and  $Q_{x|y}$ ) allowing for fast computations.

**Conditionally independent observations:**

Sparsity in  $Q$  translates into sparsity for  $Q_{x|y}$  and allows for simple Laplace approximation (no cross-derivatives).

**Few parameters ( $\leq 10$ ):**

Needed for the numerical integration.

## Using Laplace Approximation to improve MCMC

As an alternative to INLA the Laplace approximation can be used to improve the MCMC sampling and achieve blocking. I will briefly present a few ideas (and potential issues)

## Laplace Approximation — Only $\mathbf{x}$

A MCMC algorithm using proposals for  $\mathbf{x}$  based on the Laplace approximation is given by

1. Find the mode  $\hat{\mathbf{x}}(\vartheta)$ .
2. Draw a proposal,  $\mathbf{x}^*$  from  $N(\hat{\mathbf{x}}(\vartheta), \mathbf{Q}_{\mathbf{x}|\mathbf{y}}^{-1}(\vartheta))$
3. Accept the new sample with probability

$$\alpha(\mathbf{x}^*; \mathbf{x}^j) = \frac{p(\mathbf{x}^* | \tau, \mathbf{y})}{p(\mathbf{x}^j | \tau, \mathbf{y})} \cdot \frac{q(\mathbf{x}^j | \tau)}{q(\mathbf{x}^* | \tau)}$$

where  $\mathbf{x}^j$  is the current MCMC-sample and

$$q(\mathbf{x}^* | \tau) \propto \exp\left(-\frac{1}{2}(\mathbf{x}^* - \hat{\mathbf{x}}(\tau))^T \mathbf{Q}_{\mathbf{x}|\mathbf{y}}(\tau)(\mathbf{x}^* - \hat{\mathbf{x}}(\tau))\right)$$

4. Sample the parameters  $\vartheta$  using, e.g. random-walk (RW) or conjugate-priors.
5. Repeat.

## Laplace Approximation — Both $\mathbf{x}$ and $\vartheta$

To updated both  $\mathbf{x}$  and  $\vartheta$  jointly a suitable strategy is a RW proposal for  $\vartheta$  followed by sampling from the Laplace approximation obtained for the new  $\vartheta^*$ . This gives

1. Sample a new  $\vartheta^*$  using a RW
2. Sample a proposal for  $\mathbf{x}$  from the Laplace approximation given the new  $\vartheta^*$  (this includes computing the mode  $\hat{\mathbf{x}}(\vartheta^*)$ )

$$\mathbf{x}^* | \vartheta^* \sim N(\hat{\mathbf{x}}(\vartheta^*), \mathbf{Q}_{\mathbf{x}|\mathbf{y}}^{-1}(\vartheta^*))$$

The resulting proposal density can be decomposed as

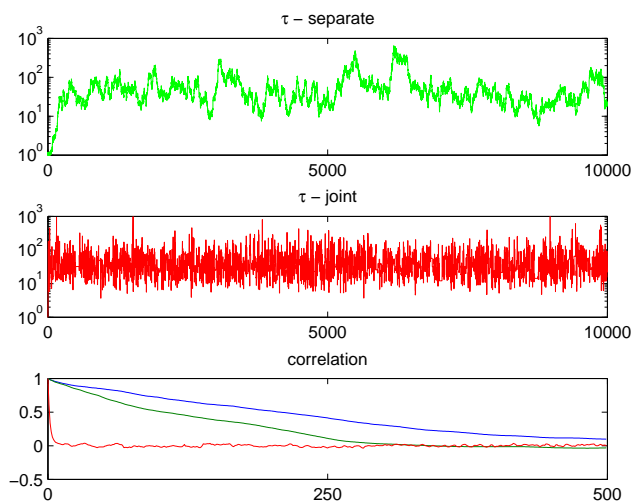
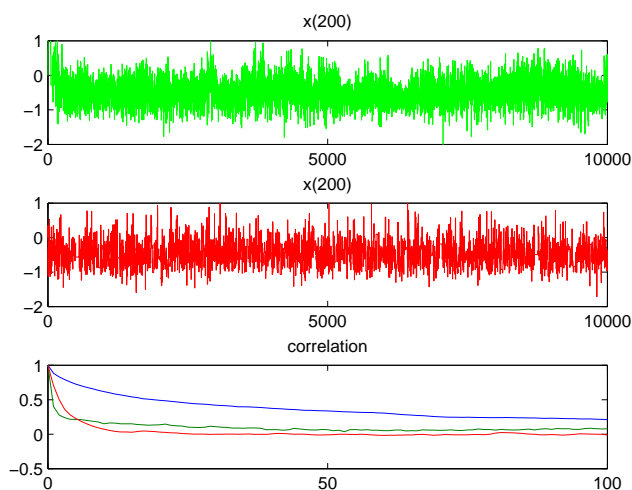
$$q(\mathbf{x}^*, \vartheta^* | \mathbf{x}^j, \vartheta^j) = q(\mathbf{x}^* | \vartheta^*) q(\vartheta^* | \vartheta^j)$$

3. Accept the new sample with probability

$$\alpha(\mathbf{x}^*, \vartheta^*; \mathbf{x}^j, \vartheta^j) = \frac{p(\mathbf{x}^*, \vartheta^* | \mathbf{y})}{p(\mathbf{x}^j, \vartheta^j | \mathbf{y})} \cdot \frac{q(\mathbf{x}^j | \vartheta^j) q(\vartheta^j | \vartheta^*)}{q(\mathbf{x}^* | \vartheta^*) q(\vartheta^* | \vartheta^j)}$$

4. Repeat.

A simple adaptive MCMC (Haario et al., 2001; Andrieu and Thoms, 2008) can be used to adjust the proposal variance for the RW.

$\tau$  — Laplace Approximation MCMC $x_{200}$  — Laplace Approximation MCMC

## Using Laplace Approximation to improve MCMC

Although popular (especially before INLA) the use of Laplace Approximations in the proposal for MCMC has several drawbacks:

- ▶ Potentially expensive mode-finding in each iteration.
- ▶ The acceptance rate decreases with increasing size of the latent field, for large applications  $\alpha(\mathbf{x}^*, \vartheta^*; \mathbf{x}^i, \vartheta^i) \ll 0.1$  are common.
- ▶ The joint proposal is dependent on good proposals for  $\vartheta^*$ .
- ▶ The separate proposals still have issues with poor mixing between  $\vartheta$  and  $\mathbf{x}$ .

Since we're computing derivatives anyway it's often better to use **Metropolis-Adjusted Langevin Algorithm (MALA)** (Girolami and Calderhead, 2011)



# Questions?

## Bibliography I

- Albert, J. H. and Chib, S. (1993), "Bayesian Analysis of Binary and Polychotomous Response Data," *J. Am. Stat. Assoc.*, 88, 669–679.
- Andrieu, C. and Thoms, J. (2008), "A tutorial on adaptive MCMC," *Stat. Comput.*, 18, 343–373.
- Campbell, Y. E. and Davis, T. A. (1995), "Computing the sparse inverse subset: an inverse multifrontal approach," Tech. rep., University of Florida, rEP-1995-183.
- Gilks, W. and Wild, P. (1992), "Adaptive Rejection Sampling for Gibbs Sampling," *J. R. Stat. Soc. C*, 41, 337–348.
- Girolami, M. and Calderhead, B. (2011), "Riemann manifold Langevin and Hamiltonian Monte Carlo methods," *J. R. Stat. Soc. B*, 73, 123–214.
- Haario, H., Saksman, E., and Tamminen, J. (2001), "An adaptive Metropolis algorithm," *Bernoulli*, 7, 223–224.
- Held, L. and Holmes, C. C. (2006), "Bayesian auxiliary variable models for binary and multinomial regression," *Bayesian Anal.*, 1, 145–168.

## Bibliography II

- Kitagawa, G. (1987), "Non-Gaussian State — Space Modeling of Nonstationary Time Series," *J. Am. Stat. Assoc.*, 82, 1032–1041.
- Knorr-Held, L. and Rue, H. (2002), "On Block Updating in Markov Random Field Models for Disease Mapping," *Scand. J. Stat.*, 29, 597–614.
- Rue, H., Martino, S., and Chopin, N. (2009), "Approximate Bayesian inference for hierarchical Gaussian Markov random field models," *J. R. Stat. Soc. B*, 71, 1–35.