

Latent Gaussian models: Approximate Bayesian inference (INLA)

Jo Eidsvik

January 30, 2018

Schedule

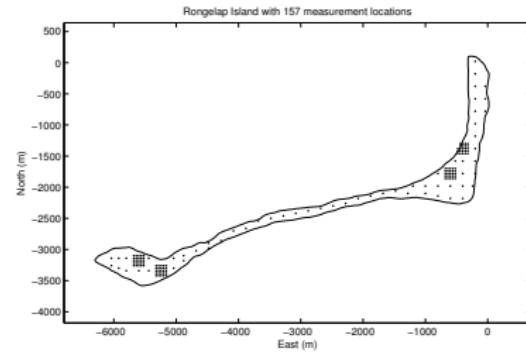
- ▶ 16 Jan: Gaussian processes (Jo Eidsvik)
- ▶ 23 Jan: Hands-on project on Gaussian processes (Team effort, work in groups)
- ▶ **30 Jan: Latent Gaussian models and INLA (Jo Eidsvik)**
- ▶ 6 Feb: Hands-on project on INLA (Team effort, work in groups)
- ▶ 12-13 Feb: Template model builder. (Guest lecturer Hans J Skaug)
- ▶ To be decided...

Plan for today

- ▶ Latent Gaussian models.
- ▶ Prior for parameters of Gaussian process. (We will be Bayesian today.)
- ▶ Laplace approximation and numerics for inference, INLA, (Rue et al., 2009)
- ▶ INLA shown for geostatistical applications. (Eidsvik et al., 2009)

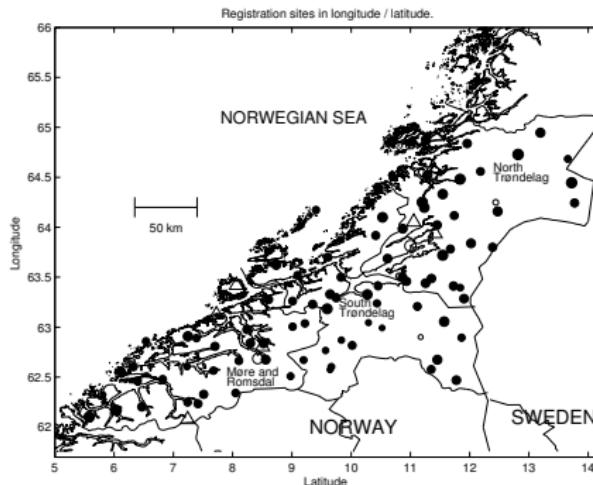
Examples of spatial latent Gaussian models

Radioactivity counts: Poisson



Example of spatial latent Gaussian models

Number of days with rain for $k = 92$ sites in September-October 2006.



Objective

Main goals:

- ▶ Fit model parameters of statistical covariance model.
- ▶ Predict latent intensity or risk at all spatial sites.

- ▶ Outlier detection.
- ▶ Spatial design.

Statistical model

Consider the following hierarchical model

1. Observed data $\mathbf{y} = (y_1, \dots, y_k)$ where

$$\pi(\mathbf{y} \mid \mathbf{x}, \boldsymbol{\eta}) = \pi(\mathbf{y} \mid \mathbf{x}) = \prod_{i=1}^k \pi(y_i \mid x_{s_i})$$

Often exponential family: Normal, Poisson, binomial, etc.

$$\log \pi(y_i | x_{s_i}) = \frac{y_i x_{s_i} - b(x_{s_i})}{a(\phi)} + c(\phi, y_i). \quad b(x) \text{ canonical link.}$$

2. Latent Gaussian field $\mathbf{x} = (x_1, \dots, x_k)$

$$\pi(\mathbf{x} \mid \boldsymbol{\eta}) = N[\mu \mathbf{1}_k, \boldsymbol{\Sigma}(\boldsymbol{\eta})]$$

3. Prior for hyperparameters $\pi(\boldsymbol{\eta})$

NOTE : Last point means we are Bayesian today!

Mixed models - Normal linear case

Common model

- ▶ $y_i = \mathbf{H}_i\boldsymbol{\beta} + x_i + \epsilon_i$
- ▶ y_i is observation. (Could be y_{ij} , individual or group i , replicate j .)
- ▶ $\boldsymbol{\beta}$ fixed effect. Prior $\pi(\boldsymbol{\beta}) \sim N(\boldsymbol{\mu}_{\boldsymbol{\beta}}, \boldsymbol{\Sigma}_{\boldsymbol{\beta}})$.
- ▶ x_i Gaussian random effect having a structured covariance model with parameter η . (Could be $\mathbf{U}_{ij}\mathbf{x}$ for group or individual i .)
- ▶ ϵ_i is random (unstructured) measurement noise.
 $\epsilon_i \sim N(0, \tau^2)$.

Mixed models - Inference

Can integrate out β .

$$\pi(\mathbf{x}|\boldsymbol{\eta}) = \int \pi(\mathbf{x}|\boldsymbol{\beta})\pi(\boldsymbol{\beta})d\boldsymbol{\beta} = N[\mathbf{H}\boldsymbol{\mu}_\beta, \mathbf{H}\boldsymbol{\Sigma}_\beta\mathbf{H}' + \boldsymbol{\Sigma}(\boldsymbol{\eta})]$$

Still a challenge to do inference on $\boldsymbol{\eta}$.

Mixed models - Inference

Common situation that has been hard to infer effectively:

- ▶ Frequentist, $\hat{\eta}$: Laplace approximations or estimating equations.
- ▶ Bayesian $\pi(\eta|\mathbf{y})$: Markov chain Monte Carlo.
- ▶ Inference not enough, wish to do model criticism, outlier detection, design, etc. Such goals require fast tools!

Mixed models - GLM

Likelihood is Poisson, binomial, or similar.

- ▶ Frequentist: Breslow and Clayton (1993).
- ▶ Bayes: Diggle, Tawn and Moyeed (1998), Christensen, Roberts and Sköld (2003), Diggle and Ribeiro (2007).

MCMC - Markov chain Monte Carlo

Around 2000 MCMC was very popular for inference and prediction in latent Gaussian models.

Today MCMC is still very popular, but not for latent Gaussian models.

These models are solved by Laplace approximations (Frequentist), or INLA (Bayes).

Typical MCMC algorithm

Initiate $\boldsymbol{\eta}^1, \mathbf{x}^1$.

Iterate for $i = 1, \dots, B$

- ▶ Propose $\boldsymbol{\eta}^* | \mathbf{x}, \mathbf{y}$.
- ▶ Accept (Set $\boldsymbol{\eta}^{i+1} = \boldsymbol{\eta}^*$) or reject (Set $\boldsymbol{\eta}^{i+1} = \boldsymbol{\eta}^i$).
- ▶ For all j . Propose $x_j^* | \mathbf{x}_{1:j-1}^{i+1}, \mathbf{x}_{j+1:k}^i, \boldsymbol{\eta}^{i+1}, \mathbf{y}$. Accept $(x_j^{i+1} = x_j^*)$ or reject $(x_j^{i+1} = x_j^i)$.

Converges to sampling from the joint distribution. All properties of distribution can be extracted from MCMC samples.

Mixing of Markov chain can be very slow. Blocking helps, but not always.

Gibbs sampler requires conjugate priors. Fast updates, but mixing not better.

Inference

Posterior

$$\pi(\mathbf{x}, \boldsymbol{\eta} \mid \mathbf{y}) \propto \pi(\boldsymbol{\eta}) \pi(\mathbf{x} \mid \boldsymbol{\eta}) \pi(\mathbf{y} \mid \mathbf{x})$$

In most cases the main tasks are:

- ▶ *PREDICTION*: Posterior marginals for x_j , $j = 1, \dots, n$

$$\pi(x_j \mid \mathbf{y})$$

- ▶ *PARAMETER ESTIMATION*: Posterior marginals for η_j

$$\pi(\eta_j \mid \mathbf{y})$$

Inference

Split the joint density

$$\pi(\mathbf{x}, \boldsymbol{\eta}, \mathbf{y}) = \pi(\boldsymbol{\eta})\pi(\mathbf{x}|\boldsymbol{\eta})\pi(\mathbf{y} | \mathbf{x}) = \pi(\mathbf{y})\pi(\boldsymbol{\eta} | \mathbf{y})\pi(\mathbf{x} | \boldsymbol{\eta}, \mathbf{y})$$

Clearly:

$$\pi(\boldsymbol{\eta} | \mathbf{y}) = \frac{\pi(\boldsymbol{\eta})\pi(\mathbf{x}|\boldsymbol{\eta})\pi(\mathbf{y} | \mathbf{x})}{\pi(\mathbf{y})\pi(\mathbf{x} | \boldsymbol{\eta}, \mathbf{y})} \propto \frac{\pi(\boldsymbol{\eta})\pi(\mathbf{x}|\boldsymbol{\eta})\pi(\mathbf{y} | \mathbf{x})}{\pi(\mathbf{x} | \boldsymbol{\eta}, \mathbf{y})}$$

Marginalization:

$$\pi(\mathbf{x}_j | \mathbf{y}) = \int_{\boldsymbol{\eta}} \pi(\boldsymbol{\eta} | \mathbf{y})\pi(\mathbf{x}_j | \boldsymbol{\eta}, \mathbf{y})d\boldsymbol{\eta}$$

Inference

Laplace approximation

$$\hat{\pi}(\boldsymbol{\eta} \mid \mathbf{y}) \propto \frac{\pi(\boldsymbol{\eta})\pi(\mathbf{x}|\boldsymbol{\eta})\pi(\mathbf{y} \mid \mathbf{x})}{\hat{\pi}(\mathbf{x} \mid \boldsymbol{\eta}, \mathbf{y})} \Big|_{\mathbf{x}=\hat{\mathbf{m}}(\boldsymbol{\eta}, \mathbf{y})}$$

Use a *Gaussian* approximation $\hat{\pi}(\mathbf{x} \mid \boldsymbol{\eta}, \mathbf{y})$.

$$\hat{\mathbf{m}} = \hat{\mathbf{m}}(\boldsymbol{\eta}, \mathbf{y}) = \operatorname{argmax}_{\mathbf{x}} [\pi(\mathbf{x}|\boldsymbol{\eta})\pi(\mathbf{y} \mid \mathbf{x})].$$

Approximate conjugacy

The Laplace approximation relies on approximate conjugacy. If the full conditional for \mathbf{x} is Gaussian, the formula is exact. When we insert a Gaussian approximation at the mode, the approximation depends on the non-Gaussian likelihood. This cannot be bimodal.

$$\hat{\pi}(\boldsymbol{\eta} \mid \mathbf{y}) \propto \frac{\pi(\boldsymbol{\eta})\pi(\mathbf{x} \mid \boldsymbol{\eta})\pi(\mathbf{y} \mid \mathbf{x})}{\hat{\pi}(\mathbf{x} \mid \boldsymbol{\eta}, \mathbf{y})} \Bigg|_{\mathbf{x}=\hat{\mathbf{m}}(\boldsymbol{\eta}, \mathbf{y})}$$

The error of the Laplace approximation (under weak regularity conditions) is *relative* and $\mathcal{O}(k^{-1})$ (Tierney and Kadane, 1986).

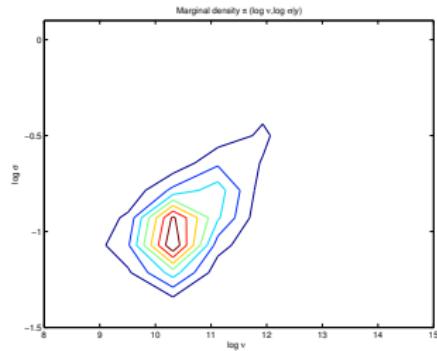
Gaussian approximation of full posterior

$$\pi(\mathbf{x} \mid \boldsymbol{\eta}, \mathbf{y}) \propto \exp \left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu} \mathbf{1}_n) \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu} \mathbf{1}_n) + \sum_{i=1}^k \log \pi(y_i | x_{s_i}) \right)$$

- ▶ $\log \pi(y_i | x_{s_i}) = \frac{y_i x_{s_i} - b(x_{s_i})}{a(\phi)} + c(\phi, y_i)$. $b(x)$ is canonical link.
- ▶ Expand GLM part $\log \pi(y_i | x_{s_i})$ to second order.
- ▶ Iterative solution to posterior mode $\hat{\mathbf{m}} = \hat{\mathbf{m}}(\boldsymbol{\eta}, \mathbf{y})$. ('Scoring').
- ▶ $\hat{\mathbf{m}} = \boldsymbol{\mu} \mathbf{1}_n - \boldsymbol{\Sigma} \mathbf{A}' [\mathbf{A} \boldsymbol{\Sigma} \mathbf{A}' + \mathbf{P}]^{-1} (\mathbf{z}(\mathbf{y}, \hat{\mathbf{m}}) - \boldsymbol{\mu} \mathbf{A} \mathbf{1}_n)$.
- ▶ Fit Gaussian approximation from Hessian at posterior mode:
 $\hat{\pi}(\mathbf{x} \mid \boldsymbol{\eta}, \mathbf{y}) = N(\hat{\mathbf{m}}, \hat{\mathbf{V}})$.
- ▶ $\mathbf{P} = \mathbf{P}(\hat{\mathbf{m}})$. Size $k \times k$ matrix inversion required.

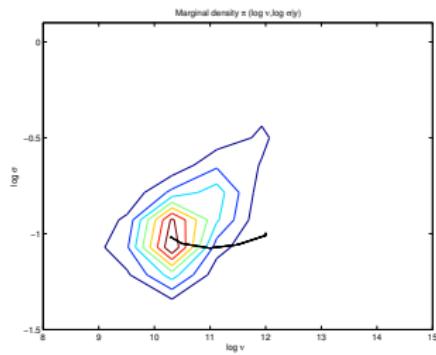
Practical implementation

Numerical approximation of $\hat{\pi}(\boldsymbol{\eta}|\mathbf{y})$



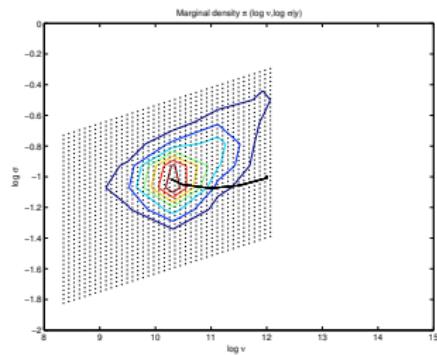
Practical implementation

Numerical approximation of $\hat{\pi}(\boldsymbol{\eta} | \mathbf{y})$, Step 1: Find mode



Each step requires $\mathbf{m}(\boldsymbol{\eta}, \mathbf{y})$, $\hat{\pi}(\mathbf{x} | \boldsymbol{\eta}, \mathbf{y})$ and Laplace.

Numerical approximation of $\hat{\pi}(\boldsymbol{\eta}|\mathbf{y})$, Step 2: Use Hessian at mode to set grid



Direct approximation of $\pi(x_j|\mathbf{y})$

Direct mixture approach for marginal prediction:

$$\hat{\pi}(x_j|\mathbf{y}) = \sum_I \hat{\pi}(x_j|\boldsymbol{\eta}_I, \mathbf{y}) \hat{\pi}(\boldsymbol{\eta}_I|\mathbf{y})$$

$$\hat{\pi}(x_j|\boldsymbol{\eta}_I, \mathbf{y}) = N(\hat{m}_j, \hat{V}_{j,j}).$$

$$\hat{m}_j = \hat{m}_j(\boldsymbol{\eta}_I, \mathbf{y}), \quad \hat{V}_{j,j} = \hat{V}_{j,j}(\boldsymbol{\eta}_I, \mathbf{y}).$$

Element j of posterior mode and j, j of full posterior covariance.

A frequentist solution would just plug in $\hat{\eta}$, the approximate MLE.

Nested approximation of $\pi(x_j|y)$

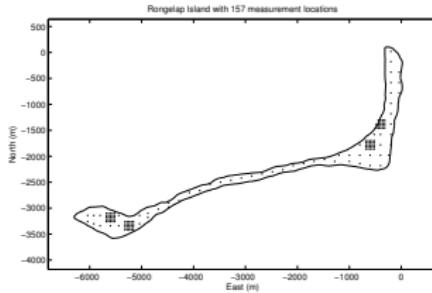
$$\pi(x_j|y, \theta) \propto \frac{\pi(y|x)\pi(x|\theta)}{\pi(x_{-j}|x_j, y, \theta)},$$

Using the Laplace approximation again, for fixed x_j .

$\hat{\pi}(x_{-j}|x_j, y, \theta)$ approximated by a Gaussian (for each x_j on a grid or design points).

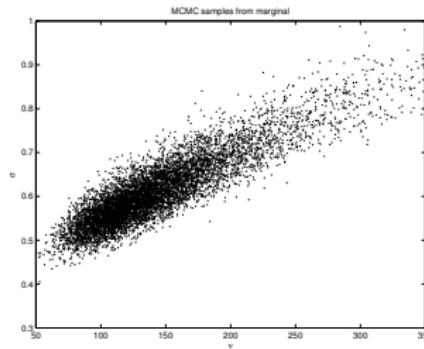
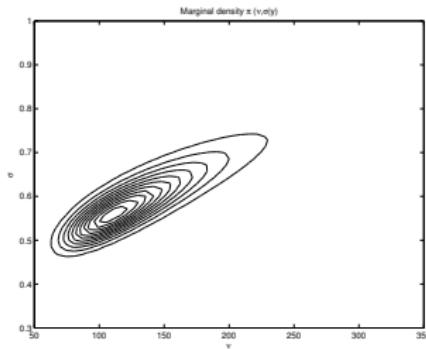
Example: Rongelap

- ▶ Radioactivity counts at 157 registration sites. Poisson counts.
- ▶ $\Sigma(\eta)$ defined from exponential covariance function.
 $\eta = (\nu, \sigma)$, range and standard deviation.



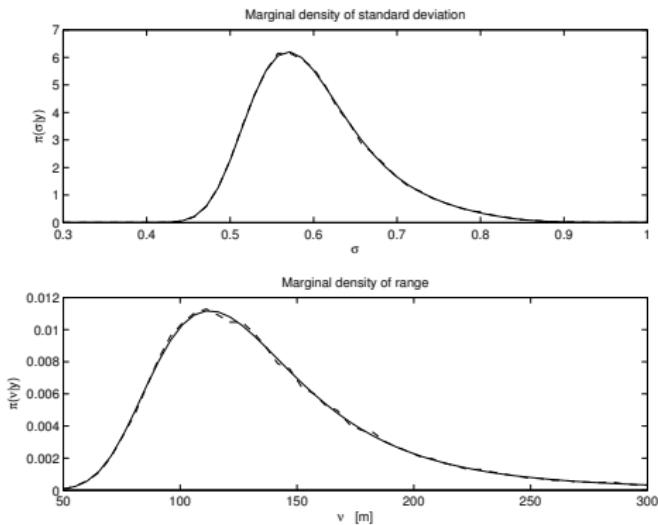
Marginals for $\hat{\pi}(\eta|y)$

Laplace approximation+numerics (left) and solutions with MCMC (right). Left) Seconds. Right) Minutes.

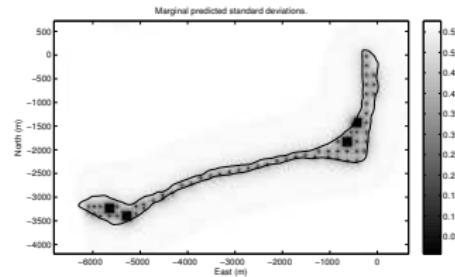
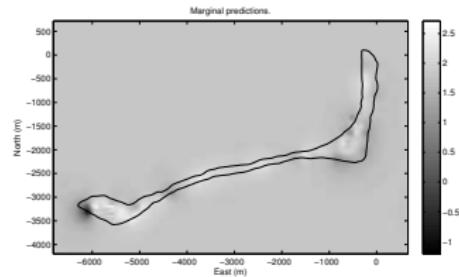


Marginals $\hat{\pi}(\eta|y)$

Laplace approximation (solid) and MCMC (dashed).

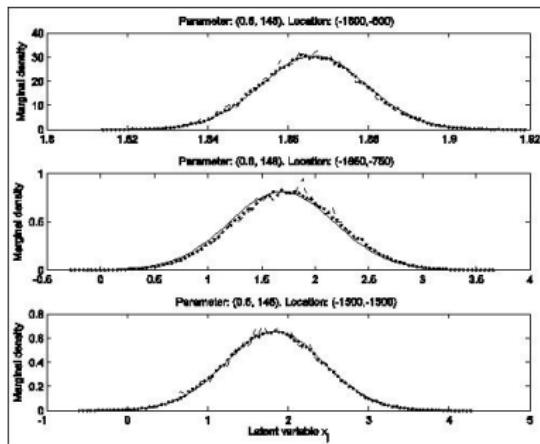
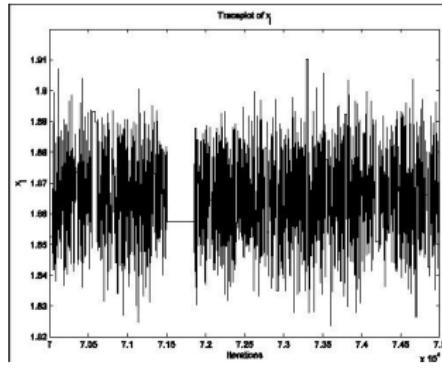


Prediction $\hat{E}(x_j|y)$ and $\hat{V}(x_j|y)$



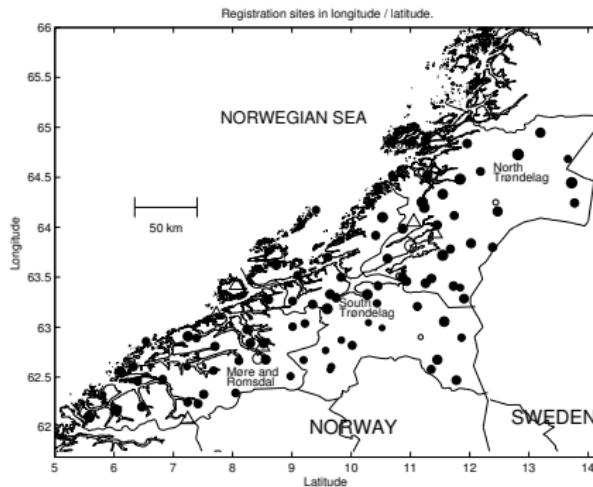
Marginals $\hat{\pi}(x_j|\eta, y)$

Conditional prediction at one spatial site MCMC (dashed),
Importance sampling (dotted) and direct Gaussian approximation
(solid).



Example: Precipitation in Middle Norway

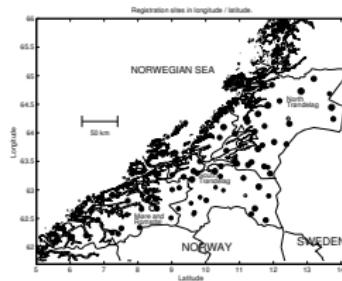
Number of days with rain for $k = 92$ sites in September-October 2006.



Example: Precipitation in Middle Norway

Binomial data $y_i = \text{Binomial}\left[\frac{e^{xs_i}}{1+e^{xs_i}}, 61\right]$.

Standard GLM gives no significance to East, North, Altitude.
Include only spatial trend.



- ▶ Outlier detection
- ▶ Spatial design

Outlier detection

Use crossvalidation $\pi(y_i | \mathbf{y}_{-i})$.

$$\hat{\pi}(y_i | \mathbf{y}_{-i}) = \int_{x_{s_i}} \sum_I \hat{\pi}(\boldsymbol{\eta}_I | \mathbf{y}_{-i}) \hat{\pi}(x_{s_i} | \boldsymbol{\eta}_I, \mathbf{y}_{-i}) \pi(y_i | x_{s_i}) dx_{s_i}$$

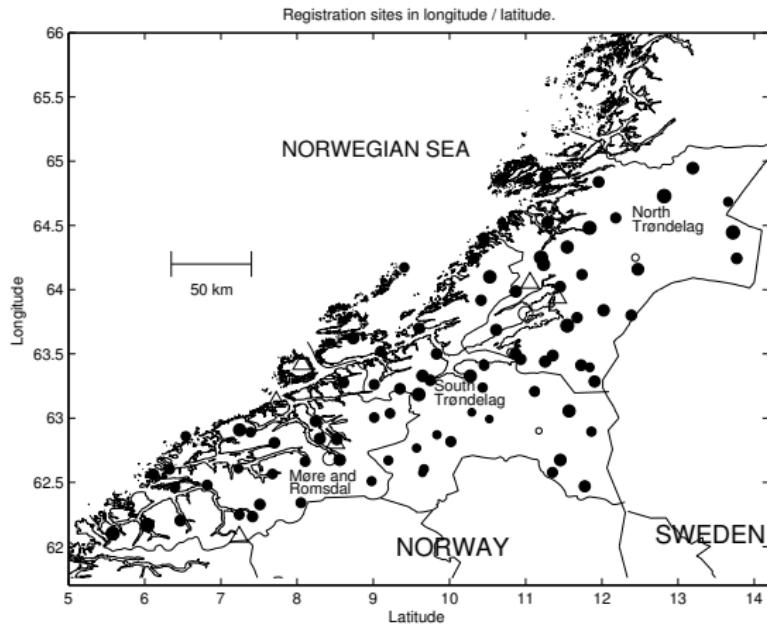
Inference separately for each y_i . I.e. k times. Approximate predictive percentiles

$$\sum_{y_i=0}^{y_{lower}} \hat{\pi}(y_i | \mathbf{y}_{-i}) = \alpha/2, \sum_{y_i=0}^{y_{upper}} \hat{\pi}(y_i | \mathbf{y}_{-i}) = 1 - \alpha/2.$$

Compare (y_{lower}, y_{upper}) with observed y_i .

Results : Outlier detection

Results $\alpha/2 = 0.01$: detect 4 outliers (open circles).



Spatial design

Prospective view: $\mathbf{y} \rightarrow (\mathbf{y}, \mathbf{y}_a)$.

\mathbf{y}_a extra data at 'new' spatial registration sites.

'Imagine' these observations - do not acquire them.

Design criterion is: Integrated prediction variance.

$$\hat{I} = \sum_{\mathbf{y}_a} \sum_j \hat{V}(x_j | \mathbf{y}, \mathbf{y}_a) \hat{\pi}(\mathbf{y}_a | \mathbf{y})$$

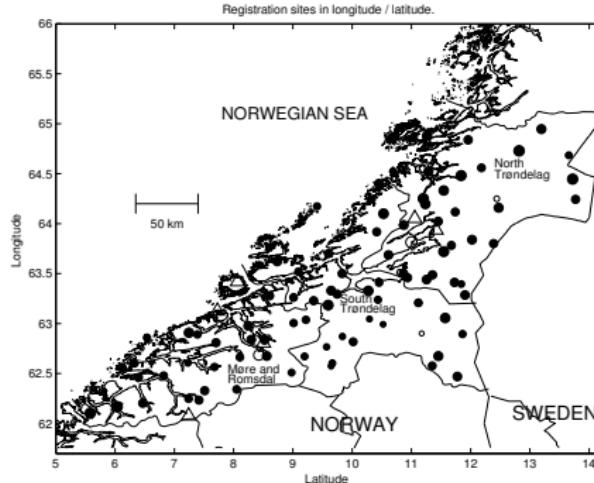
Results: Spatial design

Results of *three* design.

0: Existing design with 88 points (outliers excluded).

A: Currently installed stations, 88 plus 10 known sites (4 outlier sites and 6 sites out of service).

B: $88 + 10 = 2 \cdot 5$ new random sites around two existing sites (50km radius).



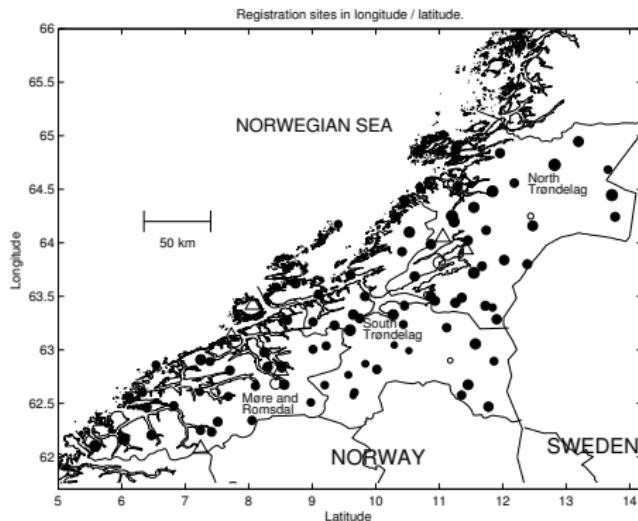
Results: Spatial design

Results of *three* designs.

0: Existing design: $\hat{I}_0 = 18.68$.

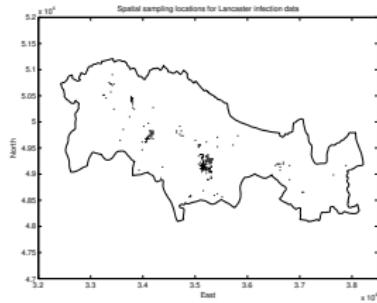
A: Currently installed stations: $\hat{I}_A = 17.94$.

B: Random around two existing sites: $\hat{I}_B = 17.85$



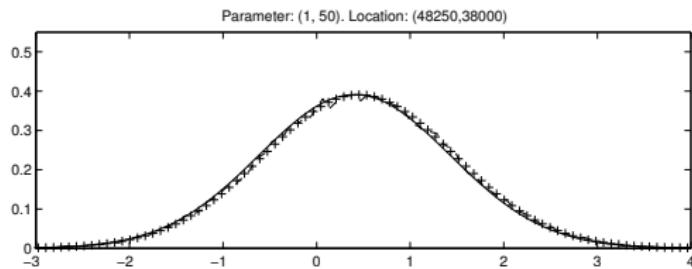
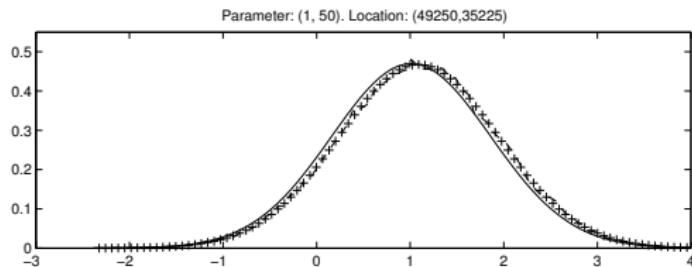
Example: Lancaster disease map

- ▶ Number of infections in different regions.
- ▶ Binomial data (with small counts).



Example: Lancaster disease map

LA, INLA and MCMC prediction at one site, for two parameter sets.



INLA contribution

- ▶ Mixed GLMs with latent Gaussian models cover wide range of applications
- ▶ The approximations work well for latent Gaussian models
- ▶ Generic routines. Software-friendly. Deterministic results (no Monte Carlo error)
- ▶ Enlarge scope of models

Conditions for INLA

- ▶ $\dim(\eta)$ is not too high
- ▶ No. of reg. sites $k < 10000$
- ▶ Marginals only. Bi-trivariate possible
- ▶ Likelihood must be well-behaved, not multimodal.

INLA vs MCMC

MCMC is very general. It explores all aspects of the joint posterior.
Approximate inference (INLA) is much faster. It is tailored to
special tasks, such as marginals.
Applicable to much more than spatial data.

INLA software

INLA software: <http://www.r-inla.org>

Easy to call:

```
inla(y ~ x + f(nu,model="iid"), family = c("poisson"), data =  
data, control.predictor=list(link=1))
```

Rue et al. (2009)

Routine runs on Gaussian Markov random fields.