



NTNU
Norwegian University of
Science and Technology

Introduction to INLA

Andrea Riebler <andrea.riebler@math.ntnu.no>

IBS Channel Network Conference 2015
Nijmegen, April 20th, 2015

2

Outline

Introduction

Bayesian hierarchical models

Latent Gaussian models

Deterministic inference

R-INLA

3

What?

The short answer:

INLA is a fast method to do Bayesian inference with latent Gaussian models and R-INLA is an R-package that implements this method with a flexible and simple interface.

A much longer answer:

Rue, Martino, and Chopin (2009) "Approximate Bayesian inference for latent Gaussian models by using integrated nested Laplace approximations." *Journal of the royal statistical society: Series B.* 319–392

4

Who?



There are
more:

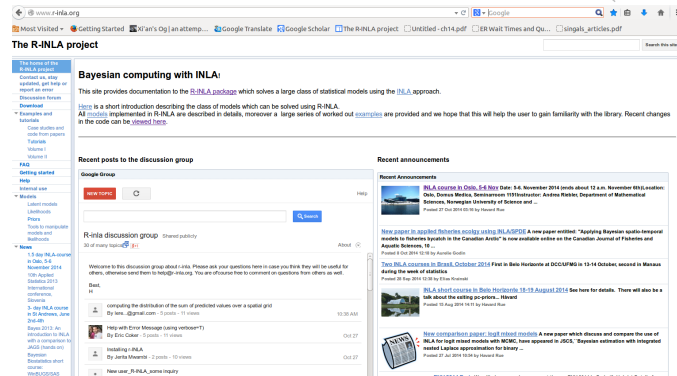


Håvard Rue, Finn Lindgren, Daniel Simpson, Andrea Riebler, Elias Teixeira Krainski, Geir-Arne Fuglstad, (Sara Martino, Thiago Guerrera Martins, Rupali Akerkar) and others (photo 2011)

5

Where?

The R-package R-INLA, some documentation, examples and help can be found at <http://www.r-inla.org>



A complete documentation about INLA is still in progress.

6

So... When can R-INLA be used?

- What type of problems can we solve?
- What type of models can we use?

7

The core

We have observed something.

We have questions.

We want answers.

8

How do we find answers?

We need to make choices:

Bayesian or frequentist?

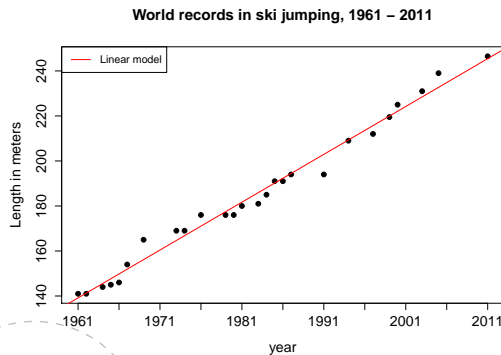
How do we model the data?

How do we compute the answer?

Example: Ski flying records

Assume a simple linear regression model with Gaussian observations $\mathbf{y} = (y_1, \dots, y_n)$, where

$$E(y_i) = \mu + \beta x_i, \quad \text{Var}(y_i) = \tau^{-1}, \quad 1, \dots, n$$



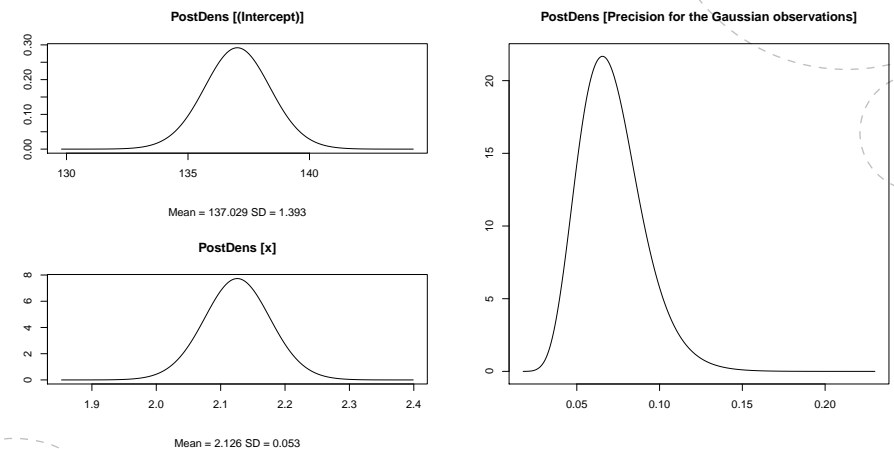
Estimates

$$\mu: 137.03 (1.42195),$$

$$\beta: 2.13 (0.05)$$

The Bayesian approach

Assign priors to the parameters α , β and τ and calculate posteriors:



Real-world datasets are usually much more complicated!

Using a Bayesian framework:

- Build (hierarchical) models to account for potentially complicated dependency structures in the data.
- Attribute uncertainty to model parameters and latent variables using priors.

Two main challenges:

- Select priors in a sensible way.
- Need computationally efficient methods to calculate posteriors.

Bayesian hierarchical models

INLA can be used with Bayesian hierarchical models where we model in different stages or levels:

- Stage 1:** What is the distribution of the responses?
- Stage 2:** What is the distribution of the underlying unobserved (latent) components?
- Stage 3:** What are our prior beliefs about the parameters controlling the components in the model?

Stage 1

How is our **data** (\mathbf{y}) generated from the underlying components (\mathbf{x}) and hyperparameters (θ) in the model:

- Gaussian response?
- Count data? (E.g. Poisson, negative binomial)
- Zero-inflation?
- Point pattern? (E.g. Log-Gaussian cox process)
- Binary data?

This information is placed into our *likelihood* $\pi(\mathbf{y}|\mathbf{x}, \theta)$

Stage 2

The underlying **unobserved components** \mathbf{x} are called **latent components** and can be:

- Covariates
- Unstructured random effects (individual effects, group effects)
- Structured random effects (AR(1), regional effects, continuously indexed spatial effects)

These are linked to the responses in the likelihood through linear predictors.

Stage 3

The likelihood and the latent model typically have hyperparameters that control their behavior. The **hyperparameters** θ can include:

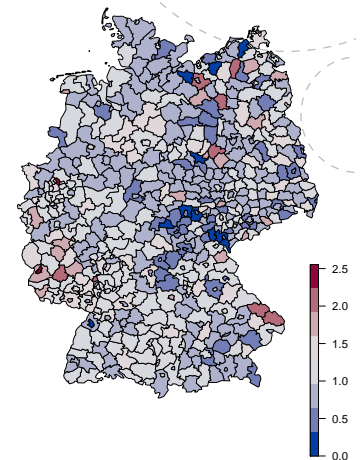
- Variance of observation noise
- Variance of unstructured effects
- Correlation of multivariate effects
- Autocorrelation parameter
- Probability of a zero (zero-inflated models)
- ...

Example: Disease mapping in Germany

We observed larynx cancer mortality counts for males in 544 district of Germany from 1986 to 1990 and want to make a model.

Information available:

- y_i : The count at location i .
- E_i : An offset; expected number of cases in district i .
- c_i : A covariate (level of smoking consumption) at location i
- s_i : spatial location i (here, district).



Stage 1: The data

First we decide on the likelihood for our data \mathbf{y}

- Our responses are counts
- We decide to model our responses as

$$y_i \mid \eta_i \sim \text{Poisson}(E_i \exp(\eta_i))$$

- η_i is a linear function of the latent components

Stage 2: The latent model

The latent field \mathbf{x} consists of two parts:

1. One fixed effect: the intercept μ
2.
 - The spatially structured effect f_s .
 - The unstructured effect \mathbf{u} which accounts for non-observed variability
 - The unknown effect $f(c_i)$ of the exposure covariate which assumes value c_i for district i .

These are combined for each location to give a linear predictor

$$\eta_i = \mu + f_s(s_i) + f(c_i) + u_i$$

The latent field is $\mathbf{x} = (\mu, \{f_s(\cdot)\}, \{f(\cdot)\}, u_1, u_2, \dots, u_n)$

Stage 3: Hyperparameters

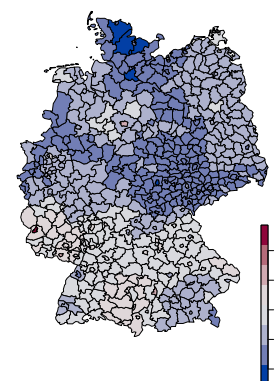
The structured and unstructured spatial effect as well as the smooth covariate effect will be each controlled by one parameter

- $\tau_c, \tau_f, \tau_\eta$: The precisions (inverse variances) of the covariate effect, spatial effect and unstructured effect, respectively.

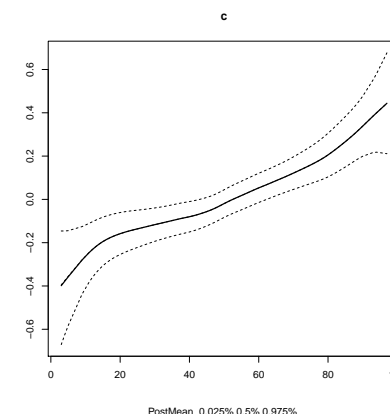
The hyperparameters are $\boldsymbol{\theta} = (\tau_c, \tau_f, \tau_\eta)$, and must be given a prior $\pi(\tau_c, \tau_f, \tau_\eta)$.

Quantities of interest

Structured spatial effect
 $\exp(f_s(s_i))$



Covariate effect $f(c_i)$



Latent Gaussian models

This example is just one example of a very useful class of models called **Latent Gaussian models**.

- The characteristic property is that the **latent part** of the hierarchical model is **Gaussian**, $\mathbf{x}|\boldsymbol{\theta} \sim N(\mathbf{0}, \mathbf{Q}^{-1})$
- The expected value is **0**
- The *precision* matrix (inverse covariance matrix) is **Q**

The general set-up

The set up contains GLMs, GLMMs, GAMs, GAMMs, and more. The mean of the observation i , μ_i , is connected to the linear predictor, η_i , through a link function g ,

$$\eta_i = g(\mu_i) = \mu + \mathbf{z}_i^\top \boldsymbol{\beta} + \sum_{\gamma} w_{\gamma,i} f_{\gamma}(\mathbf{c}_{\gamma,i}) + u_i, \quad i = 1, 2, \dots, n$$

where

μ : Intercept

$\boldsymbol{\beta}$: Fixed effects of covariates \mathbf{z}

$\{f_{\gamma}(\cdot)\}$: Non-linear/smooth effects of covariates \mathbf{c}

$\{w_{\gamma,i}\}$: Known weights defined for each observed data point

\mathbf{u} : Unstructured error terms

Loads of examples

- Generalized linear and additive (mixed) models
- Disease mapping
- Survival analysis
- Log-Gaussian Cox-processes
- Spatio and spatio-temporal models
- Stochastic volatility models
- Measurement error models
- And more!

Specification of the latent field

- Collect all parameters (random variables) in the linear predictor in a **latent field** $\mathbf{x} = \{\mu, \boldsymbol{\beta}, \{f_{\gamma}(\cdot)\}, \boldsymbol{\eta}\}$.
- A latent Gaussian model is obtained by assigning Gaussian priors to all elements of \mathbf{x} .
- Very flexible due to many different forms of the unknown functions $\{f_{\lambda}(\cdot)\}$:
- **Hyperparameters** account for variability and length/strength of dependence

Flexibility through f -functions

The functions $\{f_\gamma\}$ in the linear predictor make it possible to capture very different types of random effects in the same framework:

- $f(\text{time})$: For example, an AR(1) process, RW1 or RW2
- $f(\text{spatial location})$: For example, a Matérn field
- $f(\text{covariate})$: For example, a RW1 or RW2 on the covariate values
- $f(\text{time, spatial location})$ can be a spatio-temporal effect
- And much more

Additivity

- One of the most useful features of the framework is the additivity.
- Effects can easily be removed and added without difficulty.
- Each component might add a new latent part and might add new hyperparameters, but the modelling framework and computations stay the same.

Computations

So...

Now we have a modelling framework

But how do we get our answers?

What do we care about?

It depends on the problem!

- A single element of the latent field (e.g. the sign or quantiles of a fixed effect)
- A linear combination of elements from the latent field (the average over an area of a spatial effect, the difference of two effects)
- A single hyperparameter (the correlation)
- Predictions at unobserved locations

What do we need to compute?

Often we are interested in the posterior probability density of an element of the latent field

$$\pi(x_i | \mathbf{y})$$

or the posterior probability density of an element of the hyperparameters

$$\pi(\theta_j | \mathbf{y})$$

or some other statistics

$$\pi(f(\mathbf{x}, \boldsymbol{\theta}) | \mathbf{y})$$

But, as always in Bayesian statistics, we need to do high-dimensional integrals that cannot be computed analytically.

Traditional approach: MCMC*

Based on sampling. Construct Markov chains with the target posterior as stationary distribution.

- Extensively used within Bayesian inference since the 1980's.
- Flexible and general, sometimes the only thing we can do!
- A generic tool is available with JAGS/OpenBUGS.
- Tools for specific models are of course available, e.g. BayesX and stan.

* Markov chain Monte Carlo

Approximate inference

Bayesian inference can (almost) never be done exactly. Some form of approximation must always be done.

- MCMC “works” for everything, but it can be incredibly slow
- Is it possible to make a quicker, more specialized inference scheme which only needs to work for this limited class of models?

Recall: What is our model framework?

Latent Gaussian models

$$\mathbf{y} | \mathbf{x}, \boldsymbol{\theta} \sim \prod_i \pi(y_i | \eta_i, \boldsymbol{\theta})$$

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

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

Gaussian!

Not Gaussian

where the precision matrix $\mathbf{Q}(\boldsymbol{\theta})$ is sparse. Generally these “sparse” Gaussian distributions are called **Gaussian Markov random fields** (GMRFs).

The sparseness can be exploited for very quick computations for the Gaussian part of the model through numerical algorithms for sparse matrices.

The INLA idea

Use the posterior distribution

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

to approximate the posterior marginals

$$\pi(x_i \mid \mathbf{y}) \quad \text{and} \quad \pi(\theta_j \mid \mathbf{y})$$

directly.

Let us consider a toy example to illustrate the ideas.

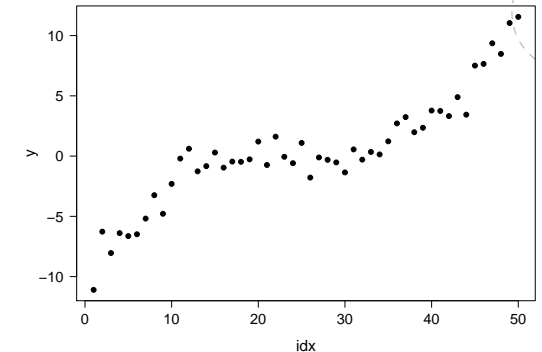
How does INLA work?

Observations

$$y_i = m(i) + \epsilon_i, \quad i = 1, \dots, n$$

Here, we assume that $m(i)$ is a smooth function wrt i and $\epsilon_i \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \tau_0)$ with *known* precision τ_0 .

```
1 n = 50
2 idx = 1:n
3 # generate something
  smooth representing m
4 fun = 100*((idx-n/2)/n)^3
5 # add some noise
6 y = fun + rnorm(n, mean
  =0, sd=1)
7 plot(idx, y)
```



Assumed hierarchical model

1. **Data:** Gaussian observations with known precision

$$y_i \mid x_i, \boldsymbol{\theta} \sim \mathcal{N}(x_i, \tau_0)$$

2. **Latent model:** A Gaussian model for the smooth function¹

$$\pi(\mathbf{x} \mid \boldsymbol{\theta}) \propto \theta^{(n-2)/2} \exp \left(-\frac{\theta}{2} \sum_{i=3}^n (x_i - 2x_{i-1} + x_{i-2})^2 \right)$$

3. **Hyperparameter:** The smoothing parameter θ which we assign a $\Gamma(a, b)$ prior

$$\pi(\theta) \propto \theta^{a-1} \exp(-b\theta), \quad \theta > 0$$

¹model="rw2"

Derivation of posterior marginals (I)

Since

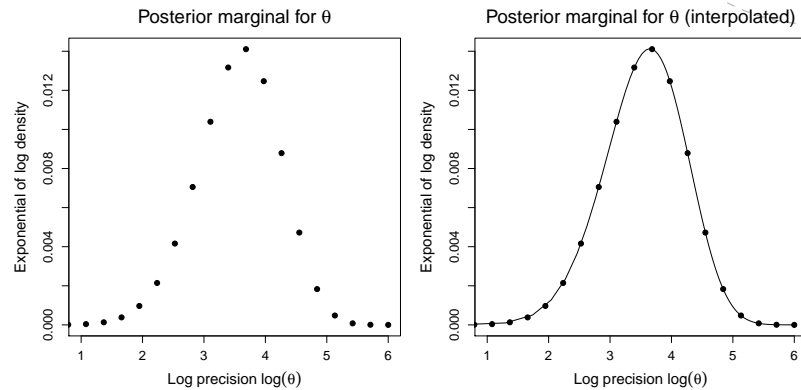
$$\mathbf{x}, \mathbf{y} \mid \boldsymbol{\theta} \sim \mathcal{N}(\cdot, \cdot)$$

(derived using $\pi(\mathbf{x}, \mathbf{y} \mid \boldsymbol{\theta}) \propto \pi(\mathbf{y} \mid \mathbf{x}, \boldsymbol{\theta}) \pi(\mathbf{x} \mid \boldsymbol{\theta})$), we can compute (numerically) all marginals, using that

$$\pi(\boldsymbol{\theta} \mid \mathbf{y}) \propto \frac{\overbrace{\pi(\mathbf{x}, \mathbf{y} \mid \boldsymbol{\theta})}^{\text{Gaussian}} \pi(\boldsymbol{\theta})}{\underbrace{\pi(\mathbf{x} \mid \mathbf{y}, \boldsymbol{\theta})}_{\text{Gaussian}}}$$

Posterior marginal for hyperparameter

Select a grid of point to represent the density $\theta | \mathbf{y}$. (Here, the points are chosen to be equi-distant).



Derivation of posterior marginals (II)

From

$$\mathbf{x} | \mathbf{y}, \theta \sim \mathcal{N}(\cdot, \cdot)$$

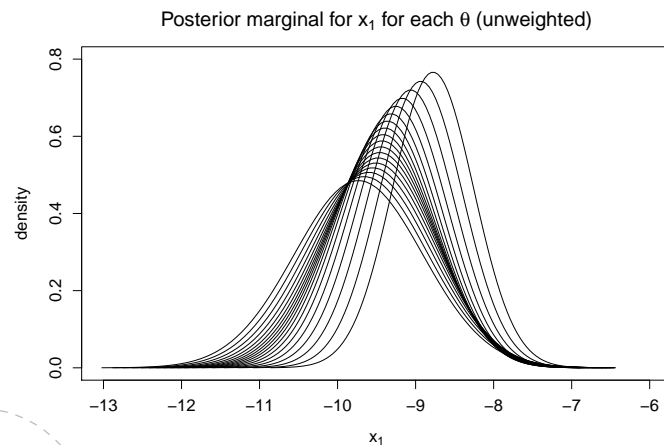
we can compute

$$\begin{aligned} \pi(x_i | \mathbf{y}) &= \int \underbrace{\pi(x_i | \theta, \mathbf{y})}_{\text{Gaussian}} \pi(\theta | \mathbf{y}) d\theta \\ &\approx \sum_k \pi(x_i | \theta_k, \mathbf{y}) \pi(\theta_k | \mathbf{y}) \Delta_k \end{aligned}$$

where $\theta_k, k = 1, \dots, K$, correspond to representative points of $\theta | \mathbf{y}$ and Δ_k are the corresponding weights.

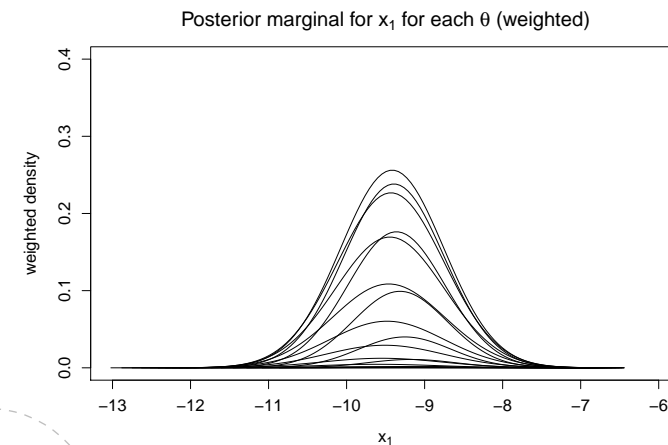
Posterior marginal for latent parameters

Compute the conditional marginal posterior for each x_i given θ_k . Here, shown for x_1 .



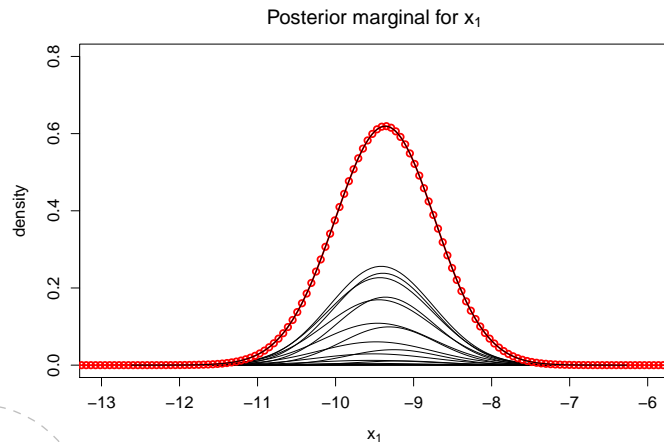
Posterior marginal for latent parameters

Weigh the resulting (conditional) marginal posterior by the density associated with each θ_k .



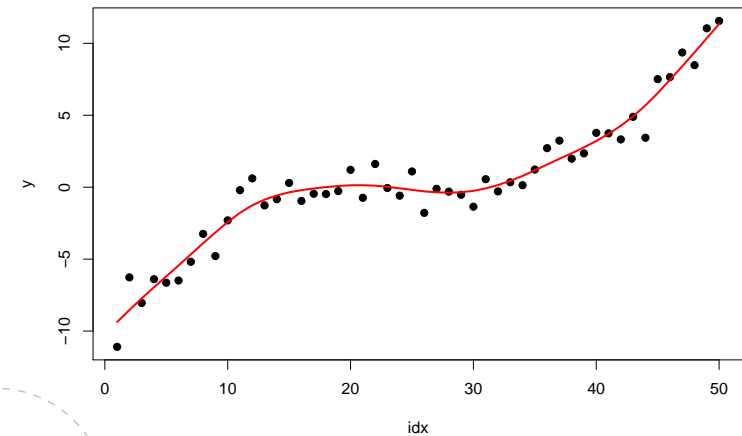
Posterior marginal for latent parameters

Numerically sum over all conditional densities to obtain the posterior marginal for each x_j .



Fitted spline

The posterior marginals are used to calculate summary statistics, like means, variances and credible intervals:



Extensions

This is the basic idea behind INLA.

However, we need to extend this basic idea so we can deal with

- More than one hyperparameter
- Non-Gaussian observations

The non-Gaussian part of the model

- In many cases $\pi(\mathbf{x} | \mathbf{y}, \boldsymbol{\theta})$ is very close to a Gaussian distribution, and can be replaced with a Laplace approximation
- This means that all the really hard, high-dimensional integrals with respect to the latent field are easy, and only the integrals with respect to the hyperparameters remain
- If the number of hyperparameters is low, these integrals can be done efficiently numerically

Limitations

- The dimension of the latent field \mathbf{x} can be large (10^2 – 10^6)
- But the dimension of the hyperparameters θ must be small (≤ 9)

In other words, each random effect can be big, but there cannot be too many random effects unless they share parameters.

How to use INLA?

INLA is implemented through the package `R-INLA` in the `R` software which

- is the most popular computing language in applied statistics
- is open source and *free*
- has a lot of packages that extend the functionality
- has a very user friendly `formula` interface