

# INLA - Introduction

Elias T. Krainski  
eliaskr@ufpr.br

62<sup>a</sup> RBras & 17<sup>o</sup> SEAGRO,  
Jul-2017, Lavras/MG

1 Tokyo example

2 On the Tokyo model

3 Bayesian inference

4 INLA overview

- 1 Tokyo example
- 2 On the Tokyo model
- 3 Bayesian inference
- 4 INLA overview

# A model for Tokyo data

Observation model

$$y_i \sim \text{Binomial}(n_i, p_i)$$

$$p_i = \frac{1}{1 + \exp(-x_i)}$$

the likelihood has no  $\theta$

$$\pi(\mathbf{y}|\mathbf{x}) = \prod_{i=1}^{366} \pi(y_i|x_i)$$

$$\pi(\mathbf{x}|\boldsymbol{\theta}) \propto \exp \left\{ -\frac{\theta}{2} \left[ (x_1 - x_{366})^2 + \sum_{i=2}^{366} (x_i - x_{i-1})^2 \right] \right\} \quad (1)$$

$$= \exp \left\{ -\frac{\theta}{2} \mathbf{x}^T \mathbf{R} \mathbf{x} \right\} \quad (2)$$

$$\pi(\mathbf{x}|\boldsymbol{\theta}) \propto \exp \left\{ -\frac{\theta}{2} \left[ (x_1 - x_{366})^2 + \sum_{i=2}^{366} (x_i - x_{i-1})^2 \right] \right\} \quad (1)$$

$$= \exp \left\{ -\frac{\theta}{2} \mathbf{x}^T \mathbf{R} \mathbf{x} \right\} \quad (2)$$

where  $\mathbf{R} =$

$$\begin{pmatrix} 2 & -1 & & & & & -1 \\ -1 & 2 & -1 & & & & \\ & -1 & 2 & -1 & & & \\ & & & \ddots & & & \\ & & & & -1 & 2 & -1 \\ & & & & & -1 & 2 & -1 \\ -1 & & & & & & -1 & 2 \end{pmatrix}$$

$$\mathbf{Q}(\boldsymbol{\theta}) = \boldsymbol{\theta} \mathbf{R}$$

$$\exp \left\{ -\frac{\theta}{2} \left[ (x_1 - x_{366})^2 + \sum_{i=2}^{366} (x_i - x_{i-1})^2 \right] \right\} \quad (3)$$

(4)

intrinsic/improper

$$\begin{array}{llll} x_i = 20, & x_{i-1} = 10 & \rightarrow & x_i - x_{i-1} = 10 \\ x_i = 10020, & x_{i-1} = 10010 & \rightarrow & x_i - x_{i-1} = 10 \end{array}$$

constraint or take the intercept out

- Tokyo example:  $Q(\theta) = \theta \mathbf{R}$ 
  - bigger  $\theta$  less variation of  $\mathbf{x}$ 
    - related to the variation of  $p_i$
- $\theta > 0$ : people usually use  $\theta \sim \text{Gamma}(a, b)$
- improper distribution:  $\theta$  values depends on  $\mathbf{R}$ 
  - hard to interpret  $\theta$  ( $a=?????$ ,  $b=?????$ )

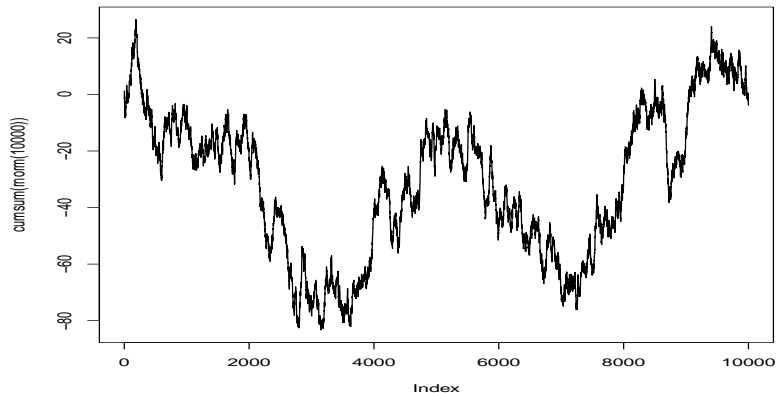


$\pi(\mathbf{x}|\theta = 1)$  and  $n$

The marginal variance and  $n$  relation

```
rw.var <- function(n, order) {  
  R <- as.matrix(INLA:::inla.rw(n, order=order))  
  mean(diag(INLA:::inla.ginv(R, rankdef=order)))  
}  
  
n <- c(10, 100, 366, 1000); names(n) <- n  
rbind(rw1=sapply(n, rw.var, order=1),  
      rw2=sapply(n, rw.var, order=2))  
  
##           10          100          366          1000  
## rw1 1.65    16.665    60.99954    166.6665  
## rw2 2.40 2381.190 116733.95702 2380955.1304
```

$\pi(\mathbf{x}|\theta = 1)$ : one realization



We need to control the marginal variance!

- 1 scale the model  $\rightarrow$  easy to interpret  $\theta$ 
  - Tutorial on `scale.option` at [www.r-inla.org/](http://www.r-inla.org/)

- ① scale the model  $\rightarrow$  easy to interpret  $\theta$ 
  - Tutorial on `scale.option` at [www.r-inla.org/](http://www.r-inla.org/)
- ② AND (new idea) Penalized complexity prior
  - P0: basic model:  $p_i = p_0$
  - P1: complex model:  $p_i$  varies
  - Kullback-Leibler divergence (KLD)
    - a distance from P1 model to P0,  $KLD(P1/P0) = 0$
  - allow variation on  $p_i$
  - AND supports the basic model
    - $\text{Gamma}(a, b)$  always overfits

- 1 Tokyo example
- 2 On the Tokyo model
- 3 Bayesian inference
- 4 INLA overview

# On our Bayesian hierarchical model

- Inference on (what we know about)  $\theta$  and  $\mathbf{x}$  given  $\mathbf{y}$ 
  - in maths:  $\pi(\mathbf{x}|\mathbf{y})$  and  $\pi(\theta|\mathbf{y})$
- considering  $\pi(\mathbf{y}|\mathbf{x}, \theta)$ ,  $\pi(\mathbf{x}|\theta)$  and  $\pi(\theta)$

# On our Bayesian hierarchical model

- Inference on (what we know about)  $\theta$  and  $\mathbf{x}$  given  $\mathbf{y}$ 
  - in maths:  $\pi(\mathbf{x}|\mathbf{y})$  and  $\pi(\theta|\mathbf{y})$
- considering  $\pi(\mathbf{y}|\mathbf{x}, \theta)$ ,  $\pi(\mathbf{x}|\theta)$  and  $\pi(\theta)$
- using the Bayes theorem,

$$\pi(\mathbf{x}|\mathbf{y}) = \int \pi(\mathbf{y}|\mathbf{x}, \theta)\pi(\mathbf{x}|\theta)\pi(\theta)d\theta$$

$$\pi(\theta|\mathbf{y}) = \int \pi(\mathbf{y}|\mathbf{x}, \theta)\pi(\mathbf{x}|\theta)\pi(\theta)d\mathbf{x}$$

# On our Bayesian hierarchical model

- Inference on (what we know about)  $\theta$  and  $\mathbf{x}$  given  $\mathbf{y}$ 
  - in maths:  $\pi(\mathbf{x}|\mathbf{y})$  and  $\pi(\theta|\mathbf{y})$
- considering  $\pi(\mathbf{y}|\mathbf{x}, \theta)$ ,  $\pi(\mathbf{x}|\theta)$  and  $\pi(\theta)$
- using the Bayes theorem,

$$\pi(\mathbf{x}|\mathbf{y}) = \int \pi(\mathbf{y}|\mathbf{x}, \theta)\pi(\mathbf{x}|\theta)\pi(\theta)d\theta$$

$$\pi(\theta|\mathbf{y}) = \int \pi(\mathbf{y}|\mathbf{x}, \theta)\pi(\mathbf{x}|\theta)\pi(\theta)d\mathbf{x}$$

- even more...
  - $\pi(\theta_j|\mathbf{y})$ ,  $j = 1, \dots, \dim(\theta)$
  - $\pi(x_i|\mathbf{y})$ ,  $i = 1, \dots, \dim(\mathbf{x})$



# The inference problem

- we have to compute

$$\pi(x_i|\mathbf{y}) \propto \int_{\mathbf{x}_{\{-i\}}} \int_{\boldsymbol{\theta}} \pi(y|\mathbf{x}, \boldsymbol{\theta}) \pi(\mathbf{x}|\boldsymbol{\theta}) \pi(\boldsymbol{\theta}) d\boldsymbol{\theta} d\mathbf{x}_{\{-i\}}$$

# The inference problem

- we have to compute

$$\pi(x_i|\mathbf{y}) \propto \int_{\mathbf{x}_{\{-i\}}} \int_{\boldsymbol{\theta}} \pi(y|\mathbf{x}, \boldsymbol{\theta}) \pi(\mathbf{x}|\boldsymbol{\theta}) \pi(\boldsymbol{\theta}) d\boldsymbol{\theta} d\mathbf{x}_{\{-i\}}$$

and

$$\pi(\theta_j|\mathbf{y}) \propto \int_{\mathbf{x}} \int_{\boldsymbol{\theta}_{\{-j\}}} \pi(y|\mathbf{x}, \boldsymbol{\theta}) \pi(\mathbf{x}|\boldsymbol{\theta}) \pi(\boldsymbol{\theta}) d\boldsymbol{\theta}_{\{-j\}} d\mathbf{x}$$

# The inference problem

- we have to compute

$$\pi(x_i|\mathbf{y}) \propto \int_{\mathbf{x}_{\{-i\}}} \int_{\boldsymbol{\theta}} \pi(y|\mathbf{x}, \boldsymbol{\theta}) \pi(\mathbf{x}|\boldsymbol{\theta}) \pi(\boldsymbol{\theta}) d\boldsymbol{\theta} d\mathbf{x}_{\{-i\}}$$

and

$$\pi(\theta_j|\mathbf{y}) \propto \int_{\mathbf{x}} \int_{\boldsymbol{\theta}_{\{-j\}}} \pi(y|\mathbf{x}, \boldsymbol{\theta}) \pi(\mathbf{x}|\boldsymbol{\theta}) \pi(\boldsymbol{\theta}) d\boldsymbol{\theta}_{\{-j\}} d\mathbf{x}$$

- remember
  - $\dim(\boldsymbol{\theta})$  is small
  - $\dim(\mathbf{x})$  is not small
  - we have to compute very high dimensional integrals

# The inference problem

- we have to compute

$$\pi(x_i|\mathbf{y}) \propto \int_{\mathbf{x}_{\{-i\}}} \int_{\boldsymbol{\theta}} \pi(y|\mathbf{x}, \boldsymbol{\theta}) \pi(\mathbf{x}|\boldsymbol{\theta}) \pi(\boldsymbol{\theta}) d\boldsymbol{\theta} d\mathbf{x}_{\{-i\}}$$

and

$$\pi(\theta_j|\mathbf{y}) \propto \int_{\mathbf{x}} \int_{\boldsymbol{\theta}_{\{-j\}}} \pi(y|\mathbf{x}, \boldsymbol{\theta}) \pi(\mathbf{x}|\boldsymbol{\theta}) \pi(\boldsymbol{\theta}) d\boldsymbol{\theta}_{\{-j\}} d\mathbf{x}$$

- remember
  - $\dim(\boldsymbol{\theta})$  is small
  - $\dim(\mathbf{x})$  is not small
  - we have to compute very high dimensional integrals
- typically they are not analytically tractable
  - $\rightarrow$  we have to approach

- single-site: compute (the expressions) for
  - $p(\theta_j | \boldsymbol{\theta}_{-j}, \mathbf{x}, \mathbf{y})$
  - $p(x_i | \mathbf{x}_{-i}, \boldsymbol{\theta}, \mathbf{y})$

- single-site: compute (the expressions) for
  - $p(\theta_j | \theta_{-j}, \mathbf{x}, \mathbf{y})$
  - $p(x_i | \mathbf{x}_{-i}, \theta, \mathbf{y})$
- draw samples from such conditionals
  - WinBUGS, OpenBUGS, JAGS, and others
- use these samples to summarize  $p(\mathbf{x})$  and  $p(\theta)$

- single-site: compute (the expressions) for
  - $p(\theta_j | \theta_{-j}, \mathbf{x}, \mathbf{y})$
  - $p(x_i | \mathbf{x}_{-i}, \theta, \mathbf{y})$
- draw samples from such conditionals
  - WinBUGS, OpenBUGS, JAGS, and others
- use these samples to summarize  $p(\mathbf{x})$  and  $p(\theta)$
- **warning**
  - sampling from  $x_i | \mathbf{x}_{-i}, \theta, \mathbf{y}$ 
    - slow convergence when strong dependence
    - **does not work for our example...**
  - better: draw joint sample from  $\mathbf{x} | \theta, \mathbf{y}$
  - best: use INLA

- 1 Tokyo example
- 2 On the Tokyo model
- 3 Bayesian inference
- 4 **INLA overview**



# What INLA does

- INLA does:
  - compute marginals of  $\pi(x_i|\mathbf{y})$  and  $\pi(\theta_j|\mathbf{y})$
- how?
  - approach  $\pi(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y})$  to approach  $\pi(\boldsymbol{\theta}|\mathbf{y})$
  - explore  $\pi(\boldsymbol{\theta}|\mathbf{y})$ 
    - approach  $\pi(\theta_j|\mathbf{y})$
  - approach  $\pi(x_i|\mathbf{x}_{-i})$

# Important ingredient

The GMRF-approximation

$$\pi(\mathbf{x} \mid \boldsymbol{\theta}, \mathbf{y}) \propto \exp \left( -\frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} + \sum_i \log \pi(y_i | x_i) \right)$$

## The GMRF-approximation

$$\begin{aligned}\pi(\mathbf{x} \mid \boldsymbol{\theta}, \mathbf{y}) &\propto \exp \left( -\frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} + \sum_i \log \pi(y_i | x_i) \right) \\ &\approx \exp \left( -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T (\mathbf{Q} + \text{diag}(\mathbf{c})) (\mathbf{x} - \boldsymbol{\mu}) \right) \\ &= \pi_G(\mathbf{x} | \boldsymbol{\theta}, \mathbf{y})\end{aligned}$$

$$c_i = -\frac{d^2 l_i}{dx_i^2} \text{ where } l_i = \log(\pi(y_i | x_i)), i = 1, \dots, \# \text{ data}$$

## The GMRF-approximation

$$\begin{aligned}\pi(\mathbf{x} \mid \boldsymbol{\theta}, \mathbf{y}) &\propto \exp \left( -\frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} + \sum_i \log \pi(y_i | x_i) \right) \\ &\approx \exp \left( -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T (\mathbf{Q} + \text{diag}(\mathbf{c})) (\mathbf{x} - \boldsymbol{\mu}) \right) \\ &= \pi_G(\mathbf{x} | \boldsymbol{\theta}, \mathbf{y})\end{aligned}$$

$c_i = -\frac{d^2 l_i}{dx_i^2}$  where  $l_i = \log(\pi(y_i | x_i))$ ,  $i = 1, \dots, \# \text{ data}$

- Markov and computational properties (on  $\mathbf{Q}$ ) are preserved

## The GMRF-approximation

$$\begin{aligned}\pi(\mathbf{x} \mid \boldsymbol{\theta}, \mathbf{y}) &\propto \exp \left( -\frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} + \sum_i \log \pi(y_i | x_i) \right) \\ &\approx \exp \left( -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T (\mathbf{Q} + \text{diag}(\mathbf{c})) (\mathbf{x} - \boldsymbol{\mu}) \right) \\ &= \pi_G(\mathbf{x} | \boldsymbol{\theta}, \mathbf{y})\end{aligned}$$

$c_i = -\frac{d^2 l_i}{dx_i^2}$  where  $l_i = \log(\pi(y_i | x_i))$ ,  $i = 1, \dots, \# \text{ data}$

- Markov and computational properties (on  $\mathbf{Q}$ ) are preserved
- $\tilde{\pi}(\mathbf{x} | \boldsymbol{\theta}, \mathbf{y})$  costs
  - temporal:  $O(n)$
  - spatial:  $O(n \log(n))$

If  $\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}$  is *Gaussian*, the “approximation” is exact.

- Considering

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

- Considering

$$\begin{aligned}\pi(\boldsymbol{\theta}|\mathbf{y}) &= \frac{\pi(\boldsymbol{\theta}, \mathbf{x}|\mathbf{y})}{\pi(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y})} \\ &\propto \frac{\pi(\boldsymbol{\theta})\pi(\mathbf{x}|\boldsymbol{\theta})\pi(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})}{\pi(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y})}\end{aligned}$$

- Considering

$$\begin{aligned}\pi(\boldsymbol{\theta}|\mathbf{y}) &= \frac{\pi(\boldsymbol{\theta}, \mathbf{x}|\mathbf{y})}{\pi(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y})} \\ &\propto \frac{\pi(\boldsymbol{\theta})\pi(\mathbf{x}|\boldsymbol{\theta})\pi(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})}{\pi(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y})}\end{aligned}$$

- Gaussian approximation to denominator

$$\pi(\boldsymbol{\theta}|\mathbf{y}) \approx \frac{\pi(\boldsymbol{\theta})\pi(\mathbf{x}|\boldsymbol{\theta})\pi(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})}{\pi_G(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y})} \Big|_{\mathbf{x}=\mathbf{x}^*(\boldsymbol{\theta})}$$



- Considering

$$\begin{aligned}\pi(\boldsymbol{\theta}|\mathbf{y}) &= \frac{\pi(\boldsymbol{\theta}, \mathbf{x}|\mathbf{y})}{\pi(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y})} \\ &\propto \frac{\pi(\boldsymbol{\theta})\pi(\mathbf{x}|\boldsymbol{\theta})\pi(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})}{\pi(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y})}\end{aligned}$$

- Gaussian approximation to denominator

$$\pi(\boldsymbol{\theta}|\mathbf{y}) \approx \frac{\pi(\boldsymbol{\theta})\pi(\mathbf{x}|\boldsymbol{\theta})\pi(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta})}{\pi_G(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y})} \Big|_{\mathbf{x}=\mathbf{x}^*(\boldsymbol{\theta})}$$

- mode of  $\tilde{\pi}(\boldsymbol{\theta}|\mathbf{y})$  (optimization)
  - explore  $\tilde{\pi}(\boldsymbol{\theta}|\mathbf{y})$ 
    - approach  $\pi(\theta_j|\mathbf{y})$  (numerical integration)

## Approaching $\pi(x_i | \mathbf{y}, \boldsymbol{\theta})$

- Problem
  - $\dim(\mathbf{x})=n$  is not small
  - $n$  marginals to compute
- Laplace approximation

$$\tilde{\pi}(x_i | \mathbf{y}, \boldsymbol{\theta}) \approx \frac{\pi(\mathbf{x}, \boldsymbol{\theta} | \mathbf{y})}{\tilde{\pi}_{GG}(\mathbf{x}_{-i} | x_i, \mathbf{y}, \boldsymbol{\theta})} \Big|_{\mathbf{x}_{-i} = \mathbf{x}_{-i}^*(x_i, \boldsymbol{\theta})}$$

## Approaching $\pi(x_i | \mathbf{y}, \boldsymbol{\theta})$

- Problem
  - $\dim(\mathbf{x})=n$  is not small
  - $n$  marginals to compute
- Laplace approximation

$$\tilde{\pi}(x_i | \mathbf{y}, \boldsymbol{\theta}) \approx \frac{\pi(\mathbf{x}, \boldsymbol{\theta} | \mathbf{y})}{\tilde{\pi}_{GG}(\mathbf{x}_{-i} | x_i, \mathbf{y}, \boldsymbol{\theta})} \Big|_{\mathbf{x}_{-i} = \mathbf{x}_{-i}^*(x_i, \boldsymbol{\theta})}$$

- simpler/cruider (fast) approximation (from  $\pi_G(\mathbf{x} | \mathbf{y}, \boldsymbol{\theta})$ )

$$\hat{\pi}(x_i | \mathbf{y}, \boldsymbol{\theta}) = N(x_i; \mu_i(\boldsymbol{\theta}), \sigma_i^2(\boldsymbol{\theta}))$$

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

- integrate  $\boldsymbol{\theta}$  out from  $\tilde{\pi}(\mathbf{x}_i | \mathbf{y}, \boldsymbol{\theta})$
- select values for  $\boldsymbol{\theta}$
- use weighted sum

$$\tilde{\pi}(\mathbf{x}_i | \mathbf{y}) \propto \sum_j \tilde{\pi}(\mathbf{x}_i|\mathbf{y}, \boldsymbol{\theta}_j) \times \tilde{\pi}(\boldsymbol{\theta}_j|\mathbf{y})$$

- 1 Expect  $\tilde{\pi}(\boldsymbol{\theta}|\mathbf{y})$  to be accurate, since
  - $\mathbf{x}|\boldsymbol{\theta}$  is a *priori* Gaussian
  - Likelihood models are 'well-behaved' so

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

is *almost* Gaussian.

- 2 There are no distributional assumptions on  $\boldsymbol{\theta}|\mathbf{y}$
- 3 Similar remarks are valid to

$$\tilde{\pi}(x_i \mid \boldsymbol{\theta}, \mathbf{y})$$

# How can we assess the error in the approximations?

**Tool 1:** Compare a sequence of improved approximations

- ① Gaussian approximation
- ② Simplified Laplace
- ③ Laplace

No big differences  $\rightarrow$  good approximation

# How can we assess the error in the approximations?

**Tool 2:** Estimate the “effective” number of parameters as defined in the Deviance Information Criteria:

$$p_D(\theta) = \overline{D}(\mathbf{x}; \theta) - D(\bar{\mathbf{x}}; \theta)$$

and compare this with the number of observations

Low ratio is good.

This criteria has theoretical justification.