

# Lecture 11: Conditional Dependency Graphs and More INLA

# Review: Bayesian Hierarchical models

Hierarchical models are an extremely useful tool in Bayesian model building.

Three parts:

- **Observation model  $\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}_1$** : Encodes information about observed data.
- **The latent model  $\mathbf{x}|\boldsymbol{\theta}_2$** : The unobserved process.
- **Hyperpriors for  $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \boldsymbol{\theta}_2)$** : Models for all of the parameters in the observation and latent processes.

Note: here we indicate the observed data by  $\mathbf{y}$  while  $\mathbf{x}$  and  $\boldsymbol{\theta}$  are parameters

# Bayesian Hierarchical models

## Unless otherwise specified or implied:

- Conditional independence is assumed
- Prior parameters,  $\theta$ , are independent except when conditioning of the responses,  $\mathbf{y}$

# Hierarchical Bayesian models - Tokyo rainfall example

Tokyo rainfall example from exercise 2

- $y_t$  number of times daily rainfall in Tokyo  $\geq 1$  mm,  $t = 1, \dots, 366$
- $\tau_t$  logit probability of exceeding 1 mm  $t = 1, \dots, 366$
- $n_t$  number of trials,  $t = 1, \dots, 366$
- $\pi(\tau_t) = \frac{1}{1 + \exp(-\tau_t)}$

# Hierarchical Bayesian models - Tokyo rainfall example

Tokyo rainfall example from exercise 2

- $y_t$  number of times daily rainfall in Tokyo  $\geq 1$  mm,  $t = 1, \dots, 366$
- $\tau_t$  logit probability of exceeding 1 mm  $t = 1, \dots, 366$
- $n_t$  number of trials,  $t = 1, \dots, 366$
- $\pi(\tau_t) = \frac{1}{1 + \exp(-\tau_t)}$

Model:

$$y_t \mid \tau_t \sim \text{Bin}(n_t, \pi(\tau_t))$$

# Hierarchical Bayesian models - Tokyo rainfall example

Tokyo rainfall example from exercise 2

- $y_t$  number of times daily rainfall in Tokyo  $\geq 1$  mm,  $t = 1, \dots, 366$
- $\tau_t$  logit probability of exceeding 1 mm  $t = 1, \dots, 366$
- $n_t$  number of trials,  $t = 1, \dots, 366$
- $\pi(\tau_t) = \frac{1}{1 + \exp(-\tau_t)}$

Model:

$$y_t \mid \tau_t \sim \text{Bin}(n_t, \pi(\tau_t))$$

Prior for  $\tau_t$ :

$$\tau_t = \tau_{t-1} + u_t, \quad u_t \stackrel{iid}{\sim} N(0, \sigma_u^2), \quad t = 2, \dots, 366.$$

# Hierarchical Bayesian models - Tokyo rainfall example

Tokyo rainfall example from exercise 2

- $y_t$  number of times daily rainfall in Tokyo  $\geq 1$  mm,  $t = 1, \dots, 366$
- $\tau_t$  logit probability of exceeding 1 mm  $t = 1, \dots, 366$
- $n_t$  number of trials,  $t = 1, \dots, 366$
- $\pi(\tau_t) = \frac{1}{1 + \exp(-\tau_t)}$

Model:

$$y_t \mid \tau_t \sim \text{Bin}(n_t, \pi(\tau_t))$$

Prior for  $\tau_t$ :

$$\tau_t = \tau_{t-1} + u_t, \quad u_t \stackrel{iid}{\sim} N(0, \sigma_u^2), \quad t = 2, \dots, 366.$$

Hyper-prior on  $\sigma_u^2$ :

$$\sigma_u^2 \sim \text{Inv-Gamma}(\alpha, \beta)$$

# Hierarchical Bayesian models - Tokyo rainfall example

Tokyo rainfall example from exercise 2

- $y_t$  number of times daily rainfall in Tokyo  $\geq 1$  mm,  $t = 1, \dots, 366$
- $\tau_t$  logit probability of exceeding 1 mm  $t = 1, \dots, 366$
- $n_t$  number of trials,  $t = 1, \dots, 366$
- $\pi(\tau_t) = \frac{1}{1 + \exp(-\tau_t)}$

Model:

$$y_t \mid \tau_t \sim \text{Bin}(n_t, \pi(\tau_t))$$

Prior for  $\tau_t$ :

$$\tau_t = \tau_{t-1} + u_t, \quad u_t \stackrel{iid}{\sim} N(0, \sigma_u^2), \quad t = 2, \dots, 366.$$

Hyper-prior on  $\sigma_u^2$ :

$$\sigma_u^2 \sim \text{Inv-Gamma}(\alpha, \beta)$$

Use conditional dependency graphs to visualize the conditional independence structure!



## Review: INLA

**What is it?** A numerical method to do fast approximate Bayesian inference

**Why?** We do not want to wait for the MCMC to converge.

**Where can it be applied?** The (wide) class of Latent Gaussian Models (a subclass of Bayesian hierarchical models)

**How does it work?** Uses GMRF and sparse matrix computations, Laplace approximation, numerical integration

**How do we use it** Already implemented in the R-INLA library

# Review: Ingredients of INLA

- Latent Gaussian Models
  - ▶ Class of models where INLA can be applied
- Gaussian Markov Random Fields
  - ▶ Sparse matrix computations
- Laplace Approximation
  - ▶ Method of approximating posterior

# Latent Gaussian Models: A Unified Framework

Observations:  $\mathbf{y}$

Latent field:  $\mathbf{x}$

Hyperparameters:  $\boldsymbol{\theta} = (\theta_1, \theta_2)$

# Latent Gaussian Models: A Unified Framework

Observations:  $\mathbf{y}$  Assumed **conditionally independent** given  $\mathbf{x}$  and  $\theta_1$

$$\mathbf{y}|\mathbf{x}, \theta_1 \sim \prod_i \pi(y_i|x_i, \theta).$$

Latent field:  $\mathbf{x}$

Hyperparameters:  $\theta = (\theta_1, \theta_2)$

# Latent Gaussian Models: A Unified Framework

Observations:  $\mathbf{y}$  Assumed **conditionally independent** given  $\mathbf{x}$  and  $\theta_1$

$$\mathbf{y}|\mathbf{x}, \theta_1 \sim \prod_i \pi(y_i|x_i, \theta).$$

Latent field:  $\mathbf{x}$  Assumed to be a **GMRF** with sparse precision matrix  $\mathbf{Q}(\theta_2)$

$$\mathbf{x}|\theta_1 \sim \mathcal{N}(0, \mathbf{Q}(\theta_2)^{-1})$$

The latent field  $\mathbf{x}$  can be large ( $10^1 - 10^6$ )

Hyperparameters:  $\theta = (\theta_1, \theta_2)$

# Latent Gaussian Models: A Unified Framework

Observations:  $\mathbf{y}$  Assumed **conditionally independent** given  $\mathbf{x}$  and  $\theta_1$

$$\mathbf{y}|\mathbf{x}, \theta_1 \sim \prod_i \pi(y_i|x_i, \theta).$$

Latent field:  $\mathbf{x}$  Assumed to be a **GMRF** with sparse precision matrix  $\mathbf{Q}(\theta_2)$

$$\mathbf{x}|\theta_1 \sim \mathcal{N}(0, \mathbf{Q}(\theta_2)^{-1})$$

The latent field  $\mathbf{x}$  can be large ( $10^1 - 10^6$ )

Hyperparameters:  $\theta = (\theta_1, \theta_2)$  Precision parameters of the Gaussian field and parameters of the likelihood

$$\theta \sim \pi(\theta)$$

The vector  $\theta$  is usually small (1-10)

## Latent Gaussian models

A very general way of specifying the problem is by modelling the mean for the  $i$ -th unit by means of an additive linear predictor, defined on a suitable scale (e.g. logistic for binomial data)

$$\eta_i = \alpha + \sum_{l=1}^L f_l(u_{li}) + \sum_{k=1}^K \beta_k z_{ki} + \epsilon_i$$

where

- $\alpha$  is the intercept
- $\boldsymbol{\beta} = (\beta_1, \dots, \beta_K)$  quantify the effect of  $\mathbf{x} = (x_1, \dots, x_K)$  on the response
- $\mathbf{f} = (f_1, \dots, f_L)$  is a set of functions defined in terms of some covariates  $\mathbf{z} = (z_1, \dots, z_K)$

And assume

$$\mathbf{x} = (\alpha, \boldsymbol{\beta}, \mathbf{f}) \sim \mathcal{N}(0, \mathbf{Q}(\theta)^{-1})$$

## Quantities of interest:

The posterior distribution is:

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

We want to approximate the posterior **marginals**

$$\pi(\theta_i|\mathbf{y}) = \int \pi(\theta|\mathbf{y})d\theta_{-i}$$

and

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

INLA strategy:

- approximate  $\pi(\theta|\mathbf{y})$  and  $\pi(x_i|\theta, \mathbf{y})$
- solve the integrals numerically



## Approximating $\pi(\boldsymbol{\theta}|\mathbf{y})$

- From  $\pi(\mathbf{x}, \boldsymbol{\theta}, \mathbf{y}) = \pi(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y}) \times \pi(\boldsymbol{\theta}|\mathbf{y}) \times \pi(\mathbf{y})$  it follows that

$$\pi(\boldsymbol{\theta}|\mathbf{y}) \propto \frac{\pi(\mathbf{x}, \boldsymbol{\theta}, \mathbf{y})}{\pi(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y})} \text{ for all } \mathbf{x}.$$

## Approximating $\pi(\boldsymbol{\theta}|\mathbf{y})$

- From  $\pi(\mathbf{x}, \boldsymbol{\theta}, \mathbf{y}) = \pi(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y}) \times \pi(\boldsymbol{\theta}|\mathbf{y}) \times \pi(\mathbf{y})$  it follows that

$$\pi(\boldsymbol{\theta}|\mathbf{y}) \propto \frac{\pi(\mathbf{x}, \boldsymbol{\theta}, \mathbf{y})}{\pi(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y})} \text{ for all } \mathbf{x}.$$

- INLA approximates  $\pi(\boldsymbol{\theta}|\mathbf{y})$  using

$$\tilde{\pi}(\boldsymbol{\theta}|\mathbf{y}) \propto \frac{\pi(\mathbf{x}, \boldsymbol{\theta}, \mathbf{y})}{\tilde{\pi}_G(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y})} \Big|_{\mathbf{x}=\mathbf{x}^*(\boldsymbol{\theta})}.$$

which is also known as **Laplace approximation**.

## Approximating $\pi(\boldsymbol{\theta}|\mathbf{y})$

- From  $\pi(\mathbf{x}, \boldsymbol{\theta}, \mathbf{y}) = \pi(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y}) \times \pi(\boldsymbol{\theta}|\mathbf{y}) \times \pi(\mathbf{y})$  it follows that

$$\pi(\boldsymbol{\theta}|\mathbf{y}) \propto \frac{\pi(\mathbf{x}, \boldsymbol{\theta}, \mathbf{y})}{\pi(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y})} \text{ for all } \mathbf{x}.$$

- INLA approximates  $\pi(\boldsymbol{\theta}|\mathbf{y})$  using

$$\tilde{\pi}(\boldsymbol{\theta}|\mathbf{y}) \propto \frac{\pi(\mathbf{x}, \boldsymbol{\theta}, \mathbf{y})}{\tilde{\pi}_G(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y})} \Big|_{\mathbf{x}=\mathbf{x}^*(\boldsymbol{\theta})}.$$

which is also known as **Laplace approximation**.

- Here  $\tilde{\pi}_G$  is the **Gaussian (GMRF) approximation** to  $\pi(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y})$  and  $\mathbf{x}^*(\boldsymbol{\theta})$  is the mode of  $\pi(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y})$ .

## The GMRF approximation

Let  $\mathbf{x}$  denote a GMRF with precision matrix  $\mathbf{Q}$  and mean  $\boldsymbol{\mu}$ .

Approximate

$$\pi(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y}) \propto \exp \left( -\frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} + \sum_{i=1}^n \log \pi(y_i|x_i) \right)$$

by using a second-order Taylor expansion of  $\log \pi(y_i|x_i)$  around  $\boldsymbol{\mu}_0$ , say.

Recall

$$f(x) \approx f(x_0) + f'(x_0)(x - x_0) + \frac{1}{2}f''(x_0)(x - x_0)^2 = a + bx - \frac{1}{2}cx^2$$

with  $b = f'(x_0) - f''(x_0)x_0$  and  $c = -f''(x_0)$ .

## The GMRF approximation (II)

Thus,

$$\begin{aligned}\tilde{\pi}(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y}) &\propto \exp\left(-\frac{1}{2}\mathbf{x}^\top \mathbf{Q}\mathbf{x} + \sum_{i=1}^n (a_i + b_i x_i - 0.5c_i x_i^2)\right) \\ &\propto \exp\left(-\frac{1}{2}\mathbf{x}^\top (\mathbf{Q} + \text{diag}(\mathbf{c}))\mathbf{x} + \mathbf{b}^\top \mathbf{x}\right)\end{aligned}$$

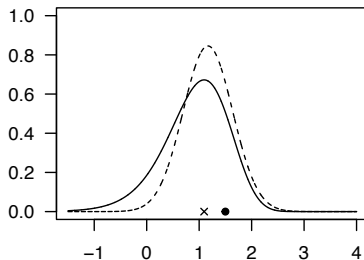
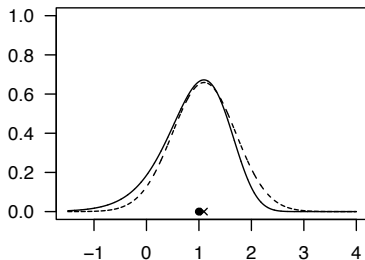
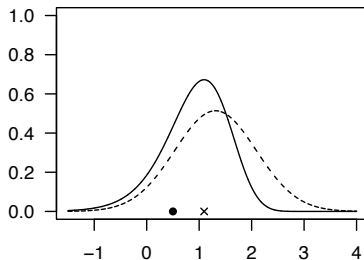
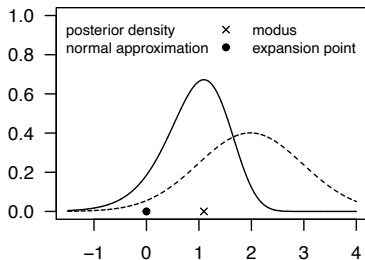
to get a Gaussian approximation with precision matrix  $\mathbf{Q} + \text{diag}(\mathbf{c})$  and mean given by the solution of  $(\mathbf{Q} + \text{diag}(\mathbf{c}))\boldsymbol{\mu} = \mathbf{b}$ . The canonical parameterization is

$$\mathcal{N}_C(\mathbf{b}, \mathbf{Q} + \text{diag}(\mathbf{c}))$$

which corresponds to

$$\mathcal{N}((\mathbf{Q} + \text{diag}(\mathbf{c}))^{-1}\mathbf{b}, (\mathbf{Q} + \text{diag}(\mathbf{c}))^{-1}).$$

# The GMRF approximation



## Exploring $\tilde{\pi}(\boldsymbol{\theta}|\mathbf{y})$

$\tilde{\pi}(\boldsymbol{\theta}|\mathbf{y})$  is “numerically explored” to find suitable support points  $\boldsymbol{\theta}_k$ .

**Main use:** Select good evaluation points  $\boldsymbol{\theta}_k$  for the numerical integration when approximating  $\tilde{\pi}(x_i|\mathbf{y})$

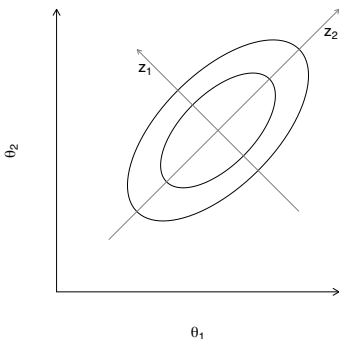
- Locate the mode

## Exploring $\tilde{\pi}(\boldsymbol{\theta}|\mathbf{y})$

$\tilde{\pi}(\boldsymbol{\theta}|\mathbf{y})$  is “numerically explored” to find suitable support points  $\boldsymbol{\theta}_k$ .

**Main use:** Select good evaluation points  $\boldsymbol{\theta}_k$  for the numerical integration when approximating  $\tilde{\pi}(x_i|\mathbf{y})$

- Locate the mode
- Compute the Hessian to construct principal components



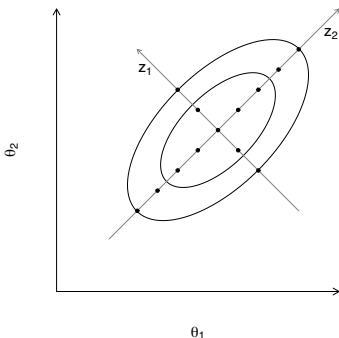


## Exploring $\tilde{\pi}(\boldsymbol{\theta}|\mathbf{y})$

$\tilde{\pi}(\boldsymbol{\theta}|\mathbf{y})$  is “numerically explored” to find suitable support points  $\boldsymbol{\theta}_k$ .

**Main use:** Select good evaluation points  $\boldsymbol{\theta}_k$  for the numerical integration when approximating  $\tilde{\pi}(x_i|\mathbf{y})$

- Locate the mode
- Compute the Hessian to construct principal components
- Grid-search to locate bulk of the probability mass

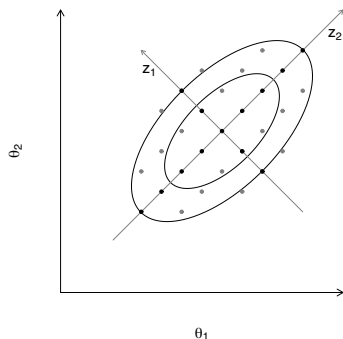


## Exploring $\tilde{\pi}(\boldsymbol{\theta}|\mathbf{y})$

$\tilde{\pi}(\boldsymbol{\theta}|\mathbf{y})$  is “numerically explored” to find suitable support points  $\boldsymbol{\theta}_k$ .

**Main use:** Select good evaluation points  $\boldsymbol{\theta}_k$  for the numerical integration when approximating  $\tilde{\pi}(x_i|\mathbf{y})$

- Locate the mode
- Compute the Hessian to construct principal components
- Grid-search to locate bulk of the probability mass



All points found have equal area weight  $\Delta_k$ .

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

For approximating the first component  $\pi(x_i|\boldsymbol{\theta}, \mathbf{y})$  we can use

- a **Gaussian approximation**, easily extractable from  $\tilde{\pi}_G(\mathbf{x}|\boldsymbol{\theta}, \mathbf{y})$ .

However, **errors in location and/or lack of skewness** possible.

## Approximating $\pi(x_i|\theta, y)$

For approximating the first component  $\pi(x_i|\theta, y)$  we can use

- a **Gaussian approximation**, easily extractable from  $\tilde{\pi}_G(\mathbf{x}|\theta, y)$ .  
However, **errors in location and/or lack of skewness** possible.
- a **Laplace approximation**

$$\tilde{\pi}_{\text{LA}}(x_i|\theta, y) \propto \frac{\pi(\mathbf{x}, \theta, y)}{\tilde{\pi}_{\text{GG}}(\mathbf{x}_{-i}|x_i, \theta, y)} \bigg|_{\mathbf{x}_{-i}=\mathbf{x}_{-i}^*(x_i, \theta)}.$$

The approximation is very accurate but very expensive.

## Approximating $\pi(x_i|\theta, y)$

For approximating the first component  $\pi(x_i|\theta, y)$  we can use

- a **Gaussian approximation**, easily extractable from  $\tilde{\pi}_G(\mathbf{x}|\theta, y)$ .  
However, **errors in location and/or lack of skewness** possible.
- a **Laplace approximation**

$$\tilde{\pi}_{LA}(x_i|\theta, y) \propto \frac{\pi(\mathbf{x}, \theta, y)}{\tilde{\pi}_{GG}(\mathbf{x}_{-i}|x_i, \theta, y)} \bigg|_{\mathbf{x}_{-i}=\mathbf{x}_{-i}^*(x_i, \theta)}.$$

The approximation is very accurate but very expensive.

- a **simplified Laplace approximation** based on fitting a skew-normal distribution to a series expansion of  $\tilde{\pi}_{LA}$ .

# INLA: Overview

**Step I** Approximate  $\pi(\boldsymbol{\theta}|\mathbf{y})$  using the Laplace approximation and select good evaluation points  $\boldsymbol{\theta}_k$ .

**Step II** For each  $\boldsymbol{\theta}_k$  and  $i$  approximate  $\pi(x_i|\boldsymbol{\theta}_k, \mathbf{y})$  using the Laplace or simplified Laplace approximation for selected values of  $x_i$ .

**Step III** For each  $i$ , sum out  $\boldsymbol{\theta}_k$

$$\tilde{\pi}(x_i|\mathbf{y}) = \sum_k \tilde{\pi}(x_i|\boldsymbol{\theta}_k, \mathbf{y}) \times \tilde{\pi}(\boldsymbol{\theta}_k|\mathbf{y}) \times \Delta_k.$$

Build a log spline corrected Gaussian to represent  $\tilde{\pi}(x_i|\mathbf{y})$ .

# INLA features

INLA fully incorporates posterior uncertainty with respect to hyperparameters  $\Rightarrow$  tool for full Bayesian inference

- Marginal posterior densities of all (hyper-)parameters
- Posterior mean, median, quantiles, std. deviation, etc.

The approach can be used for predictions, model assessment, . . .