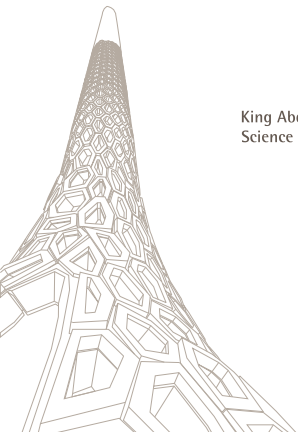


INLA 2.0 - into the future



King Abdullah University of
Science and Technology



جامعة الملك عبد الله
للعلوم والتقنية

May 2022



- 1 Introduction to INLA
- 2 Modern formulation
- 3 Linear predictor posterior inference
- 4 Examples
- 5 Discussion



Model definition

Suppose we have response data $\mathbf{y}_{n \times 1}$ with density function $\pi(y|\mathcal{X}, \boldsymbol{\theta})$ and link function $h(\cdot)$, that is linked to some covariates $\mathbf{Z} = \{\mathbf{X}, \mathbf{U}\}$ through linear predictors

$$\boldsymbol{\eta}_{n \times 1} = \beta_0 \mathbf{1} + \boldsymbol{\beta} \mathbf{X} + \sum_{k=1}^K f^k(\mathbf{u}_k)$$

The inferential aim is to estimate the latent field $\mathcal{X}_{m_* \times 1} = \{\beta_0, \boldsymbol{\beta}, \mathbf{f}\}$.
Define the augmented latent field

$$\mathcal{X}_{m \times 1} = \{\boldsymbol{\eta}, \beta_0, \boldsymbol{\beta}, \mathbf{f}\}.$$



Posterior approximations

$$\pi(\mathcal{X}, \boldsymbol{\theta}, \mathbf{y}) = \pi(\boldsymbol{\theta}) \pi(\mathcal{X} | \boldsymbol{\theta}) \prod_{i=1}^n \pi(y_i | \mathcal{X}_i, \boldsymbol{\theta})$$

$$\tilde{\pi}(\boldsymbol{\theta} | \mathbf{y}) \propto \left. \frac{\pi(\mathcal{X}, \boldsymbol{\theta}, \mathbf{y})}{\pi_G(\mathcal{X} | \boldsymbol{\theta}, \mathbf{y})} \right|_{\mathcal{X}=\boldsymbol{\mu}(\boldsymbol{\theta})}$$

$$\tilde{\pi}(\theta_j | \mathbf{y}) = \int \tilde{\pi}(\boldsymbol{\theta} | \mathbf{y}) d\boldsymbol{\theta}_{-j}$$

$$\tilde{\pi}(\mathcal{X}_j | \mathbf{y}) = \int \tilde{\pi}(\mathcal{X}_j | \boldsymbol{\theta}, \mathbf{y}) \tilde{\pi}(\boldsymbol{\theta} | \mathbf{y}) d\boldsymbol{\theta},$$

$\tilde{\pi}(\mathcal{X}_j | \boldsymbol{\theta}, \mathbf{y})$ depends on the approximation used, for Gaussian it is straightforward for the Laplace approximation we do another Gaussian approximation to $\tilde{\pi}(\boldsymbol{\mathcal{X}}_{-j} | \boldsymbol{\theta}, \mathbf{y})$.



Thoughts and ideas

- For large data, \mathcal{X} is large - can we remove η , but still produce cheap and accurate inference for η .
- For stability we can remove the "noisy" linear predictors (non-singularity)



Modern framework

The latent field is defined as

$$\boldsymbol{\mathcal{X}} = \{\beta_0, \boldsymbol{\beta}, \boldsymbol{f}\},$$

and the n linear predictors are defined as

$$\boldsymbol{\eta} = \boldsymbol{A}\boldsymbol{\mathcal{X}}, \tag{1}$$

with \boldsymbol{A} a sparse design matrix that links the linear predictors to the latent field.



Modern framework

From this formulation

$$\pi(\mathcal{X}, \boldsymbol{\theta} | \mathbf{y}) \propto \pi(\boldsymbol{\theta}) \pi(\mathcal{X} | \boldsymbol{\theta}) \prod_{i=1}^n \pi(y_i | (\mathbf{A}\mathcal{X})_i, \boldsymbol{\theta}).$$

$$\tilde{\pi}(\boldsymbol{\theta} | \mathbf{y}) = \frac{\pi(\mathcal{X}, \boldsymbol{\theta} | \mathbf{y})}{\pi_G(\mathcal{X} | \boldsymbol{\theta}, \mathbf{y})} \Big|_{\mathcal{X}=\boldsymbol{\mu}(\boldsymbol{\theta})}.$$



Modern framework

The Gaussian approximation $\pi_G(\mathcal{X}|\boldsymbol{\theta}, \mathbf{y})$ to $\pi(\mathcal{X}|\boldsymbol{\theta}, \mathbf{y})$ is calculated from a second order expansion of the likelihood around the mode of $\pi(\mathcal{X}|\boldsymbol{\theta}, \mathbf{y})$, $\boldsymbol{\mu}(\boldsymbol{\theta})$ as follows

$$\begin{aligned} \log(\pi(\mathcal{X}|\boldsymbol{\theta}, \mathbf{y})) &\propto -\frac{1}{2}\mathcal{X}^\top \mathbf{Q}(\boldsymbol{\theta})\mathcal{X} + \sum_{i=1}^n \left(b_i(\mathbf{A}\mathcal{X})_i - \frac{1}{2}c_i(\mathbf{A}\mathcal{X})_i^2 \right) \\ &= -\frac{1}{2}\mathcal{X}^\top \left(\mathbf{Q}(\boldsymbol{\theta}) + \mathbf{A}^\top \mathbf{D} \mathbf{A} \right) \mathcal{X} - \mathbf{b}^\top \mathbf{A} \mathcal{X} \end{aligned}$$

where \mathbf{b} is an n -dimensional vector with entries $\{b_i\}$ and \mathbf{D} is a diagonal matrix with n entries $\{c_i\}$. Note that both \mathbf{b} and \mathbf{D} depend on $\boldsymbol{\theta}$, so the Gaussian approximation is for a fixed $\boldsymbol{\theta}$.



Modern framework

The process is iterated to find \mathbf{b} and \mathbf{D} that gives the Gaussian approximation at the mode, $\mu(\theta)$, so that

$$\mathcal{X}|\theta, \mathbf{y} \sim N(\mu(\theta), \mathbf{Q}_{\mathcal{X}}^{-1}(\theta)) .$$

The graph of the Gaussian approximation consists of two components,

- ① \mathcal{G}_p : the graph obtained from the prior of the latent field through $\mathbf{Q}(\theta)$
- ② \mathcal{G}_d : the graph obtained from the data based on the non-zero entries of $\mathbf{A}^\top \mathbf{A}$



Modern framework

Next, the marginal conditional posteriors of the elements of $\boldsymbol{\mathcal{X}}$ is calculated from the joint Gaussian approximation as

$$\mathcal{X}_j | \boldsymbol{\theta}, \mathbf{y} \sim N \left((\boldsymbol{\mu}(\boldsymbol{\theta}))_j, (\mathbf{Q}_{\mathcal{X}}^{-1}(\boldsymbol{\theta}))_{jj} \right).$$

and the marginals

$$\tilde{\pi}(\mathcal{X}_j | \mathbf{y}) = \int \pi_G(\mathcal{X}_j | \boldsymbol{\theta}, \mathbf{y}) \tilde{\pi}(\boldsymbol{\theta} | \mathbf{y}) d\boldsymbol{\theta} \approx \sum_{k=1}^K \pi_G(\mathcal{X}_j | \boldsymbol{\theta}_k, \mathbf{y}) \tilde{\pi}(\boldsymbol{\theta}_k | \mathbf{y}) \delta_k.$$



Conditional posterior of η_i

In order to calculate $\tilde{\pi}(\eta_i|\mathbf{y})$, we first calculate $\tilde{\pi}(\eta_i|\boldsymbol{\theta}, \mathbf{y})$. We postulate a Gaussian density for $\eta_i|\boldsymbol{\theta}, \mathbf{y}$ such that $\tilde{\pi}(\eta_i|\boldsymbol{\theta}, \mathbf{y}) = \pi_G(\eta_i|\boldsymbol{\theta}, \mathbf{y})$, with mean

$$E(\boldsymbol{\eta}|\boldsymbol{\theta}, \mathbf{y}) = \mathbf{A}E(\boldsymbol{\mathcal{X}}|\boldsymbol{\theta}, \mathbf{y}) = \mathbf{A}\boldsymbol{\mu}(\boldsymbol{\theta})$$

and covariance matrix

$$\text{Cov}(\boldsymbol{\eta}|\boldsymbol{\theta}, \mathbf{y}) = \mathbf{A}\text{Cov}(\boldsymbol{\mathcal{X}}|\boldsymbol{\theta}, \mathbf{y})\mathbf{A}^\top,$$



Marginal posterior of η_j

Now let \mathbf{C} be a sparse selected inverse of $\mathbf{Q}_{\mathcal{X}} = \mathbf{Q} + \mathbf{A}^\top \mathbf{D} \mathbf{A}$ based on the graph $\mathcal{G}_{\mathcal{X}} = \{\mathcal{G}_p, \mathcal{G}_d\}$

$$\begin{aligned}\eta_j | \boldsymbol{\theta}, \mathbf{y} &\sim N(\mu_j(\boldsymbol{\theta}), \sigma_j^2(\boldsymbol{\theta})) \\ \mu_j(\boldsymbol{\theta}) &= (\mathbf{A}\boldsymbol{\mu}(\boldsymbol{\theta}))_j \\ \sigma_j^2(\boldsymbol{\theta}) &= \sum_{il} A_{ji} A_{jl} C_{il} \\ \tilde{\pi}(\eta_j | \mathbf{y}) &\approx \sum_{k=1}^K \pi_G(\eta_j | \boldsymbol{\theta}_k, \mathbf{y}) \tilde{\pi}(\boldsymbol{\theta}_k | \mathbf{y}) \delta_k,\end{aligned}$$

so that $E(\eta_j | \mathbf{y}) = \mu_j$ and $\text{Var}(\eta_j | \mathbf{y}) = \sigma_j^2$.



VB correction to latent field posteriors

VB mean correction paper: <https://arxiv.org/abs/2111.12945>

Let

$$\mu^*(\theta) = \mu(\theta) + \mathbf{M}\lambda,$$

with

$$\arg_{\lambda} \min \left(E_{\mathcal{X}|\mathbf{y}, \theta \sim N(\mu(\theta) + \mathbf{M}\lambda, \mathbf{Q}_{\mathcal{X}}^{-1}(\theta))} [-\log l(\mathcal{X}|\mathbf{y})] \right. \\ \left. + \frac{1}{2} (\mu(\theta) + \mathbf{M}\lambda)^{\top} \mathbf{Q}(\theta) (\mu(\theta) + \mathbf{M}\lambda) \right)$$



VB corrected marginal posterior of η_j

$$\begin{aligned}\eta_j|\boldsymbol{\theta}, \mathbf{y} &\sim N(\mu_j(\boldsymbol{\theta}), \sigma_j^2(\boldsymbol{\theta})) \\ \mu_j(\boldsymbol{\theta}) &= (\mathbf{A}\boldsymbol{\mu}^*(\boldsymbol{\theta}))_j \\ \tilde{\pi}(\eta_j|\mathbf{y}) &\approx \sum_{k=1}^K \pi_G(\eta_j|\boldsymbol{\theta}_k, \mathbf{y}) \tilde{\pi}(\boldsymbol{\theta}_k|\mathbf{y}) \delta_k.\end{aligned}$$



Cox proportional hazards model

We simulate survival data for n patients using the following very simple Cox proportional hazards model

$$h_i(t) = h_0(t) \exp(\beta x_i) = 1.2t^{0.2} \exp(0.1x_i), \quad i = 1, 2, \dots, n,$$

where x is a scaled and centered continuous covariate, and the baseline hazard, $h_0(t)$ is estimated using a scaled random walk order one model with 50 bins. We also consider four different values of n which are $n = 10^2$, to 10^5 .



Cox proportional hazards model

n	Augmented size	classic INLA (s)	modern INLA (s)
10^2	1 327	1.6	0.1
10^3	12 657	1.3	0.4
10^4	131 807	10.2	2.3
10^5	1 302 413	113.3	22.5

Table: Results from simulation of Cox proportional hazards model



cs-fMRI model

Functional magnetic resonance imaging (fMRI) is a noninvasive neuro-imaging technique used to localize regions of specific brain activity during certain tasks.

For T timepoints and N vertices per hemisphere resulting in data $\mathbf{y}_{TN \times 1}$ with the latent Gaussian model as follows:

$$\begin{aligned} \mathbf{y} | \boldsymbol{\beta}, \mathbf{b}, \boldsymbol{\theta} &\sim N(\boldsymbol{\mu}_y, \mathbf{V}), \quad \boldsymbol{\mu}_y = \sum_{k=0}^K \mathbf{x}_k \boldsymbol{\beta}_k + \sum_{j=1}^J \mathbf{z}_j \mathbf{b}_j \\ \boldsymbol{\beta}_k &= \boldsymbol{\Psi}_k \mathbf{w}_k \quad (\text{SPDE prior on } \boldsymbol{\beta}_k) \\ \mathbf{w}_k | \boldsymbol{\theta} &\sim N(\mathbf{0}, \mathbf{Q}_{\mathcal{T}_k, \kappa_k}^{-1}) \\ \mathbf{b}_j &\sim N(\mathbf{0}, \delta I) \quad (\text{Diffuse priors for } \mathbf{b}_j) \\ \boldsymbol{\theta} &\sim \pi(\boldsymbol{\theta}), \end{aligned}$$

where we have K task signals and J nuisance signals.



cs-fMRI model

The data consists of a 3.5-min fMRI for each subject, consisting of 284 volumes, where each subject performs 5 different motor tasks interceded with a 3 second visual cue. Each hemisphere of the brain contained 32492 surface vertices. From these, 5000 are resampled to use for the analysis. This results in a response data vector \mathbf{y} of size **2 523 624**, with an SPDE model defined on a mesh with 8795 triangles.

The inference based on the modern formulation of INLA was computed in 148 seconds.



cs-fMRI model

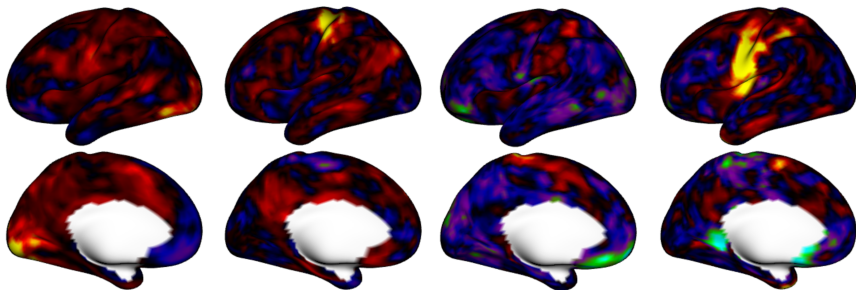


Figure: Activation areas for the different tasks in the left hemisphere - visual cue, right hand motor, right foot motor, tongue motor task (from left to right)



SPDE model

Consider an SPDE model with 948 mesh nodes and data n .

n	Classic INLA (s)	Modern INLA (s)
10^2	2.6	2.7
10^3	2.9	2.5
10^4	13.4	3.7

Table: Number of data points and the computing time (in seconds) considering the classic and modern INLA formulations



Stable prediction for SPDE model

Consider an SPDE model on a mesh with 946 nodes, conditional on 10^4 observations.

grid layout	size	Classic INLA (s)	Modern INLA (s)
250×150	37500	5.39	2.59
500×300	150000	17.70	4.68
1000×600	600000	156.48	13.69

Table: Number of predictions (grid layout and size) and the computing time (in seconds).



Discussion

`inla(..., inla.mode = "experimental")`

- INLA 2.0
- Remove the linear predictors from the latent field → accurate posterior inference with VB correction
- New applications that aren't feasible with INLA 1.0

Thank you • شكرا



جامعة الملك عبد الله
للعلوم والتقنية
King Abdullah University of
Science and Technology