

INTEGRATED NESTED LAPLACE APPROXIMATIONS FOR EXTENDED LATENT GAUSSIAN MODELS

BY ADAM HOWES¹, ALEX STRINGER²
SETH R. FLAXMAN³, JEFFREY W. EATON⁴

¹*Department of Mathematics, Imperial College London, ath19@ic.ac.uk*

²*Department of Statistics and Actuarial Science, University of Waterloo, alex.stringer@uwaterloo.ca*

³*Department of Computer Science, University of Oxford, seth.flaxman@cs.ox.ac.uk*

⁴*MRC Centre for Global Infectious Disease Analysis, School of Public Health, Imperial College London, jeffrey.eaton@imperial.ac.uk*

Naomi is a spatial evidence synthesis model used to produce district-level HIV epidemic indicators in sub-Saharan Africa. Multiple outcomes of interest, including HIV prevalence, HIV incidence and antiretroviral therapy treatment coverage are jointly modelled using both household survey data and routinely reported health system data. We propose a new inference method which combines the simplified integrated nested Laplace approximation approach of Wood (2020) with adaptive Gauss-Hermite quadrature to enable fast and accurate inference for Naomi and other extended latent Gaussian models. Using data from Malawi, our method provides substantially more accurate inferences than the empirical Bayes Gaussian approximation approach used currently, and is comparable to Hamiltonian Monte Carlo with the No-U-Turn sampler. By extending the `aghq` R package we facilitate flexible and easy use of our method when provided a TMB C++ template for the model's log-posterior.

1. Introduction. Mounting an effective public health response to the HIV epidemic requires accurate, timely HIV indicator estimates at a sufficiently fine-scale resolution. Producing these estimates is a challenging task, as all available data sources have shortcomings. Nationally-representative household surveys are the most statistically reliable data source, but due to their high cost to run, in most countries they only occur infrequently. Other data sources, such as routine health surveillance of antenatal care clinics, are more real-time but based on a biased sample of the population. To meet these challenges, the Naomi small-area estimation model (Eaton et al., 2021) synthesises data from multiple sources to estimate HIV prevalence, HIV incidence, and coverage of antiretroviral treatment (ART) at a district-level. Software (<https://naomi.unaids.org>) has been developed for Naomi, which allows countries to input their data and generate estimates in a yearly process supported by UNAIDS.

The complexity of the model presents a difficult Bayesian inference problem. Any inferential strategy must be fast, as well as easy to run in production by country teams, ruling out prohibitively slow Markov chain Monte Carlo (MCMC) approaches. Inference is currently conducted using an empirical Bayes approach, with a Gaussian approximation to the latent field, via the Template Model Builder (TMB) R package (Kristensen et al., 2016). Owing to its speed and flexibility, TMB is gaining popularity spatial statistics (Osgood-Zimmerman and Wakefield, 2021). Inference in TMB is based on optimisation of a C++ template function, with the option available to use a Laplace approximation to integrate out any subset of the function arguments. For the Naomi model, we use this option to integrate out the latent field parameters. Taking inspiration from the AD Model Builder (ADMB) package (Fournier et al.,

Keywords and phrases: spatial statistics, small-area estimation, INLA, AGHQ, HIV epidemiology.

2012), TMB uses automatic differentiation (Baydin et al., 2017) to calculate the derivatives required for numerical optimisation routines and the Laplace approximation. Although this approach has favourable computational properties, we have found the inferences generated for the Naomi model to sometimes be inaccurate.

To obtain fast, accurate inferences for the Naomi model we develop a new inference methodology which combines the simplified integrated nested Laplace approximation (INLA) approach of Wood (2020) with adaptive Gauss-Hermite quadrature (AGHQ). INLA is an approach to approximate Bayesian inference based on nested Laplace approximations and numerical quadrature. The central innovation of Rue, Martino and Chopin (2009) is a way to approximate accurate latent field posterior marginals without explicitly computing the full Laplace approximation for each element. Simplified INLA (Wood, 2020) extends INLA by relaxing the sparsity assumptions on the latent field required for this approximation to be accurate. This extension facilitates inference for models like Naomi, which fall within the extended latent Gaussian model (ELGM) (Stringer, Brown and Stafford, 2022) class, and were not previously amenable to inference with INLA. ELGMs build on latent Gaussian models (LGMs) by allowing each element of the linear predictor to depend on any subset of elements from the latent field. We combine simplified INLA with AGHQ, a quadrature rule based on the theory of polynomial interpolation which adapts to the integrand based on the Hessian at the mode. Though no theory yet exists for the nested case, the first stochastic convergence results for adaptive quadrature rules were recently obtained by Bilodeau, Stringer and Tang (2021) using AGHQ. We implement our method as an extension of the `aghq` R package (Stringer, 2021). As `aghq` is designed to naturally interface with TMB, use of the method is easy when provided a C++ user template for the log-posterior.

The remainder of this paper is organised as follows. In Section 2 we describe the Naomi model. Section 3 outlines our approach to fast, accurate Bayesian inference using simplified INLA and AGHQ. As a case-study, we fit the Naomi model on data from Malawi, and compare the accuracy of inferences in Section 2. In this section, we also demonstrate a Bayesian workflow. Finally, in Section 5 we discuss our conclusions, how our method might be used in other models, and directions for future research.

2. The Naomi model. Eaton et al. (2019) specify a joint model linking small-area estimation models of HIV prevalence from household surveys, HIV prevalence from antenatal care clinics, and antiretroviral therapy (ART) coverage from routine health data collection. This model forms the basis of the Naomi small-area estimation model, described fully in Eaton et al. (2021). Modelling data from multiple sources concurrently is attractive as it increases statistical power, mitigates the biases of any single source, and prompts investigation into any data conflicts. The model is comprised of three components, described as follows.

2.1. Household survey component. Consider a country partitioned into n areas indexed by i . Suppose a simple random household survey of m_i^{HS} people is conducted in each area, and y_i^{HS} HIV positive cases are observed. Then, cases may be modelled using a binomial logistic regression model

$$y_i^{\text{HS}} \sim \text{Bin}(m_i^{\text{HS}}, \rho_i^{\text{HS}}),$$

$$\text{logit}(\rho_i^{\text{HS}}) \sim \mathcal{N}(\beta_\phi, \sigma_\phi^2),$$

where ρ_i^{HS} is the HIV prevalence, modelled on the logit scale by a Gaussian with mean $\beta_\phi \sim p(\beta_\phi)$ and standard deviation $\sigma_\phi \sim p(\sigma_\phi)$.

2.2. *ANC component.* Routinely collected data from pregnant women attending antenatal care clinics (ANCs) is another important source of information about the HIV epidemic. Suppose that of m_i^{ANC} women attending ANC, y_i^{ANC} are HIV positive. Then an analogous binomial logistic regression model

$$\begin{aligned} y_i^{\text{ANC}} &\sim \text{Bin}(m_i^{\text{ANC}}, \rho_i^{\text{ANC}}), \\ \text{logit}(\rho_i^{\text{ANC}}) &= \text{logit}(\rho_i^{\text{HS}}) + b_i, \\ b_i &\sim \mathcal{N}(\beta_b, \sigma_b^2), \end{aligned}$$

may be used to describe HIV prevalence amongst the sub-population of women attending ANCs. Reflecting the fact that prevalence in ANCs is related, but importantly different, to prevalence in the general population, bias terms b_i are used to offset ANC prevalence from HIV prevalence on the logit scale. We use a Gaussian distribution with mean $\beta_b \sim p(\beta_b)$ and standard deviation $\sigma_b \sim p(\sigma_b)$ for these bias random effects.

2.3. *ART component.* The number of people receiving treatment at district health facilities A_i provides additional information about HIV prevalence. Districts with high prevalence are likely to have a greater number of people receiving treatment, and vice versa districts with low prevalence are likely to have fewer people receiving treatment. ART coverage, defined to be the proportion of people living with HIV (PLHIV) currently on ART on district i , is given by $\alpha_i = A_i / \rho_i^{\text{HS}} N_i$, where N_i is the total population of district i and assumed to be constant. As such, ART coverage may also be modelled using a binomial logistic regression model

$$\begin{aligned} A_i &\sim \text{Bin}(N_i, \rho_i^{\text{HS}} \alpha_i), \\ \text{logit}(\alpha_i) &\sim \mathcal{N}(\beta_\alpha, \sigma_\alpha^2), \end{aligned}$$

where the proportion of people receiving ART is $\rho_i^{\text{HS}} \alpha_i$. We again use a Gaussian distribution with mean $\beta_\alpha \sim p(\beta_\alpha)$ and standard deviation $\sigma_\alpha \sim p(\sigma_\alpha)$. Here we assume no travel between districts to receive treatment.

2.4. *Joint model.* The three components above may be combined to a single model as follows

$$\begin{aligned} y_i^{\text{HS}} &\sim \text{Bin}(m_i^{\text{HS}}, \rho_i^{\text{HS}}), \\ y_i^{\text{ANC}} &\sim \text{Bin}(m_i^{\text{ANC}}, \rho_i^{\text{ANC}}), \\ A_i &\sim \text{Bin}(N_i, \rho_i^{\text{HS}} \alpha_i), \\ \text{logit}(\rho_i^{\text{HS}}) &\sim \mathcal{N}(\beta_\phi, \sigma_\phi^2), \\ \text{logit}(\rho_i^{\text{ANC}}) &= \text{logit}(\rho_i^{\text{HS}}) + b_i, \\ b_i &\sim \mathcal{N}(\beta_b, \sigma_b^2), \\ \text{logit}(\alpha_i) &\sim \mathcal{N}(\beta_\alpha, \sigma_\alpha^2). \end{aligned}$$

3. Fast inference methods. Consider a latent Gaussian model (LGM) of the form

$$\begin{aligned} \text{(Observations)} \quad & \mathbf{y} \sim p(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}), \\ \text{(Latent field)} \quad & \mathbf{x} \sim \mathcal{N}(\mathbf{x} | \mathbf{0}, \mathbf{Q}(\boldsymbol{\theta})^{-1}), \\ \text{(Parameters)} \quad & \boldsymbol{\theta} \sim p(\boldsymbol{\theta}), \end{aligned}$$

where $\dim(\mathbf{y}) = \dim(\mathbf{x}) = n$ and $\dim(\boldsymbol{\theta}) = m$, and $m < n$. The joint posterior of $(\mathbf{x}, \boldsymbol{\theta})$ is given by

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

We consider approximations to the posterior marginals of each latent random variable x_i and parameter θ_j given by

$$(3.1) \quad \tilde{p}(x_i | \mathbf{y}) \approx p(x_i | \mathbf{y}) = \int p(x_i, \boldsymbol{\theta} | \mathbf{y}) d\boldsymbol{\theta} = \int p(x_i | \boldsymbol{\theta}, \mathbf{y}) p(\boldsymbol{\theta} | \mathbf{y}) d\boldsymbol{\theta}, \quad i = 1, \dots, n,$$

$$(3.2) \quad \tilde{p}(\theta_j | \mathbf{y}) \approx p(\theta_j | \mathbf{y}) = \int p(\boldsymbol{\theta} | \mathbf{y}) d\boldsymbol{\theta}_{-j} \quad j = 1, \dots, m.$$

3.1. Algorithm. Given a C++ user template `model.cpp` for the negative unnormalised log posterior $-\log p(\mathbf{y}, \mathbf{x}, \boldsymbol{\theta})$, we obtain the posterior marginal approximations $\{\tilde{p}(x_i | \mathbf{y})\}_{i=1}^n$ and $\tilde{p}(\theta_j | \mathbf{y})_{j=1}^m$ via the following algorithm, comprised of nested applications of Laplace approximation and adaptive Gauss-Hermite quadrature.

1. Use a Laplace approximation to obtain the unnormalised $\tilde{p}_{\text{LA}}(\boldsymbol{\theta}, \mathbf{y})$

$$\tilde{p}_{\text{LA}}(\boldsymbol{\theta}, \mathbf{y}) = \frac{p(\mathbf{y}, \mathbf{x}, \boldsymbol{\theta})}{\tilde{p}_{\text{G}}(\mathbf{x} | \boldsymbol{\theta}, \mathbf{y})} \Big|_{\mathbf{x}=\hat{\mathbf{x}}(\boldsymbol{\theta})}$$

where $\tilde{p}_{\text{G}}(\mathbf{x} | \boldsymbol{\theta}, \mathbf{y}) = \mathcal{N}(\mathbf{x} | \hat{\mathbf{x}}(\boldsymbol{\theta}), \mathbf{H}(\boldsymbol{\theta})^{-1})$ is a Gaussian approximation to $p(\mathbf{x} | \boldsymbol{\theta}, \mathbf{y})$ with mode and precision matrix given by

$$\hat{\mathbf{x}}(\boldsymbol{\theta}) = \arg \min_{\mathbf{x}} -\log p(\mathbf{y}, \mathbf{x}, \boldsymbol{\theta}),$$

$$\mathbf{H}(\boldsymbol{\theta}) = \frac{\partial^2}{\partial \mathbf{x} \partial \mathbf{x}^\top} -\log p(\mathbf{y}, \mathbf{x}, \boldsymbol{\theta}) \Big|_{\mathbf{x}=\hat{\mathbf{x}}(\boldsymbol{\theta})}$$

2. Normalise $\tilde{p}_{\text{LA}}(\boldsymbol{\theta}, \mathbf{y})$ using adaptive Gauss-Hermite quadrature to obtain

$$\tilde{p}_{\text{AQ}}(\boldsymbol{\theta} | \mathbf{y}) = \frac{\tilde{p}_{\text{LA}}(\boldsymbol{\theta}, \mathbf{y})}{\tilde{p}_{\text{AQ}}(\mathbf{y})},$$

where the normalising constant is calculated using nodes from a Gauss-Hermite quadrature rule $\mathbf{z} \in \mathcal{Q}(m, k)$ with $m = \dim(\boldsymbol{\theta})$, k nodes per dimension, and weights $\omega : \mathbf{z} \in \mathcal{Q}(m, k) \rightarrow \mathbb{R}$ as

$$\tilde{p}_{\text{AQ}}(\mathbf{y}) = \sum_{\mathbf{z} \in \mathcal{Q}(m, k)} \tilde{p}_{\text{LA}}(\boldsymbol{\theta}(\mathbf{z}), \mathbf{y}) \omega(\mathbf{z}).$$

The nodes \mathbf{z} are adapted based on the mode and curvature at the mode of the Laplace approximation as follows

$$\boldsymbol{\theta}(\mathbf{z}) = \hat{\boldsymbol{\theta}} + \mathbf{L}\mathbf{z},$$

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} \tilde{p}_{\text{LA}}(\boldsymbol{\theta}, \mathbf{y}),$$

$$\mathbf{H} = \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} -\log \tilde{p}_{\text{LA}}(\boldsymbol{\theta}, \mathbf{y}) \Big|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}}$$

$$\mathbf{H}^{-1} = \mathbf{L}\mathbf{L}^\top.$$

We typically set $k = 3$ such that there are 3^m nodes in total.

3. Obtain an unnormalised nested approximation to the posterior marginal of the i th latent effect by

$$\tilde{p}_{\text{LA}}(x_i, \mathbf{y}) = \sum_{\mathbf{z} \in \mathcal{Q}(m, k)} \tilde{p}_{\text{LA}}(x_i, \boldsymbol{\theta}(\mathbf{z}), \mathbf{y}) \omega(\mathbf{z}).$$

The nodes and weights $\{\mathcal{Q}(m, k), \omega\}$ used to obtain $\tilde{p}_{\text{AQ}}(\mathbf{y})$ are reused to perform integration with respect to the hyperparameters above. For each of the k^m values of $\boldsymbol{\theta}(\mathbf{z})$ we obtain $\tilde{p}_{\text{AQ}}(x_i | \boldsymbol{\theta}(\mathbf{z}), \mathbf{y})$ by setting $\boldsymbol{\theta} = \boldsymbol{\theta}(\mathbf{z})$ in the following Laplace approximation

$$\tilde{p}_{\text{LA}}(x_i, \boldsymbol{\theta}, \mathbf{y}) = \frac{p(x_i, \mathbf{x}_{-i}, \boldsymbol{\theta}, \mathbf{y})}{p_{\text{G}}(\mathbf{x}_{-i} | x_i, \boldsymbol{\theta}, \mathbf{y})} \Big|_{\mathbf{x}_{-i} = \hat{\mathbf{x}}_{-i}(x_i, \boldsymbol{\theta})}$$

where $\tilde{p}_{\text{G}}(\mathbf{x} | \boldsymbol{\theta}, \mathbf{y}) = \mathcal{N}(\mathbf{x} | \hat{\mathbf{x}}_{-i}(x_i, \boldsymbol{\theta}), \mathbf{H}_{-i, -i}(x_i, \boldsymbol{\theta})^{-1})$ is a Gaussian approximation to $p(\mathbf{x} | \boldsymbol{\theta}, \mathbf{y})$ with mode and precision matrix given by

$$\begin{aligned} \hat{\mathbf{x}}_{-i}(x_i, \boldsymbol{\theta}) &= \arg \min_{\mathbf{x}_{-i}} -\log p(\mathbf{y}, x_i, \mathbf{x}_{-i}, \boldsymbol{\theta}), \\ \mathbf{H}_{-i, -i}(x_i, \boldsymbol{\theta}) &= \frac{\partial^2}{\partial \mathbf{x}_{-i} \partial \mathbf{x}_{-i}^{\top}} -\log p(\mathbf{y}, x_i, \mathbf{x}_{-i}, \boldsymbol{\theta}) \Big|_{\mathbf{x}_{-i} = \hat{\mathbf{x}}_{-i}(x_i, \boldsymbol{\theta})} \end{aligned}$$

4. Normalise $\tilde{p}_{\text{LA}}(x_i, \mathbf{y})$ using $\tilde{p}_{\text{AQ}}(\mathbf{y})$ to obtain

$$\tilde{p}_{\text{AQ}}(x_i | \mathbf{y}) = \frac{\tilde{p}_{\text{LA}}(x_i, \mathbf{y})}{\tilde{p}_{\text{AQ}}(\mathbf{y})}.$$

4. Application to the Naomi model.

- The R ([R Core Team, 2021](https://www.r-project.org/)) code used to produce all results we describe is available at github.com/athowes/elgm-inf. The inference method is available in versions 0.5.0. onwards of the `aghq` package
- Using the TMB template, we fit the model using four inferential approaches: (1) empirical Bayes combined with a Gaussian approximation, (2) AGHQ combined with a Gaussian approximation, (3) AGHQ combined with a Laplace approximation and (4) the Hamiltonian Monte Carlo (HMC) algorithm No-U-Turn Sampling (NUTS). We treat results from NUTS as the gold-standard
- We used the Kolmogorov-Smirnov test for the maximum difference between two empirical cumulative distribution functions to compare posterior marginal distributions
- We performed posterior predictive checks to assess the coverage of our estimates via the uniformity of the data within each posterior marginal distribution

5. Discussion.

- We developed an approximate Bayesian inference algorithm to solve a challenging problem in the small-area estimation of HIV in low resource settings
- The flexibility of our method implementation, including compatibility with any TMB C++ template, allows broader use, as well as investigation of, deterministic inference methods than had previously been possible
- We demonstrated a Bayesian workflow for deterministic inference methods

Acknowledgements. AH was supported by the EPSRC Centre for Doctoral Training in Modern Statistics and Statistical Machine Learning (EP/S023151/1).

REFERENCES

- BAYDIN, A. G., PEARLMUTTER, B. A., RADUL, A. A. and SISKIND, J. M. (2017). Automatic differentiation in machine learning: a survey. *The Journal of Machine Learning Research* **18** 5595–5637.
- BILODEAU, B., STRINGER, A. and TANG, Y. (2021). Stochastic convergence rates and applications of adaptive quadrature in Bayesian inference. *arXiv preprint arXiv:2102.06801*.
- EATON, J. W., BAJAJ, S., JAHN, A., KALUA, T., MGANGA, A., AULD, A. F., KIM, E., PAYNE, D., SHIRAIISHI, R. W., GUTREUTER, S., HALLETT, T. B. and JOHNSON, L. F. (2019). Joint small-area estimation of HIV prevalence, ART coverage and HIV incidence. *Working paper*.
- EATON, J. W., DWYER-LINDGREN, L., GUTREUTER, S., O'DRISCOLL, M., STEVENS, O., BAJAJ, S., ASHTON, R., HILL, A., RUSSELL, E., ESRA, R., DOLAN, N., ANIFOWOSHE, Y. O., WOODBRIDGE, M., FELLOWS, I., GLAUBIUS, R., HAEUSER, E., OKONEK, T., STOVER, J., THOMAS, M. L., WAKEFIELD, J., WOLOCK, T. M., BERRY, J., SABALA, T., HEARD, N., DELGADO, S., JAHN, A., KALUA, T., CHIMPANDULE, T., AULD, A., KIM, E., PAYNE, D., JOHNSON, L. F., FITZJOHN, R. G., WANYEKI, I., MAHY, M. I. and SHIRAIISHI, R. W. (2021). Naomi: a new modelling tool for estimating HIV epidemic indicators at the district level in sub-Saharan Africa. *Journal of the International AIDS Society* **24** e25788.
- FOURNIER, D. A., SKAUG, H. J., ANCHETA, J., IANELLI, J., MAGNUSSON, A., MAUNDER, M. N., NIELSEN, A. and SIBERT, J. (2012). AD Model Builder: using automatic differentiation for statistical inference of highly parameterized complex nonlinear models. *Optimization Methods and Software* **27** 233–249.
- KRISTENSEN, K., NIELSEN, A., BERG, C. W., SKAUG, H., BELL, B. M. et al. (2016). TMB: Automatic Differentiation and Laplace Approximation. *Journal of Statistical Software* **70**.
- OSGOOD-ZIMMERMAN, A. and WAKEFIELD, J. (2021). A Statistical Introduction to Template Model Builder: A Flexible Tool for Spatial Modeling.
- RUE, H., MARTINO, S. and CHOPIN, N. (2009). Approximate Bayesian inference for latent Gaussian models by using integrated nested Laplace approximations. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **71** 319–392.
- STRINGER, A. (2021). Implementing Approximate Bayesian Inference using Adaptive Quadrature: the aghq Package. *arXiv preprint arXiv:2101.04468*.
- STRINGER, A., BROWN, P. and STAFFORD, J. (2022). Fast, scalable approximations to posterior distributions in extended latent Gaussian models. *Journal of Computational and Graphical Statistics* 1–15.
- R CORE TEAM (2021). R: A Language and Environment for Statistical Computing R Foundation for Statistical Computing, Vienna, Austria.
- WOOD, S. N. (2020). Simplified integrated nested Laplace approximation. *Biometrika* **107** 223–230.