

INTEGRATED NESTED LAPLACE APPROXIMATIONS FOR EXTENDED LATENT GAUSSIAN MODELS WITH APPLICATION TO THE NAOMI HIV MODEL

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 policy 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 which is currently in use, 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 to make targeted interventions. Producing these estimates is a challenging task because all available data sources have shortcomings which must be overcome. Nationally-representative household surveys provide the most statistically reliable data, but due to their high cost to run, in most countries they occur only infrequently. Other data sources, such as routine health surveillance of antenatal care clinics, are more real-time but based on limited or biased samples 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. Simultaneous modelling of the data mitigates the limitations of any single source, increases statistical power, and prompts investigation into any conflicting information. Software (<https://naomi.unaids.org>) has been developed for Naomi, allowing countries to input their data and interactively generate estimates in a yearly process supported by UNAIDS. The creation of estimates by country teams, rather than external agencies, is a noteworthy feature of the HIV response. Drawing on expertise closest to the data generating process, improves the accuracy of the process, as well as strengthening trust and ownership of the estimates.

The use case requirements for the model, combined with its relative complexity, present a difficult Bayesian inference problem. Any inferential strategy must be fast, as well as easy to

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

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 has recently been gaining popularity more broadly in 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, which for the Naomi model, we use for the high-dimensional latent field parameters. Taking inspiration from the AD Model Builder (ADMB) package (Fournier et al., 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. This has motivated us to look for a more accurate approach, which is flexible enough to be compatible with the model and fast enough to be run in production by country teams.

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 an approximation which enables 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). Since `aghq` is designed to naturally interface with TMB, use of our method is simple when provided a C++ user template for the log-posterior.

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

2. A simplified Naomi model. Eaton et al. (2021) 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. The model is defined over three time points: T_1 the time of the most recent household survey with HIV testing; T_2 , the current time period; and T_3 , a short term projection period. We consider a simplified version of the model which is defined only at T_1 omitting temporal projection to T_2 and T_3 . Below we provide an overview of the three components in

the simplified model, highlighting the aspects which make it a challenge for existing inferential approaches. A more complete mathematical description of the simplified model, as well as a C++ template for the log-posterior, are provided in the appendix.

2.1. Household survey component. Consider a country in sub-Saharan Africa where a household survey has taken place during the quarter T_1 . Let x index district, a five-year age band, s sex and t time. For ease of notation, let i index the coarsest district-age-sex division included in the model. The data we observe may be aggregated over indices i , so we let \mathcal{I} be a set of i for which observations are reported. Let $N_i \in \mathbb{N}$ be the population size, $\rho_i \in [0, 1]$ be HIV prevalence, $\alpha_i \in [0, 1]$ be ART coverage, and $\lambda_i > 0$ be the annual HIV incidence rate. We specify independent mixed effects models for HIV prevalence and ART coverage in the general population on the logit scale such that

$$\begin{aligned}\text{logit}(\rho_i) &= \eta_i^\rho, \\ \text{logit}(\alpha_i) &= \eta_i^\alpha,\end{aligned}$$

for certain choice of linear predictors η_i^ρ and η_i^α . For the HIV incidence rate we use a mixed effects model on the log scale

$$\log(\lambda_i) = \eta_i^\lambda(\{\rho_i, \alpha_i\}_{i \in \mathcal{I}}),$$

where the linear predictor depends on $\{\rho_i, \alpha_i\}_{i \in \mathcal{I}}$ for some \mathcal{I} . Let κ_i be the proportion recently infected among HIV positive persons. For each set of observed strata indices \mathcal{I} , we calculate the weighted observations $\hat{\theta}_{\mathcal{I}}$ for $\theta \in \{\rho, \alpha, \kappa\}$ with respective Kish effective sample sizes

$$M_{\mathcal{I}}^{\hat{\theta}} = \frac{\left(\sum_j w_j\right)^2}{\sum_j w_j^2},$$

where j indexes individuals across all strata $i \in \mathcal{I}$, with corresponding survey weights w_j . The observed number of indicator cases is then

$$Y_{\mathcal{I}}^{\hat{\theta}} = M_{\mathcal{I}}^{\hat{\theta}} \cdot \hat{\theta}_{\mathcal{I}}.$$

For $\theta \in \{\rho, \alpha, \kappa\}$ we model these aggregate observations using a binomial working likelihood

$$Y_{\mathcal{I}}^{\hat{\theta}} \sim \text{xBin}(M_{\mathcal{I}}^{\hat{\theta}}, \theta_{\mathcal{I}}),$$

where $\theta_{\mathcal{I}}$ are the following appropriately weighted aggregates

$$\begin{aligned}\rho_{\mathcal{I}} &= \frac{\sum_{i \in \mathcal{I}} N_i \rho_i}{\sum_{i \in \mathcal{I}} N_i}, \\ \alpha_{\mathcal{I}} &= \frac{\sum_{i \in \mathcal{I}} N_i \rho_i \alpha_i}{\sum_{i \in \mathcal{I}} N_i \rho_i}, \\ \kappa_{\mathcal{I}} &= \frac{\sum_{i \in \mathcal{I}} N_i \rho_i \kappa_i}{\sum_{i \in \mathcal{I}} N_i \rho_i}.\end{aligned}$$

To link the proportion recently infected among HIV positive persons κ_i to HIV incidence λ_i we use

$$\kappa_i = 1 - \exp\left(-\lambda \cdot \frac{1 - \rho_i}{\rho_i} \cdot (\Omega_T - \beta_T) - \beta_T\right),$$

where the mean duration of recent infection Ω_T and is the false recent ratio β_T are strongly informed by priors for the particular survey.

2.2. ANC testing component. We model HIV prevalence ρ_i^{ANC} and ART coverage α_i^{ANC} among pregnant women as being offset on the logit scale from the general population indicator as follows

$$\begin{aligned}\text{logit}(\rho_i^{\text{ANC}}) &= \text{logit}(\rho_i) + \eta_i^{\rho^{\text{ANC}}}, \\ \text{logit}(\alpha_i^{\text{ANC}}) &= \text{logit}(\alpha_i) + \eta_i^{\alpha^{\text{ANC}}}.\end{aligned}$$

We inform these processes using the following aggregate ANC data from the year of the most recent survey: the number of ANC clients with ascertained status $X_{\mathcal{I}}^{\text{ANC}}$, the number of those with positive status $Y_{\mathcal{I}}^{\text{ANC}}$, and the number of ANC clients already on ART prior to their first ANC visit $Z_{\mathcal{I}}^{\text{ANC}}$. We use the binomial working likelihoods

$$\begin{aligned}Y_{\mathcal{I}}^{\text{ANC}} &\sim \text{Bin}(X_{\mathcal{I}}^{\text{ANC}}, \alpha_{\mathcal{I}}^{\text{ANC}}) \\ Z_{\mathcal{I}}^{\text{ANC}} &\sim \text{Bin}(Y_{\mathcal{I}}^{\text{ANC}}, \alpha_{\mathcal{I}}^{\text{ANC}}),\end{aligned}$$

where again we use weighted aggregates

$$\begin{aligned}\rho_{\mathcal{I}}^{\text{ANC}} &= \frac{\sum_{i \in \mathcal{I}} \Psi_i \rho_i^{\text{ANC}}}{\sum_{i \in \mathcal{I}} \Psi_i}, \\ \alpha_{\mathcal{I}}^{\text{ANC}} &= \frac{\sum_{i \in \mathcal{I}} \Psi_i \rho_i^{\text{ANC}} \alpha_i^{\text{ANC}}}{\sum_{i \in \mathcal{I}} \Psi_i \rho_i^{\text{ANC}}},\end{aligned}$$

with Ψ_i the number of pregnant women.

2.3. ART attendance component. Let $\gamma_{x,x'} \in [0, 1]$ be the probability that a person on ART residing in district x receives ART in district x' . We assume that $\gamma_{x,x'} = 0$ unless $x = x'$ or the two districts are adjacent $x \sim x'$.

Let $A_{\mathcal{I}}$ be the number of people receiving ART.

3. Fast approximate 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.

3.2. *Algorithm.*

1. Calculate

$$\begin{aligned}\hat{\boldsymbol{\theta}} &= \arg \max_{\boldsymbol{\theta}} \tilde{p}_{\text{LA}}(\boldsymbol{\theta}, \mathbf{y}), \\ \mathbf{H} &= \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} - \log \tilde{p}_{\text{LA}}(\boldsymbol{\theta}, \mathbf{y})|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}} \\ \mathbf{H}^{-1} &= \mathbf{L}\mathbf{L}^\top,\end{aligned}$$

where

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

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

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

2. Generate a set of nodes $\mathbf{z} \in \mathcal{Q}(m, k)$ and weights $\omega : \mathbf{z} \in \mathcal{Q}(m, k) \rightarrow \mathbb{R}$ from a Gauss-Hermite quadrature rule with $m = \dim(\boldsymbol{\theta})$, k nodes per dimension, as follows

$$\begin{aligned}\boldsymbol{\theta}(\mathbf{z}) &= \hat{\boldsymbol{\theta}} + \mathbf{L}\mathbf{z}, \\ \hat{\boldsymbol{\theta}} &= \arg \max_{\boldsymbol{\theta}} \tilde{p}_{\text{LA}}(\boldsymbol{\theta}, \mathbf{y}), \\ \mathbf{H} &= \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^\top} - \log \tilde{p}_{\text{LA}}(\boldsymbol{\theta}, \mathbf{y})|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}} \\ \mathbf{H}^{-1} &= \mathbf{L}\mathbf{L}^\top.\end{aligned}$$

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

3. Use this quadrature rule to calculate $\tilde{p}_{\text{AQ}}(\mathbf{y})$ as follows

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

4. For $x_i \in \{\dots\}$ and $\boldsymbol{\theta} \in \{\boldsymbol{\theta}(\mathbf{z})\}_{\mathbf{z} \in \mathcal{Q}(m, k)}$ calculate the modes and Hessians

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

where optimisation to obtain $\hat{\mathbf{x}}_{-i}(x_i, \boldsymbol{\theta})$ may be initialised at $\hat{\mathbf{x}}(\boldsymbol{\theta})_{-i}$.

	Run-time	Memory
Empirical Bayes, Gaussian	0	0
AGHQ, Gaussian	0	0
AGHQ, Laplace	0	0
NUTS	0	0

5. For each x_i calculate

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

where

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

and

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

4. Application to data from Malawi. Using a single TMB template, we fit the simplified Naomi model to data from Malawi 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). The R ([R Core Team, 2021](#)) code used to produce all results we describe is available at github.com/athowes/elgm-inf.

4.1. NUTS convergence. In order to treat results from NUTS as the gold-standard we assessed MCMC convergence using. We ran the algorithm

4.2. Model assessment. We performed posterior predictive checks to assess the coverage of our estimates via the uniformity of the data within each posterior marginal distribution.

4.3. Inference comparison. We used three methods to assess the accuracy of posterior distributions produced by each inferential method: (1) Kolmogorov-Smirnov tests, (2) maximum mean discrepancy, and (3) Pareto-smoothed importance sampling.

Let $\{\theta_i\}_{i=1}^n$ be posterior marginal samples with empirical cumulative distribution (ECDF) function $F(\vartheta) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{\theta_i \leq \vartheta}$. The two-sample Kolmogorov-Smirnov (KS) test statistic is given by the maximum absolute difference between two ECDFs. We compare the two KS statistics

$$D_{\text{aghq}} = \sup_{\vartheta} |F_{\text{tmbstan}}(\vartheta) - F_{\text{aghq}}(\vartheta)|,$$

$$D_{\text{TMB}} = \sup_{\vartheta} |F_{\text{tmbstan}}(\vartheta) - F_{\text{TMB}}(\vartheta)|.$$

See results in Table [4.3](#).

	KS	MMD	PSIS
Empirical Bayes, Gaussian	0	0	0
AGHQ, Gaussian	0	0	0
AGHQ, Laplace	0	0	0
NUTS	0	0	0

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
- Following further testing, we anticipate adding our inference method as an option for
- 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 are excited about exploration of its accuracy for the following challenging model structures which could not previously be fit: list of model types here.
- 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., 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 SHIRAISHI, 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.