

Deterministic Bayesian inference methods for the Naomi model

HIV Inference Lab Group Meeting

Adam Howes

Imperial College London

April 2023

Bayes

- As a statistical modeller, the bulk of our job is in constructing a generative model for data y using parameters ϑ
- This is the joint distribution $p(y, \vartheta) = p(y | \vartheta)p(\vartheta)$
- What we actually want is the posterior $p(\vartheta | y)$ which is **just**¹

$$p(\vartheta | y) = \frac{p(y, \vartheta)}{p(y)} = \frac{p(y | \vartheta)p(\vartheta)}{p(y)}$$

- The central problem of Bayesian inference is that we can't compute $p(y)$

$$p(y) = \int p(y, \vartheta) d\vartheta$$

¹I've bolded this with sarcasm in mind

How might you do it?

- If you want to integrate something deterministically, you could use numerical integration methods
- Pick nodes $\vartheta \in \mathcal{Q} \subset \Theta$ and weights $\omega : \Theta \rightarrow \mathbb{R}$ then compute the sum

$$\tilde{p}(y) = \sum_{\vartheta \in \mathcal{Q}} p(y, \vartheta) \omega(\vartheta)$$

Monte Carlo as an example of numerical integration

- Suppose we can sample $\vartheta_i \sim p(y, \vartheta)$ for $i = 1, \dots, N$
- If we set $\omega(\vartheta_i) = 1/N$ for all i then we get a **Monte Carlo** (MC) estimate

$$\tilde{p}(y) = \frac{1}{N} \sum_i p(y, \vartheta_i)$$

- For complicated models² it's not possible to sample directly from $p(y, \vartheta)$, but we can usually sample from a Markov chain which if you squint a bit is good enough (MCMC)

²Or not even that complicated

Monte Carlo is fundamentally unsound

- “Monte Carlo ignores information” according to O’Hagan (1987)
- Suppose $N = 3$ and we sample $\vartheta_1, \vartheta_2, \vartheta_3$ with $\vartheta_2 = \vartheta_3$ then our MC estimate is

$$\tilde{p}(y) = \frac{1}{3} (p(y, \vartheta_1) + p(y, \vartheta_2) + p(y, \vartheta_3))$$

- This is despite the fact that nothing new about the function has been learned by adding $\{\vartheta_3, p(y, \vartheta_3)\}$

Application to HIV survey sampling

- This is a digression but. . .
- Say we're running a household survey, and sample the same individual twice
- We didn't learn anything new about HIV by surveying them again!
- This doesn't just bite for nodes or individuals which are exactly the same: an analogous argument can be made if they are close together and we expect their function evaluations to be similar

⇒ Bayesian quadrature, Bayesian survey design

For some half-baked thoughts, see athowes.github.io/fourth-gen/paper.pdf

Latent variables and hyperparameters

- Quadrature doesn't work very well when $\dim(\vartheta)$ gets even moderately sized
- Previously I had all of the parameters under the symbol ϑ
- What if we split them up as being $\vartheta = (x, \theta)$
- The key part about this is that $\dim(x)$ is big and $\dim(\theta)$ is small

Names for x	Names for θ
Latent variables, random effects, latent field	Hyperparameters, fixed effects

Spatio-temporal statistics

- There is nothing inherently special about spatio-temporal statistics
- We have observations indexed by space $s \in \mathcal{S}$ and time $t \in \mathcal{T}$
- Usually we associate parameters to spatio-temporal locations as well as observations
- This ends up with us having something like $\{x_{s,t}\}$

What's important about this?

1. There might be **a lot** of spatio-temporal locations, so $\dim(x)$ might be pretty big! If you have 100 districts and 10 years, that's already $100 \times 10 = 1000$ parameters
2. Perhaps we're willing to make assumptions about how things vary over space-time³

³Are there any slides about spatial statistics that don't describe Tobler's first law of geography?

Latent Gaussian models

A latent Gaussian model (LGM) (Rue, Martino, and Chopin 2009) looks along these lines:

(Observations)	$y \sim p(y \mid x, \theta),$
(Latent field)	$x \sim \mathcal{N}(x \mid \mu(\theta), Q(\theta)^{-1}),$
(Hyperparameters)	$\theta \sim p(\theta).$

Laplace approximation

- Remember that we wanted to compute

$$p(y) = \int p(y, \vartheta) d\vartheta$$

- One trick for doing this is to pretend $p(\vartheta | y)$ is Gaussian

- Mode $\hat{\vartheta} = \arg \max_{\vartheta} \log p(y, \vartheta)$
- Hessian $H(\hat{\vartheta}) = -\partial_{\vartheta}^2 \log p(y, \vartheta)|_{\vartheta=\hat{\vartheta}}$
- Gaussian approximation $\implies \tilde{p}_G(\vartheta | y) = \mathcal{N}(\vartheta | \hat{\vartheta}, H(\hat{\vartheta})^{-1})$

Laplace approximation

Include a figure here showing what it looks like to approximate some distribution by a Gaussian.

Laplace approximation

- Now

$$p(y) = \frac{p(\vartheta, y)}{p(\vartheta | y)} \approx \frac{p(\vartheta, y)}{p_G(\vartheta | y)}$$

and we can evaluate RHS where we would like, so let's pick the point at which the Gaussian is most accurate, which is $\hat{\vartheta}$

$$p_{\text{LA}}(y) = \frac{p(\vartheta, y)}{p_G(\vartheta | y)} \Big|_{\vartheta=\hat{\vartheta}} = (2\pi)^{\dim(\vartheta)/2} \det(H(\hat{\vartheta}))^{-1/2} p(\hat{\vartheta}, y)$$

Marginal Laplace approximation

- Hey wait a second, is it reasonable to just assume $p(\vartheta | y)$ is Gaussian?
 - No, not in general. But...
1. We just described a class of models (LGMs) where some subset of the parameters (the latent field x) have a Gaussian prior \implies it's a lot more reasonable to think that they would have a marginal posterior which is close to Gaussian
 2. We just talked about how big x is in comparison to θ ! \implies most of the work in our integral can be done using a **marginal Laplace** approximation to get rid of x

Dichotomy in statistical inference methods:

1. Ones which aim to be completely general
2. Ones which aim to "exploit" properties of the problem at hand

We are taking approach 2.

Marginal Laplace approximation

- What does this look like? Instead of assuming $p(\vartheta | y) = p(x, \theta | y)$ is Gaussian we assume $p(x | \theta, y)$ is

$$\tilde{p}_G(x | \theta, y) = \mathcal{N}(x | \hat{x}, H(\hat{x}))^{-1}$$

where $\hat{x} = \hat{x}(\theta)$

- Now the marginal Laplace approximation is

$$p_{\text{LA}}(\theta, y) = \frac{p(x, \theta, y)}{\tilde{p}_G(x | \theta, y)} \Big|_{x=\hat{x}} = (2\pi)^{\dim(x)/2} \det(H(\hat{x}))^{-1/2} p(\hat{x}, \theta, y)$$

Integrated nested Laplace approximation

- Now we can compute $p_{\text{LA}}(\theta, y)$ but what we really want is still $p(y)$
- But hopefully⁴ the dimension of θ is small enough that we can now tackle this with quadrature
- So pick some nodes \mathcal{Q} and a weighting function ω and away we go

$$p(y) \approx \sum_{\theta \in \mathcal{Q}} p_{\text{LA}}(\theta, y) \omega(\theta)$$

- This is the famous integrated nested Laplace approximation (INLA)

⁴Really: hopefully

*My main comment is that several aspects of the computational machinery that is presented by Rue and his colleagues **could benefit from the use of a numerical technique known as automatic differentiation (AD)** ... By the use of AD one could obtain a system that is automatic from a user's perspective... the benefit would be a fast, flexible and easy-to-use system for doing Bayesian analysis in models with Gaussian latent variables*

- Hans J. Skaug (coauthor of TMB), RSS discussion of Rue, Martino, and Chopin (2009)

Thanks for listening!

- Working on a paper “Fast approximate Bayesian inference for small-area estimation of HIV indicators using the Naomi model” based on this work, joint with Alex Stringer (Waterloo) and my PhD supervisors Seth Flaxman (Oxford) and Jeff Eaton (Imperial)
- Let me know if you'd be up for being an early reader!
- Code for this project is at athowes.github.io/elgm-inf

References I

- O'Hagan, Anthony. 1987. "Monte Carlo Is Fundamentally Unsound." *The Statistician*, 247–49.
- Rue, Håvard, Sara Martino, and Nicolas Chopin. 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 (2): 319–92.