

On automating Markov chain Monte Carlo for a class of spatial models

Murali Haran

School of Statistics

The Pennsylvania State University

`mharan@stat.psu.edu`

Luke Tierney

Department of Statistics and Actuarial Science

University of Iowa

`luke@stat.uiowa.edu`

Abstract

Markov chain Monte Carlo (MCMC) algorithms provide a very general recipe for estimating properties of complicated distributions. While their use has become commonplace and there is a large literature on MCMC theory and practice, MCMC users still have to contend with several challenges with each implementation of the algorithm. These challenges include determining how to construct an efficient algorithm, finding reasonable starting values, deciding whether the sample-based estimates are accurate, and determining an appropriate length (stopping rule) for the Markov chain. We describe an approach for resolving these issues in a theoretically sound fashion in the context of spatial generalized linear models, an important class of models that result in challenging posterior distributions. Our approach combines analytical approximations for constructing provably fast mixing MCMC algorithms, and takes advantage of recent developments in MCMC theory. We apply our methods to real data examples, and

find that our MCMC algorithm is automated and efficient. Furthermore, since starting values, rigorous error estimates and theoretically justified stopping rules for the sampling algorithm are all easily obtained for our examples, our MCMC-based estimation is practically as easy to perform as Monte Carlo estimation based on independent and identically distributed draws.

1 Introduction

Markov chain Monte Carlo (MCMC) methods have become standard tools in Bayesian inference and other areas in statistics. When inference is based on some distribution π , the Metropolis-Hastings algorithm provides a general recipe for constructing a Markov chain with π as its stationary distribution. From a careful practitioner's standpoint, however, MCMC-based inference poses several challenges. In addition to developing and fine-tuning an MCMC algorithm that produces accurate sample-based inference quickly, MCMC users also need to determine an appropriate length for the Markov chain. These issues pose a challenge to non-experts since, even for a specific class of models, the MCMC algorithm needs to be tuned carefully to the posterior distribution resulting from each new data set. Also, the commonly used MCMC "convergence diagnostics" used to determine stopping rules for the algorithm may be unreliable (Cowles and Carlin, 1996) and may not always be directly connected to the central goal of MCMC-based inference, which is to estimate properties of the posterior distribution π up to some desired level of accuracy. Furthermore, users may not always have a reasonable approach for finding starting values for the algorithm.

In this paper, we consider approaches for resolving these closely related issues in a rigorous manner. We describe how we can construct provably efficient MCMC algorithms where the MCMC standard errors, which represent the accuracy of sample-based estimates, can be estimated consistently. These MCMC standard errors can, in turn, be used in a theoretically sound approach to determine an appropriate length for the Markov chain using new developments in MCMC output analysis (Flegal et al., 2008; Jones et al., 2006). Our approach also automatically provides reasonable starting values for the MCMC algorithm. Hence, the re-

sulting MCMC-based inference is automated, theoretically sound, and practical because: (i) initial values for the chain are obtained automatically, (ii) the algorithm produces accurate answers quickly and is therefore useful in practice, (iii) an appropriate Monte Carlo sample size (length for the Markov chain) is determined automatically, and (iv) the accuracy of the sample-based estimates can be assessed rigorously.

A key tool we develop for the samplers described here is an accurate heavy-tailed approximation to the posterior distribution of interest. In this paper we describe how one can construct such an approximation and investigate its use in constructing two samplers: a fast mixing independence Metropolis-Hastings chain (Tierney, 1994) and the classic rejection (‘accept-reject’) sampler. We study these algorithms in the context of a popular class of spatial generalized linear models. Naive MCMC samplers are known to mix poorly for these models, and little is known about the theoretical properties of the samplers. These models are closely related in form to the geostatistical models described in Diggle et al. (1998), which are Bayesian versions of generalized linear models (McCullagh and Nelder, 1999) for spatial data. Hence, our approach, or variants of it, may be applicable to several other important statistical models. Our MCMC algorithm is virtually as simple to use as simple Monte Carlo using independent and identically distributed draws. An underlying assumption of our work here is that sample-based inference is of interest, for instance when propagating uncertainties associated with inferences based on one model or model-component to other models, or when modelers are interested in reporting estimates at a level of accuracy that they would like to control. When sample-based inference is not critical, we note that purely analytical approximations (cf. Rue et al., 2009; Tierney and Kadane, 1986) may, of course, be completely ‘automatic’ approaches for approximate inference since they avoid the MCMC issues outlined above.

In Section 2 we provide a general description of the samplers we consider, along with a discussion of relevant MCMC theory. In Section 3 we discuss the class of models to which we apply our methods. In Section 4 we describe a general approach for deriving heavy-tailed approximations in hierarchical models, discuss how such approximations can be used to construct fast mixing MCMC algorithms, and provide theoretical details. Section 5 describes the application of these methods to several real data sets. We conclude with a discussion of

our results in Section 6.

2 Background

This section provides a brief overview of the samplers we consider, along with some relevant theoretical background. In Section 2.1 we provide some basic theory on Markov chain Monte Carlo, which sets the stage for the fast mixing MCMC sampler. In Section 2.2 we briefly review rejection sampling since we are later able to use our heavy-tailed approximation to construct effective rejection samplers in some cases.

2.1 Markov chain Monte Carlo basics

Our goal is typically to estimate expectations with respect to a distribution π . That is, we are interested in $E_\pi g(x) = \int_\Omega g(x)\pi(dx)$, where g is a real-valued π -integrable function on Ω . Consider a Harris-ergodic Markov chain $\mathbf{X} = \{X_1, X_2, \dots\}$ with state space Ω and stationary distribution π (for definitions see Meyn and Tweedie, 1993). If $E_\pi |g(x)| < \infty$, we can appeal to the ergodic theorem

$$\bar{g}_n = \frac{1}{n} \sum_{i=1}^n g(X_i) \rightarrow E_\pi g(x) \text{ with probability 1.}$$

Now let $P^n(x, \cdot)$ be the n -step transition kernel for this chain, so that $P^n(x, A) = P(X_{i+n} \in A \mid X_i = x)$ for $x \in \Omega$ and any measurable set A . Then it also follows that

$$\lim_{n \rightarrow \infty} \|P^n(x, \cdot) - \pi(\cdot)\|_{TV} = 0 \quad \text{for any } x \in \Omega,$$

where the convergence is in total variation distance, which implies convergence in distribution (Billingsley, 1999). However, to estimate standard errors, we need to appeal to a Central Limit Theorem (CLT), which does not hold in general. For a CLT to hold, the rate of convergence of the Markov chain to π is critical. A Markov chain is said to be *geometrically ergodic* if for some constant $t \in (0, 1)$ and π -almost surely finite function $M : \Omega \rightarrow \mathbb{R}^+$ and $n \in \mathbb{N}$

$$\|P^n(x, \cdot) - \pi(\cdot)\|_{TV} \leq M(x) t^n \text{ for any } x \in \Omega.$$

If $M(x)$ is bounded, then \mathbf{X} is *uniformly ergodic*. If \mathbf{X} is uniformly ergodic and $E_\pi g^2 < \infty$, we have a CLT

$$\sqrt{n}(\bar{g}_n - E_\pi g) \xrightarrow{d} \mathbf{N}(0, \sigma_g^2)$$

as $n \rightarrow \infty$ with $\sigma_g^2 = \text{var}_\pi\{g(X_1)\} + 2 \sum_{i=2}^{\infty} \text{cov}_\pi\{g(X_1), g(X_i)\}$. For a review of Markov chain CLTs under other sets of sufficient conditions see Jones (2004) and Roberts and Rosenthal (2004). For a general review of MCMC theory refer to Tierney (1994).

If the CLT holds as above then Jones et al. (2006) provide a consistent estimate of σ_g^2 , the MCMC standard error in \bar{g}_n , which can then be used to determine how long to run the Markov chain. In the ‘fixed width’ approach to determining chain length, the user stops the simulation when the width of a confidence interval based on an ergodic average is less than a user-specified value (Flegal et al., 2008; Jones et al., 2006). Thus, *for a uniformly ergodic chain, sample-based inference is comparable to i.i.d. Monte Carlo in many ways: a CLT holds under similar conditions, a consistent estimate of Monte Carlo standard errors is easily obtained, and the error estimate can be used to determine a stopping rule for the sampler.* Note that it is challenging to construct uniformly ergodic Markov chains for real problems, and it is generally not easy to prove that a given Markov chain is uniformly or geometrically ergodic; most such proofs have relied on establishing drift and minorization conditions (cf. Hobert and Geyer, 1998).

We briefly describe the consistent batch means approach to calculating Monte Carlo standard errors. Suppose the Markov chain \mathbf{X} is run for a total of $n = ab$ iterations (hence a and b are implicit functions of n) and define

$$\bar{Y}_j := \frac{1}{b} \sum_{i=(j-1)b+1}^{jb} g(X_i) \quad \text{for } j = 1, \dots, a.$$

The batch means estimate of σ_g^2 is

$$\hat{\sigma}_g^2 = \frac{b}{a-1} \sum_{j=1}^a (\bar{Y}_j - \bar{g}_n)^2.$$

Jones et al. (2006) showed that if the batch size and the number of batches are allowed to increase as the overall length of the simulation increases by setting $b_n = \lfloor \sqrt{n} \rfloor$ and $a_n = \lfloor n/b_n \rfloor$ then $\hat{\sigma}_g^2 \rightarrow \sigma_g^2$ with probability 1 as $n \rightarrow \infty$. This is called consistent batch means

(CBM) to distinguish it from the standard (fixed number of batches) version. This is an attractive approach to estimating standard errors since it is easy to compute and holds under the regularity conditions that the chain is uniformly ergodic and $E_\pi |g|^2 < \infty$ (though these are not the only set of sufficient conditions; see Jones et al. (2006) for details).

2.2 Rejection sampling

Rejection or ‘accept-reject sampling’ (von Neumann, 1951) is a well established, simple but powerful method for generating random variates from a given distribution π with support Ω (also see Robert and Casella, 2005). Assume we have a proposal distribution r so that we can draw random samples from r , and we know B such that

$$\text{ess sup}_{x \in \Omega} \frac{\pi(x)}{r(x)} < B, \text{ for some } B < \infty, \quad (1)$$

where ‘ess sup’, the essential supremum (the supremum over all but a set of measure zero), is taken with respect to π . The accept-reject sampling algorithm is as follows:

- Draw $X \sim r$ and draw $U \sim \text{Uniform}(0,1)$.
- If $U \leq \frac{\pi(X)}{r(X)B}$ return X , else do not return X .

Values returned by the above algorithm are distributed according to π . Note that we only need to know both π and r up to a constant of proportionality, that is, we could replace $\pi(x)$ and $r(x)$ with unnormalized functions $h(x)$ and $q(x)$ where $h(x) \propto \pi(x)$ and $q(x) \propto r(x)$. However, in addition to satisfying (1), we also need a specific value of B that satisfies this condition. These are stringent requirements, and explains why rejection sampling is rarely considered a practical option for sample-based inference for realistic Bayesian models.

3 Generalized linear models for spatial data

Spatial generalized linear models are very convenient models for spatial data when the sampling mechanism is known to be non-Gaussian. The spatial dependence can be modeled via Gaussian processes (Diggle et al., 1998) or Gaussian Markov random fields (GMRFs) (cf. Rue and Held, 2005). For brevity, we only describe GMRF-based models for count data as

this is the example used later in this paper. Other models such as Gaussian process-based models may be specified in analogous fashion.

Consider the following hierarchical spatial model for areal data (data arising as sums or averages over geographic regions): the count in region i , denoted by Y_i for $i = 1, \dots, N$, is modeled as a Poisson random variable with mean $E_i \exp(\mu_i)$. E_i , assumed fixed and known, is the count in region i when assuming constant rates for all regions and is typically the product of the population of the i th region and the overall rate in the entire study region. μ_i is the log-relative risk specific to the i th region. Hence the Y_i s are modeled as conditionally independent random variables,

$$Y_i \mid \mu_i \sim \text{Poi}(E_i e^{\mu_i}), \quad i = 1, \dots, N, \quad (2)$$

with μ_i modeled linearly as $\mu_i = \theta_i + \phi_i$, and $\boldsymbol{\theta} = (\theta_1, \dots, \theta_N)^T$ and $\boldsymbol{\phi} = (\phi_1, \dots, \phi_N)^T$ are vectors of random effects. The θ_i 's are independent and identically distributed normal variables, while the ϕ_i 's are assumed to follow a GMRF. In this way, each θ_i captures the i th region's extra-Poisson variability due to area-wide heterogeneity, while each ϕ_i captures the i th region's excess variability attributable to regional clustering. These distributions are specified as follows:

$$\theta_i \mid \tau_h \sim N(0, 1/\tau_h), \text{ and } \phi_i \mid \phi_{j \neq i} \sim N(\mu_{\phi_i}, \sigma_{\phi_i}^2), \quad i = 1, \dots, N,$$

$$\text{where } \mu_{\phi_i} = \frac{\sum_{j \neq i} w_{ij} \phi_j}{\sum_{j \neq i} w_{ij}} \text{ and } \sigma_{\phi_i}^2 = \frac{1}{\tau_c \sum_{j \neq i} w_{ij}}.$$

The μ_{ϕ_i} for a region i is thus a weighted average of the clustering parameters in other regions. Here we use the most common weighting, where w_{ij} is taken as 1 if regions i and j are immediate neighbours, and 0 otherwise. This improper prior can be written as

$$\boldsymbol{\phi} \mid \tau_c \propto \tau_c^{(N-1)/2} \exp\left(-\frac{\tau_c}{2} \boldsymbol{\phi}^T Q \boldsymbol{\phi}\right)$$

where

$$Q_{ij} = \begin{cases} n_i & \text{if } i = j \\ 0 & \text{if } i \text{ is not adjacent to } j \\ -1 & \text{if } i \text{ is adjacent to } j \end{cases}$$

The model specification is completed by specifying priors on the precision (inverse variance) parameters, $\tau_h \sim \text{Gamma}(\alpha_h, \beta_h)$ and $\tau_c \sim \text{Gamma}(\alpha_c, \beta_c)$. Inference for this model is based

on the $2N + 2$ dimensional posterior distribution $\pi(\boldsymbol{\theta}, \boldsymbol{\phi}, \tau_h, \tau_c \mid \mathbf{Y})$, where $\mathbf{Y} = (Y_1, \dots, Y_n)^T$. The posterior distribution is provided in Appendix A.

This model was used by Besag et al. (1991) for Bayesian image restoration, and has since become popular in disease mapping (Mollié, 1996). GMRFs are very widely used in hierarchical Bayes models (see for example Banerjee et al., 2004). The distributions arising from such models are challenging enough so standard MCMC samplers for the posterior distributions are known to mix poorly. A variety of improved MCMC approaches for these models have been studied in Knorr-Held and Rue (2002) and Haran et al. (2003) but their theoretical properties are not known. The issues described in Section 1 are hence unresolved: determining an appropriate chain length is difficult, and MCMC standard errors of estimates are not estimated rigorously.

4 Constructing automated samplers

We describe an approach for constructing samplers that are virtually automated for the class of spatial generalized linear models we consider in this paper. To achieve this, we first derive an approximation for posterior distributions arising from hierarchical models with latent Gaussian random fields. A heavy-tailed version of this approximation is then used to create a proposal distribution that allows for the construction of the samplers described in Section 2.

We derive the approximation by first obtaining an approximation for the joint distribution in a form that allows for a large number of parameters to be integrated out analytically. The method of integrating out model parameters to sample from a lower dimensional marginal distribution has been explored by several authors, including Blinded (2003); Everson (2001); Everson and Morris (2000); Gamerman et al. (2003); Wolfinger and Kass (2000). In each of the hierarchical models considered above, it is possible to analytically obtain the exact marginal distribution of the variance components. Furthermore, any MCMC algorithms used when sampling from the marginal distributions neither have known theoretical properties, nor rigorous estimates of error associated with them. For hierarchical models in general and the spatial models considered here in particular, exact formulae are not unavailable for any

marginal distributions.

4.1 A heavy-tailed approximation for hierarchical models

We now describe our general approach for deriving an approximation for a posterior distribution of a variance component model. Denote the precision parameters by Λ , the random effects by Θ and the data by \mathbf{Y} . Suppose the distribution of interest is $\pi(\Theta, \Lambda | \mathbf{Y})$ and we can derive an approximation $\hat{\pi}(\Theta, \Lambda | \mathbf{Y})$ for which we can analytically obtain the marginal distribution of the variance components, $s_1(\Lambda | \mathbf{Y})$, and the conditional distribution, $s_2(\Theta | \Lambda, \mathbf{Y})$. Approaches for deriving such a $\hat{\pi}$ involve using a Gaussian approximation to the likelihood, as described in Section 4.2; in general, the Laplace approximations may be useful for this as well (cf. Robert and Casella, 2005). Our general approach, first described in Blinded (2003), can be summarized as follows:

- Step 1: Approximate the target posterior, $\pi(\Theta, \Lambda | \mathbf{Y})$, as arising from a generalized linear model, by $\hat{\pi}(\Theta, \Lambda | \mathbf{Y})$.
- Step 2: Analytically integrate $\hat{\pi}(\Theta, \Lambda | \mathbf{Y})$ with respect to Θ to obtain the approximate marginal posterior $s_1(\Lambda | \mathbf{Y})$. $\hat{\pi}(\Theta, \Lambda | \mathbf{Y})$ can then be written as the product $s_1(\Lambda | \mathbf{Y})s_2(\Theta | \Lambda, \mathbf{Y})$ where s_2 is the easily obtained approximate conditional distribution of the random effects.
- Step 3: Find r_1 and r_2 such that they are heavier tailed distributions with similar shapes to s_1 and s_2 respectively. $r(\Theta, \Lambda | \mathbf{Y}) = r_1(\Lambda | \mathbf{Y})r_2(\Theta | \Lambda, \mathbf{Y})$ is then a heavy-tailed approximation to π . If we can demonstrate that r is heavy-tailed with respect to π , that is, it satisfies (1), it can be used to construct samplers with desirable properties.

We believe the approach outlined above can be used to derive useful approximations for several interesting and important Bayesian models. We focus our attention here on showing how it can be applied to the models described in Section 3.

4.2 A heavy-tailed approximation for a spatial generalized linear model

This subsection provides an outline of the derivation of the approximate marginal distributions (s_1) of the precision parameters for the GMRF model described in Section 3. Details are provided in Appendix A. A Gaussian approximation for the likelihood (2) is

$$Y_i \sim N(E_i e^{\mu_i}, E_i e^{\mu_i}) \quad (3)$$

Let $\hat{\mu}_i$ be $\log(Y_i/E_i)$. The delta method gives us the approximation

$$\hat{\mu}_i \sim N(\theta_i + \phi_i, 1/Y_i). \quad (4)$$

Note that if Y_i is 0, we replace it with 0.5 when constructing the approximation. A similar strategy was described in Haran et al. (2003) in the context of block MCMC sampling. If we denote $\boldsymbol{\theta} = (\theta_1, \dots, \theta_N)^T$, $\boldsymbol{\phi} = (\phi_1, \dots, \phi_N)^T$ and $\Theta = (\boldsymbol{\theta}^T, \boldsymbol{\phi}^T)^T$ we can analytically integrate the approximate joint posterior distribution $\hat{\pi}(\Theta, \tau_h, \tau_c | \mathbf{Y})$ with respect to Θ to obtain $s_1(\tau_h, \tau_c | \mathbf{Y})$, an approximation to the marginal distribution of (τ_h, τ_c) .

If we let $V^{-1} = \text{Diag}(Y_1, \dots, Y_N)$ and $\hat{\boldsymbol{\mu}} = (\hat{\mu}_1, \dots, \hat{\mu}_N)^T$, the approximate distribution of the random effects parameters, conditional on the precision parameters is

$$s_2(\Theta | \tau_h, \tau_c, \mathbf{Y}) \sim N \left(C^{-1} \left(-\frac{1}{2} D^T \right), C^{-1} \right), \quad (5)$$

where

$$C_{2N \times 2N} = \begin{bmatrix} V^{-1} + \tau_h I & +V^{-1} \\ +V^{-1} & V^{-1} + \tau_c Q \end{bmatrix},$$

and

$$D_{1 \times 2N} = (-2\hat{\boldsymbol{\mu}}^T V^{-1}, -2\hat{\boldsymbol{\mu}}^T V^{-1}).$$

We use heavy-tailed distributions r_1 and r_2 that have roughly the same shape and scale as s_1 and s_2 respectively. We let $r_1 = r_{1a} r_{1b}$, where r_{1a} and r_{1b} are independent log-t distributions. Among candidate distributions considered were the gamma, Weibull and log-normal, but the log-t was found to be the most appropriate due to its tail behavior and flexibility. r_2 was then obtained by using a multivariate-t with the mean and covariance of $s_2(\boldsymbol{\theta}, \boldsymbol{\phi} | \tau_h, \tau_c, \mathbf{Y})$. We can now state the following result.

Proposition 1. Let $r(\tau_h, \tau_c, \Theta) = r_1(\tau_h, \tau_c)r_2(\Theta | \tau_h, \tau_c)$. Let π be the posterior distribution corresponding to the Poisson-Markov random field model in Section 3. Then,

$$ess \sup \frac{\pi(\tau_h, \tau_c, \Theta)}{r(\tau_h, \tau_c, \Theta)} < B,$$

for some $B < \infty$, with $\tau_h > 0, \tau_c > 0, \Theta \in \mathbb{R}^{2n}$.

Proof. See Appendix B. □

To specify the parameter values for r_1 and r_2 , we find r_{1a} and r_{1b} that best match s_1 and s_2 . This can be done in a number of ways, but a simple and effective approach we used was as follows: The profiles of the log-transformed function $s_1(\tau_h, \tau_c | \mathbf{Y})$ along τ_h and τ_c were plotted. The heavy-tailed proposals for the precision parameters were found on the log-scale by matching the mode and variance of the t-distributions to the approximate log-scale profiles. The corresponding log-t distributions $r_{1a}(\tau_h)$ and $r_{1b}(\tau_c)$ were used jointly as the proposal $r_1(\tau_h, \tau_c | \mathbf{Y})$ for (τ_h, τ_c) . It is easy to draw $(\tau_h^*, \tau_c^*, \Theta^*)$ from $r_1 r_2$:

1. Draw $\tau_h^* \sim r_{1a}$ and $\tau_c^* \sim r_{1b}$.
2. Conditional on the values drawn above, simulate $\Theta^* \sim r_2(\Theta | \tau_h^*, \tau_c^*, \mathbf{Y})$.

4.3 Sampling algorithms

We now provide some details of the samplers constructed based on the approximations obtained in the previous sections. The ease of producing samples from r , combined with Proposition 1, allows us to construct two samplers:

1. Rejection sampler: r can be used as a proposal distribution in a classical rejection sampling algorithm.
2. Independence Metropolis-Hastings: r can be used as a proposal in an independence chain algorithm, where every proposed update to the Markov chain is obtained from r , regardless of the current value of the chain (Tierney, 1994). In particular, we consider an independence chain where the entire state space is updated in a single block. The starting value for the Markov chain is obtained by drawing a sample from r .

An important result corresponding to the above independence chain is obtained.

Proposition 2. *Consider an independence Metropolis-Hastings algorithm with target density $\pi(\tau_h, \tau_c, \Theta)$ and proposal density $r(\tau_h, \tau_c, \Theta)$. The resulting Markov chain is uniformly ergodic.*

Proof. Follows directly from Proposition 1 and Theorem 2.1 in Mengersen and Tweedie (1996). \square

An important consequence of Proposition 2 is that a Central Limit Theorem can be shown to hold for the independence Markov chain, and a consistent estimate of Monte Carlo standard errors is easily obtained via consistent batch means (Jones et al., 2006), as discussed in Section 2.1. Furthermore, the Markov chain can be stopped based on the estimated standard errors attaining a desired threshold (Flegal et al., 2008). Of course, the fact that the sampler is uniformly ergodic does not, on its own, imply that it is an efficient sampler in practice. We therefore study the efficiency of our independence Metropolis-Hastings algorithm in a variety of real data situations in Section 5.

We note that r also satisfies conditions for the construction of a perfect tempering sampler (Møller and Nicholls, 2009), an algorithm that utilizes simulated tempering to construct a variant of the Propp-Wilson perfect sampler (Propp and Wilson, 1996). This presents an intriguing possibility for future research since Møller and Nicholls (2009) report an increase in efficiency by using perfect tempering. However, as also seen in an in depth study of perfect tempering for such models in Blinded (2003), it is challenging to construct a perfect tempering algorithm that is consistently more efficient than the much simpler rejection sampler for the examples considered here.

5 Applications to Cancer and Infant Mortality Data

5.1 Description of data sets

We consider a total of three data sets: two on cancer in the U.S. state of Minnesota, and one on infant mortalities in the United States. The first two Minnesota cancer data sets were obtained from the Minnesota Cancer Surveillance System (MCSS), a cancer registry

maintained by the Minnesota Department of Health. The MCSS is population-based for the state of Minnesota, and collects information on geographic location and stage at detection for colorectal, prostate, lung, and female breast cancers. We illustrate our computational approaches by analyzing the MCSS data for two of the cancers, breast and colorectal. Each of the 87 counties in the data set has associated with it the total number of cancer cases recorded between 1995 and 1997, and the number of these detected late. We then take the expected number of late detections for that county as the number of cancer cases for that county multiplied by the statewide rate of late detections. The question of interest is whether there are clusters of counties in the state of Minnesota with much higher than expected late detection rates for either cancer. The spatial model provides smoothed estimates of the relative risk of cancer cases being detected late in each county. For these data, the posterior distribution based on the model in Section 3 has 176 dimensions.

We also consider a larger data set on infant mortalities in the United States. These data are derived from the Bureau of Health Professions Area Resource File, which is a county-level database for health analysis. Total number of births and deaths are obtained from 1998 to 2000. The geographic units used in this study include all counties of the following five states: Alabama, Georgia, Mississippi, North Carolina and South Carolina (Health Resources and Services Administration, 2003). The substantive problem of interest with this data is to determine spatial trends in infant mortality, and finding clusters of regions with unusually elevated levels in order to study possible socio-economic contributing factors. This data set, resulting in a posterior distribution of 910 dimensions, affords an opportunity to study the performance of our algorithms in a different and potentially more challenging high dimensional setting.

5.2 Setting up the algorithm

To find appropriate heavy-tailed distributions, it is useful to look at profiles of the approximate marginal posterior distributions (on log-scale) and find t-distributions that match them reasonably well. The log-t distributions are obtained by exponentiating these t-distributions. Once the log-t distributions are specified, the joint proposal distribution is obtained easily according to (5), and both rejection and independence Metropolis-Hastings samplers can be

constructed. For the rejection sampler, we find an upper bound B that satisfies (1) by numerically optimizing the ratio of the target and candidate densities. This worked well for our examples, though an alternative would be to use the “empirical-sup” rejection sampler (Caffo et al., 2002) to allow for the value of B to be adaptively estimated by the maximum based on the simulated candidates. We note that an advantage of the independence Metropolis-Hastings algorithm is the fact that this upper bound B is not needed.

The rejection sampler based on the heavy-tailed approximation produces a small number of accepted samples relative to the number of samples proposed. Similarly, the independence M-H algorithm may have low acceptance rates, especially for high dimensional distributions. It is therefore important to use efficient computational methods both for drawing the proposals, and for evaluating the Metropolis-Hastings ratios. The steps that take up most of the computational time involve operations on large matrices. Thus, for these samplers to be useful in practice it is necessary to exploit the sparseness of these matrices. We follow Rue (2001) by minimizing the bandwidth of the matrices involved and running fast band-matrix algorithms to speed up matrix computation. We find that this dramatically speeds up computation time and makes rejection sampling practical in some cases, and makes independence Metropolis-Hastings efficient and practical in all the examples we consider. The computer code was entirely implemented in R (Ihaka and Gentleman, 1996), utilizing appropriate sparse matrix routines written in FORTRAN from the well established online resource Netlib (<http://www.netlib.org>).

5.3 Results

The efficiency of the algorithms were compared both in terms of the number of samples required for the estimates to attain a desired level of accuracy, as well as the total time taken by the algorithm before stopping. The same desired level of accuracy was specified for each of the parameters in both algorithms. For example, for the breast cancer data set standard errors for each of the random effects was set to be no more than 0.01, while it was set to be no more than 2 for each of the precision parameters. These results are summarized in Table 1. Clearly, the independence M-H sampler is most time efficient for all three data sets. The efficiency of the independence M-H algorithm when compared to rejection sampling is

Table 1: Comparison of rejection sampler and independence MH (I-MH) algorithms

data set	samples required before stopping		total time taken (seconds)	
	rejection	I-MH	rejection	I-MH
breast cancer	4,118	29,241	2,663s	183s
colo-rectal cancer	4,735	27,225	543s	170s
infant mortality	—	97,721	—	10,066s

in accordance with results in Liu (1996). For the breast cancer and colo-rectal cancer data sets, exact sampling is viable. Note that while the number of samples required is similar for the breast cancer and colo-rectal cancer data sets, the time taken to produce samples is much longer for the breast cancer data. This is because fewer exact samples were returned per unit time for the breast cancer data set. For the infant mortality data set, the rejection sampler was unable to attain the desired accuracy level, even after the algorithm was run for well over 60 hours. The independence sampler, however, was able to provide estimates within 3 hours. Hence, the independence sampler may be practical even when the exact sampler is not. We feel it is important to note that the timings above assume that we do not utilize any form of parallel computing. Both the rejection sampler and the independence M-H sampler are ‘embarrassingly’ parallel in that each of the draws from the proposal as well as the target distribution and proposal evaluations can be done entirely in parallel. Hence, with relatively little effort the computational effort can be reduced by a factor corresponding to the number of available processors. Hence, for instance, a 3 hour independence M-H algorithm would take well under 5 minutes if 50 processors are available. This is of particular interest since computer clusters with large numbers of available processors are becoming increasingly accessible for scientific computing.

A motivation behind exploring rejection sampling is the simplicity and rigor of sample-based inference: easily computable consistent estimates of standard errors, avoidance of issues about determining starting values and not having to rely on ad-hoc approaches for determining Markov chain length. While the classic rejection sampler we have constructed based on our heavy-tailed approximation works surprisingly well in some cases, it is impractical

for challenging, high dimensional examples. We find our independence Metropolis-Hastings sampler is still efficient and feasible in such cases, while retaining much of the simplicity of exact sampling. In particular, to obtain starting values for our independence chain, we simply simulate from our heavy-tailed approximation r . From Proposition 2, this chain is uniformly ergodic, and we can easily see from results described in Section 2.1 that Monte Carlo standard errors for different parameters can be estimated consistently, and these standard errors can be used to provide a stopping rule for this sampler based on sound principles, just as one would for exact samplers. Hence, from the user’s perspective, sample-based inference is no more complicated than for inference based on the exact samplers. For the data sets and models considered here, this sampler resolves all MCMC issues originally raised in Section 1.

6 Discussion

We have demonstrated that it is feasible to implement samplers for realistic hierarchical Bayesian models in a manner that permits rigorous estimation of standard errors, while avoiding the usual issues regarding the determination of simulation lengths. We focused on hierarchical models where we are able to construct accurate heavy-tailed analytical approximations. We described a class of models to which we can apply our approximation techniques and derive theoretical results. Our general approximation strategy (in Section 4.1) is more broadly applicable, as are the central ideas behind automating the starting and stopping of the sampler in rigorous fashion. We do not claim that our approach to constructing the approximation is the best for any given problem since, for example, other approaches may improve upon the proposals constructed for the example outlined in Section 4. Better approximations will naturally lead to more efficient samplers. Rather, we believe the purpose of this paper is to show that using carefully derived heavy-tailed proposals along with recent developments in MCMC theory, sample-based inference can be carried out in a simple, fairly automated and rigorous fashion for some models.

For the examples in this paper we have explored exact sampling via the rejection sampler and have successfully demonstrated the applicability of our methods to some real data sets. These exact samplers, however, are generally not feasible for higher dimensional problems. A

fast mixing independence Metropolis-Hastings algorithm is a more efficient alternative, while retaining the rigor and simplicity of inference with exact samplers by guaranteeing that a Central Limit Theorem holds, and allowing for simple but consistent estimates of standard errors and intuitive and theoretically justified stopping rules. The increased efficiency of the independence chain becomes critical in higher dimensional problems, where it is viable even when the exact samplers are not. A potential weakness of our independence chain algorithm is that it involves block updates of high dimensions, which can lead to low acceptance rates as the dimensions of the problem increases. While this is certainly of concern for very large dimensional problems, we believe that for moderately high dimensional problems (thousands of dimensions), a combination of sparse matrix approaches and the latest high speed parallel computing, holds much promise for generating and evaluating high dimensional proposals extremely quickly. Such approaches are being explored elsewhere.

Using a combination of analytical work and modern computing power, our results suggest that it may be possible to construct rigorous, nearly automated approaches to sample-based inference for some classes of models that are of practical importance. In situations where multimodality may be an issue, the heavy-tailed approximation also provides a simple and reasonable approach for generating starting values for simulating multiple chains on several different processors, since the approximation is genuinely over-dispersed with respect to the target distribution. While obtaining theoretical results for each new model using our approach may prove to be non-trivial in general, we believe that the methodology outlined for constructing accurate heavy-tailed approximations may be generally useful, and our fixed width stopping rules may still be valuable in cases where analytical results are hard to obtain.

A Approximate marginal and conditional distributions

The full joint posterior distribution from Section 3 is:

$$\begin{aligned} \pi(\boldsymbol{\theta}, \boldsymbol{\phi}, \tau_h, \tau_c) \propto & \exp\left(\sum_{i=1}^N ((\theta_i + \phi_i)Y_i - E_i e^{\theta_i + \phi_i})\right) \exp\left(-\frac{1}{2}\boldsymbol{\theta}^T(\tau_h I)\boldsymbol{\theta}\right) \\ & \times \tau_h^{N/2 + \alpha_h - 1} \tau_c^{M/2 + \alpha_c - 1} \exp\left(-\frac{1}{2}\boldsymbol{\phi}^T(\tau_c Q)\boldsymbol{\phi}\right) \exp\left(-\frac{\tau_h}{\beta_h} - \frac{\tau_c}{\beta_c}\right). \end{aligned}$$

The approximate joint posterior distribution based on (3), (4) is:

$$\hat{\pi}(\Theta, \tau_h, \tau_c | \mathbf{Y}) \propto \exp \left(-\frac{1}{2}(\hat{\boldsymbol{\mu}} - (\boldsymbol{\theta} + \boldsymbol{\phi}))^T V^{-1}(\hat{\boldsymbol{\mu}} - (\boldsymbol{\theta} + \boldsymbol{\phi})) - \frac{1}{2}\boldsymbol{\theta}^T(\tau_h I)\boldsymbol{\theta} - \frac{1}{2}\boldsymbol{\phi}^T(\tau_c Q)\boldsymbol{\phi} \right) \\ \times \tau_h^{N/2+\alpha_h-1} \tau_c^{M/2+\alpha_c-1} \exp(-\tau_h/\beta_h - \tau_c/\beta_c)$$

where $\hat{\boldsymbol{\mu}} = (\log(Y_1/E_1), \dots, \log(Y_N/E_N))^T$, $M = N - 1$, $V^{-1} = \text{Diag}(Y_1, \dots, Y_N)$, I is an $N \times N$ identity matrix, Q is the adjacency matrix described in Section 3, and $\alpha_h, \alpha_c, \beta_h, \beta_c > 0$ are hyperparameters of the Gamma density (as in Section 4.2). To obtain the approximate marginal posterior distribution $s_1(\tau_h, \tau_c | \mathbf{Y})$ up to a constant of proportionality, we integrate $\hat{\pi}(\Theta, \tau_h, \tau_c | \mathbf{Y})$ with respect to Θ to obtain

$$s_1(\tau_h, \tau_c | \mathbf{Y}) = \tau_h^{N/2+\alpha_h-1} \tau_c^{M/2+\alpha_c-1} \exp(-\tau_h/\beta_h - \tau_c/\beta_c) \\ \times (\det(\tau_h V^{-1} + \tau_c V^{-1} Q + \tau_h \tau_c Q))^{-1/2} \\ \times \exp \left\{ -\frac{1}{2} \left(\hat{\boldsymbol{\theta}}^T C \hat{\boldsymbol{\theta}} + D \hat{\boldsymbol{\theta}} + k \right) \right\}$$

with

$$C_{2N \times 2N} = \begin{bmatrix} V^{-1} + \tau_h I & +V^{-1} \\ +V^{-1} & V^{-1} + \tau_c Q \end{bmatrix},$$

and

$$D_{1 \times 2N} = (-2\hat{\boldsymbol{\mu}}^T V^{-1}, -2\hat{\boldsymbol{\mu}}^T V^{-1}), \\ \hat{\boldsymbol{\theta}}^T = (\hat{\boldsymbol{\theta}}^T, \hat{\boldsymbol{\phi}}^T),$$

where $\hat{\boldsymbol{\theta}} = \frac{\tau_c}{\tau_h} Q (I + \frac{\tau_c}{\tau_h} Q + \tau_c V Q)^{-1} \hat{\boldsymbol{\mu}}$ and $\hat{\boldsymbol{\phi}} = (I + \frac{\tau_c}{\tau_h} Q + \tau_c V Q)^{-1} \hat{\boldsymbol{\mu}}$. For convenience, we have denoted C by $C(\tau_h, \tau_c)$. Our heavy-tailed approximation for the joint marginal of (τ_h, τ_c) is a product of log-t distributions,

$$r_1(\tau_h, \tau_c | \mathbf{Y}) \propto \frac{1}{\tau_h \tau_c} \left[1 + \frac{1}{\nu_h} \left(\frac{\log(\tau_h) - \mu_h}{\sigma_h} \right)^2 \right]^{-(\nu_h+1)/2} \left[1 + \frac{1}{\nu_c} \left(\frac{\log(\tau_c) - \mu_c}{\sigma_c} \right)^2 \right]^{-(\nu_c+1)/2}$$

where $\mu_h, \mu_c \in \mathbb{R}$, $\sigma_h, \sigma_c > 0$ and $\nu_h, \nu_c \in \mathbb{Z}^+$ are tuned to match the approximate marginal posterior density $s_1(\tau_h, \tau_c | \mathbf{Y})$. Our approximation for the conditional distribution of the random effects parameters, $\pi(\Theta | \mathbf{Y})$ is a multivariate-t version of (5):

$$r_2(\Theta | \tau_h, \tau_c, \mathbf{Y}) = \frac{|C|^{0.5} \Gamma(\frac{\nu_r+2N}{2})}{(\nu_r \pi)^{N/2} \Gamma(\nu_r/2)} \left(1 + \frac{1}{\nu_r} (\Theta - \mu_N)^T C (\Theta - \mu_N) \right)^{-(\nu_r+2N)/2}$$

where $\mu_N = -C^{-1}D^T/2$ and $\nu_r \in \mathbb{Z}^+$. Note that for technical reasons (see Lemma 2) Q in the matrix C is replaced by a positive definite \tilde{Q} matrix. Our heavy-tailed approximation for the joint density $\pi(\Theta, \tau_h, \tau_c|\mathbf{Y})$ is

$$r(\Theta, \tau_h, \tau_c|\mathbf{Y}) = r_1(\tau_h, \tau_c|\mathbf{Y})r_2(\Theta|\tau_h, \tau_c, \mathbf{Y}).$$

B Proofs

We begin with several lemmas that will be helpful for proving Proposition 1. We follow the notation used in Sections 3 and 4.

Lemma 1.

$$\frac{\tau_h^{N/2}\tau_c^{M/2}}{|C(\tau_h, \tau_c)|^{1/2}} \leq \left(\frac{\min \mathbf{Y} + \tau_h}{\min \mathbf{Y}}\right)^{1/2} \frac{1}{\tau_h^{1/2} \prod_{i=1}^{N-1} \lambda_i^{1/2}},$$

where $\lambda_1, \dots, \lambda_{N-1}$ are the non-zero eigen-values of Q .

Proof.

$$|C| = \left| \begin{bmatrix} V^{-1} + \tau_h I & V^{-1} \\ V^{-1} & V^{-1} + \tau_c Q \end{bmatrix} \right| = \left| \begin{bmatrix} V^{-1} + \tau_h I & V^{-1} \\ 0 & V^{-1} + \tau_c Q - V^{-1}(V^{-1} + \tau_h I)^{-1}V^{-1} \end{bmatrix} \right|, \quad (6)$$

by subtracting a matrix multiple of the first row block from the second. The determinant of the right hand side is the product of the determinants of the two diagonal blocks. The determinant of the first block is bounded by

$$|V^{-1} + \tau_h I| \geq \tau_h^N \quad (7)$$

The second diagonal block is

$$V^{-1} + \tau_c Q - V^{-1}(V^{-1} + \tau_h)^{-1}V^{-1} = \text{diag}(\mathbf{Y} - \mathbf{Y}^2/(\mathbf{Y} + \tau_h)) + \tau_c Q = \text{diag}\left(\frac{\mathbf{Y}\tau_h}{\mathbf{Y} + \tau_h}\right) + \tau_c Q$$

The determinant can be bounded below by replacing the Y_i by their minimum value ($\min \mathbf{Y}$).

We can then write $Q = U\Lambda U^T$ with $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_N)$. This gives the lower bound

$$|V^{-1} + \tau_c Q - V^{-1}(V^{-1} + \tau_h)^{-1}V^{-1}| \geq \prod_{i=1}^N \left(\frac{\min \mathbf{Y} \tau_h}{\min \mathbf{Y} + \tau_h} + \tau_c \lambda_i \right) \quad (8)$$

$$\geq \left(\frac{\min \mathbf{Y}}{\min \mathbf{Y} + \tau_h} \right) \tau_h \tau_c^M \prod_{i=1}^{N-1} \lambda_i. \quad (9)$$

Combining (6),(7),(8), we obtain the bound,

$$\frac{\tau_h^{N/2} \tau_c^{M/2}}{|C(\tau_h, \tau_c)|^{1/2}} \leq \left(\frac{\min \mathbf{Y} + \tau_h}{\min \mathbf{Y}} \right)^{1/2} \frac{1}{\tau_h^{1/2} \prod_{i=1}^{N-1} \lambda_i^{1/2}}$$

□

Lemma 2. *Let the center for the multivariate- t approximation, r_2 , be*

$$\mu_N = \begin{bmatrix} \mu_{N,\theta} \\ \mu_{N,\phi} \end{bmatrix} = 2 \begin{bmatrix} V^{-1} + \tau_h I & V^{-1} \\ V^{-1} & V^{-1} + \tau_c Q \end{bmatrix} \begin{bmatrix} V^{-1} \hat{\mu} \\ V^{-1} \hat{\mu} \end{bmatrix}.$$

Then, the quadratic terms $\|\tau_c \mu_{N,\phi}\|$ and $\|\tau_h \mu_{N,\theta}\|$ are bounded,

$$\|\tau_c \mu_{N,\phi}\| \leq 2 \|Q^{-1}\| \|\text{diag}(\mathbf{Y}^2)\| \|\hat{\mu}\|$$

$$\|\tau_h \mu_{N,\theta}\| \leq 2 \|\hat{\mu}\|.$$

Proof. We will use a non-singular version of the Q matrix. Let $\tilde{Q} = Q + \delta I$, for a small $\delta > 0$ in computing μ_N . This does not change the model, and is necessary only for the center, not the spread of the approximation. Using a partitioned form of the inverse, we can write

$$C^{-1} = \begin{bmatrix} A & -AV^{-1}(V^{-1} + \tau_c \tilde{Q})^{-1} \\ -BV^{-1}(V^{-1} + \tau_h I)^{-1} & B \end{bmatrix}$$

with

$$A = [V^{-1} + \tau_h I - V^{-1}(V^{-1} + \tau_c \tilde{Q})^{-1} V^{-1}]^{-1}$$

$$B = [\tau_c \tilde{Q} + V^{-1} - V^{-1}(V^{-1} + \tau_h I)^{-1} V^{-1}]^{-1}$$

Since $V^{-1} = \text{diag}(\mathbf{Y})$, B simplifies to

$$B = \left[\tau_c \tilde{Q} + \text{diag} \left(\mathbf{Y} - \frac{\mathbf{Y}^2}{\mathbf{Y} + \tau_h} \right) \right]^{-1} = \left[\tau_c \tilde{Q} + \text{diag} \left(\frac{\tau_h \mathbf{Y}}{\mathbf{Y} + \tau_h} \right) \right]^{-1}$$

and

$$C^{-1} = \begin{bmatrix} A & -AV^{-1}(V^{-1} + \tau_c \tilde{Q})^{-1} \\ -B \text{diag} \left(\frac{\mathbf{Y}}{\mathbf{Y} + \tau_h} \right) & B \end{bmatrix}$$

Now

$$\mu_{N,\phi} = 2B \left(I - \text{diag} \left(\frac{\mathbf{Y}}{\mathbf{Y} + \tau_h} \right) \right) \text{diag}(\mathbf{Y}) \hat{\mu} = B \text{diag} \left(\frac{\tau_h \mathbf{Y}^2}{\mathbf{Y} + \tau_h} \right) \hat{\mu}$$

To see that $\|\mu_{N,\phi}\|$ is bounded, note that

$$B^{-1} \geq \text{diag} \left(\frac{\tau_h \mathbf{Y}}{\mathbf{Y} + \tau_h} \right)$$

in non-negative definite matrix order and therefore

$$\|\mu_{N,\phi}\| \leq 2 \left\| \text{diag} \left(\frac{\tau_h \mathbf{Y}}{\mathbf{Y} + \tau_h} \right)^{-1} \text{diag} \left(\frac{\tau_h \mathbf{Y}^2}{\mathbf{Y} + \tau_h} \right) \hat{\mu} \right\| = 2 \|\text{diag}(\mathbf{Y}) \hat{\mu}\|$$

To bound $\|\tau_c \mu_{N,\phi}\|$, we can similarly use

$$B^{-1} \geq \tau_c \tilde{Q}$$

and

$$\text{diag} \left(\frac{\tau_h \mathbf{Y}^2}{\mathbf{Y} + \tau_h} \right) \leq \text{diag}(\mathbf{Y}^2)$$

to conclude that

$$\|\tau_c \mu_{N,\phi}\| \leq 2 \|\tilde{Q}^{-1}\| \|\text{diag}(\mathbf{Y}^2)\| \|\hat{\mu}\|.$$

The component $\mu_{N,\theta}$ is given by

$$\mu_{N,\theta} = 2A[I - V^{-1}(V^{-1} + \tau_c Q)^{-1}]V^{-1}\hat{\mu} = 2A[V^{-1} - V^{-1}(V^{-1} + \tau_c Q)^{-1}V^{-1}]\hat{\mu}$$

To bound $\|\mu_{N,\theta}\|$ and $\|\tau_h \mu_{N,\theta}\|$ we can use the inequalities

$$A \leq [V^{-1}V^{-1}(V^{-1} + \tau_c Q)^{-1}V^{-1}]^{-1}$$

and

$$A \leq \frac{1}{\tau_c} I$$

Using the first inequality

$$\|\mu_{N,\theta}\| \leq 2 \|[V^{-1}V^{-1}(V^{-1} + \tau_c Q)^{-1}V^{-1}]^{-1}[V^{-1} - V^{-1}(V^{-1} + \tau_c Q)^{-1}V^{-1}]\hat{\mu}\| = 2\|\hat{\mu}\|$$

and using the second inequality

$$\|\tau_h \mu_{N,\theta}\| \leq 2\tau_h \left\| \frac{1}{\tau_h} \hat{\mu} \right\| = 2\|\hat{\mu}\|.$$

□

Lemma 3.

$$\exp\left(\sum_{i=1}^N((\theta_i + \phi_i)Y_i - E_i e^{\theta_i + \phi_i} - \frac{\tau_h}{2}\theta^T\theta - \frac{\tau_c}{2}\phi^T Q\phi)\right) \left(1 + \frac{1}{\nu_r}(\Theta - \mu_N)C(\Theta - \mu_N)\right)^{(\nu_r + 2N)/2}$$

is bounded.

Proof. First note that for $A, B \geq 0$ we have

$$1 + A + B = (1 + A) \left(1 + \frac{B}{1 + A}\right) \leq (1 + A)(1 + B)$$

It is therefore sufficient to bound each of the three terms

$$\exp\left(\sum_{i=1}^N((\theta_i + \phi_i)Y_i - E_i e^{\theta_i + \phi_i})\right) \left(1 + \frac{1}{\nu_r}(\Theta - \mu_N)C^*(\Theta - \mu_N)\right)^{(\nu_r + 2N)/2} \quad (10)$$

$$\exp\left(-\frac{\tau_h}{2}\theta^T\theta\right) \left(1 + \frac{1}{\nu_r}(\Theta - \mu_N)C^{**}(\Theta - \mu_N)\right)^{(\nu_r + 2N)/2} \quad (11)$$

$$\exp\left(-\frac{\tau_c}{2}\phi^T Q\phi\right) \left(1 + \frac{1}{\nu_r}(\Theta - \mu_N)C^{***}(\Theta - \mu_N)\right)^{(\nu_r + 2N)/2} \quad (12)$$

where

$$C^* = \begin{bmatrix} V^{-1} & V^{-1} \\ V^{-1} & V^{-1} \end{bmatrix} \quad C^{**} = \begin{bmatrix} \tau_h I & 0 \\ 0 & 0 \end{bmatrix} \quad C^{***} = \begin{bmatrix} 0 & 0 \\ 0 & \tau_c Q \end{bmatrix}$$

To bound (10), reparameterize in terms of $w_i = \theta_i + \phi_i$ and $v_i = \theta_i - \phi_i$. Then the quadratic form in (10) can be written as

$$\begin{aligned} (\Theta - \mu_N)^T C^* (\Theta - \mu_N) &= [(w - \mu_w^*)^T \quad (v - \mu_v^*)^T] \begin{bmatrix} \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} \end{bmatrix} \begin{bmatrix} V^{-1} & V^{-1} \\ V^{-1} & V^{-1} \end{bmatrix} \begin{bmatrix} \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} \end{bmatrix} \begin{bmatrix} w - \mu_w^* \\ v - \mu_v^* \end{bmatrix} \\ &= (w - \mu_w^*)^T V^{-1} (w - \mu_w^*) \\ &= \sum_{i=1}^N Y_i (w_i - \mu_i^*)^2 \end{aligned}$$

with

$$\mu^* = \begin{bmatrix} \mu_w^* \\ \mu_v^* \end{bmatrix} = \begin{bmatrix} \frac{1}{2}I & \frac{1}{2}I \\ \frac{1}{2}I & -\frac{1}{2}I \end{bmatrix} \mu_N$$

This does not depend on v . So the quadratic form in (10) is bounded by

$$(\Theta - \mu_N)^T C^* (\Theta - \mu_N) \leq K_1 + K_2 \sum_{i=1}^N Y_i w_i^2$$

for some constants K_1 and K_2 . The term (10) can thus be bounded by

$$(1 + K_1)^{(\nu_r+2N)/2} \prod_{i=1}^N \left[\exp(w_i Y_i - E_i e^{w_i}) \left(1 + \frac{K_2}{\nu_r} Y_i w_i^2 \right)^{(\nu_r+2N)/2} \right]$$

None of these terms involves τ_h or τ_c and since $Y_i > 0$ each term

$$\exp(w_i Y_i - E_i e^{w_i}) \left(1 + \frac{K_2}{\nu_r} Y_i w_i^2 \right)^{(\nu_r+2N)/2}$$

is bounded in w_i . So (10) is bounded. Using Lemma 2, the quadratic forms in (11) and (12) are bounded by

$$(\Theta - \mu_N)^T C^{**} (\Theta - \mu_N) = \tau_h (\theta - \mu_{N,\theta})^T (\theta - \mu_{N,\theta}) \leq K_3 + K_4 \tau_h \theta^T \theta$$

and

$$(\Theta - \mu_N)^T C^{***} (\Theta - \mu_N) = \tau_c (\phi - \mu_{N,\phi})^T Q (\phi - \mu_{N,\phi}) \leq K_5 + K_6 \tau_c \phi^T Q \phi$$

for some constants K_3, K_4, K_5 , and K_6 . Thus term (11) is bounded by

$$(1 + K_3)^{(\nu_r+2N)/2} \sup_{z \geq 0} \left\{ \exp\left(-\frac{z}{2}\right) \left(1 + \frac{K_4}{\nu_r} z \right)^{(\nu_r+2N)/2} \right\} < \infty$$

and term (12) is bounded analogously. \square

We can now utilize the above lemmas to prove Proposition 1.

Proposition 1. The ratio $\pi(\tau_h, \tau_c, \Theta)/r(\tau_h, \tau_c, \Theta)$ can be written as

$$\begin{aligned} & \exp\left(\sum_{i=1}^N ((\theta_i + \phi_i) Y_i - E_i e^{\theta_i + \phi_i} - \frac{\tau_h}{2} \theta^T \theta - \frac{\tau_c}{2} \phi^T Q \phi)\right) \left(1 + \frac{1}{\nu_r} (\Theta - \mu_N)^T C (\Theta - \mu_N) \right)^{(\nu_r+2N)/2} \\ & \times \exp\left(-\frac{\tau_h}{\beta_h} - \frac{\tau_c}{\beta_c}\right) \left[1 + \frac{1}{\nu_h} \left(\frac{\log(\tau_h) - \mu_h}{\sigma_h} \right)^2 \right]^{(\nu_h+1)/2} \left[1 + \frac{1}{\nu_c} \left(\frac{\log(\tau_c) - \mu_c}{\sigma_c} \right)^2 \right]^{(\nu_c+1)/2} \\ & \times \tau_h^{N/2 + \alpha_h} \tau_c^{M/2 + \alpha_c} \det(C)^{-0.5}. \end{aligned}$$

The first term is bounded by Lemma 3. We can apply Lemma 1 to bound the determinant in the third term. Hence, it follows that the above ratio (ignoring constants) is bounded by

$$\begin{aligned} & \exp\left(-\frac{\tau_h}{\beta_h} - \frac{\tau_c}{\beta_c}\right) \left[1 + \frac{1}{\nu_h} \left(\frac{\log(\tau_h) - \mu_h}{\sigma_h} \right)^2 \right]^{(\nu_h+1)/2} \left[1 + \frac{1}{\nu_c} \left(\frac{\log(\tau_c) - \mu_c}{\sigma_c} \right)^2 \right]^{(\nu_c+1)/2} \\ & \times \tau_h^{\alpha_h} \tau_c^{\alpha_c} \left(\frac{\min \mathbf{Y} + \tau_h}{\min \mathbf{Y}} \right)^{1/2} \frac{1}{\tau_h^{1/2}}, \end{aligned}$$

which is bounded as long as $\alpha_h \geq 1$. \square

References

- Banerjee, S., Carlin, B., and Gelfand, A. (2004). *Hierarchical modeling and analysis for spatial data*. Chapman & Hall Ltd.
- Besag, J., York, J., and Mollie, A. (1991). Bayesian image restoration, with two applications in spatial statistics (Disc: p21-59). *Annals of the Institute of Statistical Mathematics*, 43:1–20.
- Billingsley, P. (1999). *Convergence of probability measures*. John Wiley and Sons.
- Blinded, B. (2003). Blinded Ph.D. thesis. Technical report, Blinded University, Blinded.
- Caffo, B. S., Booth, J. G., and Davison, A. C. (2002). Empirical supremum rejection sampling. *Biometrika*, 89(4):745–754.
- Cowles, M. K. and Carlin, B. P. (1996). Markov chain Monte Carlo convergence diagnostics: A comparative review. *Journal of the American Statistical Association*, 91:883–904.
- Diggle, P. J., Tawn, J. A., and Moyeed, R. A. (1998). Model-based geostatistics (Disc: p326-350). *Applied Statistics*, 47:299–326.
- Everson, P. J. (2001). Exact Bayesian inference for normal hierarchical models. *Journal of Statistical Computation and Simulation*, 68(3):223–241.
- Everson, P. J. and Morris, C. N. (2000). Inference for multivariate normal hierarchal models. *Journal of the Royal Statistical Society, Series B, Methodological*, 62(2):399–412.
- Flegal, J., Haran, M., and Jones, G. (2008). Markov chain Monte Carlo: Can we trust the third significant figure? *Statistical Science*, 23:250–260.
- Gamerman, D., Moreira, A. R., and Rue, H. (2003). Monte Carlo EM with importance reweighting and its applications in random effects models. *Computational Statistics and Data Analysis*, 42:513–533.

- Haran, M., Hodges, J. S., and Carlin, B. P. (2003). Accelerating computation in Markov random field models for spatial data via structured MCMC. *Journal of Computational and Graphical Statistics*, 12:249–264.
- Health Resources and Services Administration (2003). Health professions, area resource file (arf) system. Technical report, Quality Resource Systems, Inc., Fairfax, VA.
- Hobert, J. P. and Geyer, C. J. (1998). Geometric ergodicity of Gibbs and block Gibbs samplers for a hierarchical random effects model. *Journal of Multivariate Analysis*, 67:414–430.
- Ihaka, R. and Gentleman, R. (1996). R: A language for data analysis and graphics. *Journal of Computational and Graphical Statistics*, 5:299–314.
- Jones, G. L. (2004). On the Markov chain central limit theorem. *Probability Surveys*, 1:299–320.
- Jones, G. L., Haran, M., Caffo, B. S., and Neath, R. (2006). Fixed-width output analysis for Markov Chain Monte Carlo. *Journal of the American Statistical Association*, 101(476):1537–1547.
- Knorr-Held, L. and Rue, H. (2002). On block updating in Markov random field models for disease mapping. *Scandinavian Journal of Statistics*, 29(4):597–614.
- Liu, J. S. (1996). Metropolized independent sampling with comparisons to rejection sampling and importance sampling. *Statistics and Computing*, 6:113–119.
- McCullagh, P. and Nelder, J. A. (1999). *Generalized Linear Models*. Chapman & Hall Ltd.
- Mengersen, K. L. and Tweedie, R. L. (1996). Rates of convergence of the Hastings and Metropolis algorithms. *The Annals of Statistics*, 24:101–121.
- Meyn, S. P. and Tweedie, R. L. (1993). *Markov chains and stochastic stability*. Springer-Verlag Inc.
- Møller, J. and Nicholls, G. (2009). Perfect simulation for sample-based inference. Technical report, University of Aarhus, Department of Mathematical Sciences.

- Mollié, A. (1996). Bayesian mapping of disease. In *Markov chain Monte Carlo in practice* Eds. Gilks, W.R., Richardson, S. and Spiegelhalter, D.J., pages 359–379. London: Chapman and Hall.
- Propp, J. G. and Wilson, D. B. (1996). Exact sampling with coupled Markov chains and applications to statistical mechanics. *Random Structures and Algorithms*, 9(1-2):223–252.
- Robert, C. P. and Casella, G. (2005). *Monte Carlo Statistical Methods, Second Edition*. Springer, New York.
- Roberts, G. O. and Rosenthal, J. S. (2004). General state space Markov chains and MCMC algorithms. *Probability Surveys*, 1:20–71.
- Rue, H. (2001). Fast sampling of Gaussian Markov random fields. *Journal of the Royal Statistical Society, Series B, Methodological*, 63(2):325–338.
- Rue, H. and Held, L. (2005). *Gaussian Markov Random Fields: Theory and Applications*. Chapman and Hall.
- 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, Methodological*, 71(2):1–35.
- Tierney, L. (1994). Markov chains for exploring posterior distributions (Disc: P1728-1762). *The Annals of Statistics*, 22:1701–1728.
- Tierney, L. and Kadane, J. B. (1986). Accurate approximations for posterior moments and marginal densities. *Journal of the American Statistical Association*, 81:82–86.
- von Neumann, J. (1951). Various techniques used in connection with random digits. *Applied Math Series*, 12:36–38.
- Wolfinger, R. D. and Kass, R. E. (2000). Nonconjugate Bayesian analysis of variance component models. *Biometrics*, 56(3):768–774.