



## Bayesian Inference in the Presence of Intractable Normalizing Functions

Jaewoo Park & Murali Haran

To cite this article: Jaewoo Park & Murali Haran (2018) Bayesian Inference in the Presence of Intractable Normalizing Functions, Journal of the American Statistical Association, 113:523, 1372-1390, DOI: [10.1080/01621459.2018.1448824](https://doi.org/10.1080/01621459.2018.1448824)

To link to this article: <https://doi.org/10.1080/01621459.2018.1448824>



View supplementary material [↗](#)



Accepted author version posted online: 14 Mar 2018.  
Published online: 14 Mar 2018.



Submit your article to this journal [↗](#)



Article views: 322



View Crossmark data [↗](#)



# Bayesian Inference in the Presence of Intractable Normalizing Functions

Jaewoo Park and Murali Haran

Department of Statistics, Pennsylvania State University, Pennsylvania, PA

## ABSTRACT

Models with intractable normalizing functions arise frequently in statistics. Common examples of such models include exponential random graph models for social networks and Markov point processes for ecology and disease modeling. Inference for these models is complicated because the normalizing functions of their probability distributions include the parameters of interest. In Bayesian analysis, they result in so-called doubly intractable posterior distributions which pose significant computational challenges. Several Monte Carlo methods have emerged in recent years to address Bayesian inference for such models. We provide a framework for understanding the algorithms, and elucidate connections among them. Through multiple simulated and real data examples, we compare and contrast the computational and statistical efficiency of these algorithms and discuss their theoretical bases. Our study provides practical recommendations for practitioners along with directions for future research for Markov chain Monte Carlo (MCMC) methodologists. Supplementary materials for this article are available online.

## ARTICLE HISTORY

Received January 2017  
Revised September 2017

## KEYWORDS

Doubly intractable distributions; Exponential random graph models; Importance sampling; Markov chain Monte Carlo; Markov point processes

## 1. Introduction

Markov chain Monte Carlo (MCMC) has been used to routinely carry out Bayesian inference for an enormous variety of complicated models (see Brooks et al. 2011). However, inference for Bayesian models with intractable normalizing functions, where the normalizing constant of the model is itself a function of parameters of interest, is still far from routine. There are many well-known models that have intractable functions, for instance the class of exponential family random graph models and its variants (see Robins et al. 2007; Hunter and Handcock 2012) which are popular models for social networks. Non-Gaussian Markov random field models in spatial statistics (see Besag 1974; Hughes, Haran, and Caragea 2011, for a review) also have intractable normalizing functions. In fact, it is worth noting that the Ising model (Lenz 1920; Ising 1925), which is a non-Gaussian Markov random field, appears in the landmark article on the Metropolis algorithm (Metropolis et al. 1953). While the Metropolis algorithm provides an elegant way to simulate from this model for a given parameter value, the article did not consider the more difficult problem of performing inference for this model.

Consider  $h(\mathbf{x}|\theta)$ , an unnormalized probability model for a random variable  $\mathbf{x} \in \mathcal{X}$  given a parameter vector  $\theta \in \Theta$ , with a normalizing function  $Z(\theta) = \int_{\mathcal{X}} h(\mathbf{x}|\theta) d\mathbf{x}$ . Let  $p(\theta)$  be the prior density for  $\theta$ . The likelihood function,  $L(\theta|\mathbf{x})$  is  $h(\mathbf{x}|\theta)/Z(\theta)$  and the posterior density of  $\theta$  is

$$\pi(\theta|\mathbf{x}) \propto p(\theta) \frac{h(\mathbf{x}|\theta)}{Z(\theta)}. \quad (1)$$

The problem in constructing an MCMC algorithm stems from the fact that  $Z(\theta)$  cannot be easily evaluated and the acceptance probability at each step of the Metropolis–Hastings (MH) algorithm (Metropolis et al. 1953; Hastings 1970) involves evaluating  $Z(\theta)$  both at the current and proposed value of  $\theta$ .

There have been several proposals to address the intractable normalizing function problem in a maximum likelihood context. Besag (1974) proposed the maximum pseudolikelihood estimate (MPLE) which maximizes the pseudolikelihood, a particular likelihood approximation that does not require  $Z(\theta)$ . MPLE does not, in general, define a likelihood function. Furthermore, MPLE may be a reasonable estimator when there is weak dependence among data points relative to the size of the dataset but in other cases its performance is unsatisfactory. Younes (1988) proposed a stochastic gradient algorithm to solve normal equations to find the maximum likelihood estimate (MLE) for models with intractable normalizing functions. However, the step size and starting point must be selected carefully, otherwise the algorithm becomes slow and does not converge (Ibáñez and Simó 2003). Geyer and Thompson (1992) proposed MCMC-MLE which is based on maximizing a Monte Carlo approximation to the likelihood; this approximation is based on an importance sampling approximation of  $Z(\theta)$ . This is an elegant and theoretically justified algorithm. In practice, MCMC-MLE suffers from some of the usual challenges faced by importance sampling approaches, namely, that for an accurate approximation to MLE, the initial value for the algorithm should be reasonably close to the MLE. Some of these issues may be partially addressed by umbrella sampling (Torrìe and Valleau 1977; Geyer 2011), which involves sampling from mixtures of

importance sampling distributions. However approximating standard errors can also be difficult in situations where analytical gradients for the unnormalized likelihood are unavailable and using bootstrap techniques may be computationally infeasible (see Goldstein et al. 2015). In such cases and also in situations where there is an interest in incorporating prior information about the parameters or avoiding model degeneracies (see Handcock 2003), Bayesian alternatives may be preferable. In this manuscript, we focus on Bayesian inference; we refer readers to Geyer (2011) for a general overview and Hunter, Krivitsky, and Schweinberger (2012) for a review of recent MCMC-MLE methods for exponential random graph models.

Several MCMC algorithms have recently been proposed for Bayesian inference in the presence of intractable normalizing functions. These algorithms may be broadly classified into two general if somewhat overlapping categories: (1) algorithms where the introduction of a well-chosen auxiliary variable results in the normalizing function (or a ratio of normalizing functions) canceling out in the Metropolis–Hastings acceptance probability, and (2) directly approximating the normalizing function (or a ratio of normalizing functions), and substituting the approximation into the acceptance probability. Here, we will refer to the first as an *auxiliary variable approach* and the second as a *likelihood approximation approach* (see also a discussion in Liang et al. 2016). In addition, MCMC algorithms may also be classified as “asymptotically exact” or “asymptotically inexact.” For asymptotically exact algorithms, the Markov chain’s stationary distribution is exactly equal to the desired posterior distribution. On the other hand, asymptotically inexact or “noisy” algorithms generate Markov chains without this property; even asymptotically the samples generated only follow the target distribution approximately.

In what follows we discuss several MCMC algorithms. We provide an explanation of the ideas underpinning each algorithm along with figures that summarize them and make it easier to see how they are related. We also discuss theoretical justifications and practical implementation issues. We carry out a comparative study by using three different examples: an Ising model, a social network model, and a spatial point process. We provide some guidance about potential advantages and disadvantages of each algorithm along with connections among them, providing some future avenues for research. The remainder of this article is organized as follows. In Section 2, we discuss several auxiliary variable algorithms. In Section 3, we cover several likelihood approximation algorithms. In Section 4, we describe the application of the algorithms in the context of three different case studies and provide some insights based on our results. In particular, we discuss in detail the computational complexity of the algorithms in Section 5. We point out connections between the algorithms in Section 6 along with general guidelines and recommendations based on our study. We conclude with a summary in Section 7.

## 2. Auxiliary Variable Approaches

In this section, we review several auxiliary variable approaches. Here, the target distribution includes both the parameter of interest as well as an auxiliary variable. By a clever choice of the auxiliary variable proposal the intractable functions get cancelled in the acceptance probability of the Metropolis–Hastings algorithm. What distinguishes the different auxiliary variable algorithms from each other is how the auxiliary variable is sampled.

### 2.1. Auxiliary Variable MCMC

Møller et al. (2006) introduced an auxiliary variable  $\mathbf{y}$  with the conditional density  $f(\mathbf{y}|\theta, \mathbf{x})$  so that the intractable terms are cancelled in the acceptance probability of the Metropolis–Hastings algorithm. Suppose the original target density is  $\pi(\theta|\mathbf{x}) \propto p(\theta)h(\mathbf{x}|\theta)/Z(\theta)$ . Then the augmented target density is  $\pi(\theta, \mathbf{y}|\mathbf{x}) \propto f(\mathbf{y}|\theta, \mathbf{x})p(\theta)h(\mathbf{x}|\theta)/Z(\theta)$  whose marginal density becomes  $\int_{\mathcal{X}} \pi(\theta, \mathbf{y}|\mathbf{x})d\mathbf{y} = \pi(\theta|\mathbf{x})$ , the original target density. Now consider a joint proposal density for  $\{\theta, \mathbf{y}\}$  updates which can be factorized as  $q(\theta', \mathbf{y}'|\theta, \mathbf{y}) = q(\mathbf{y}'|\theta')q(\theta'|\theta)$ . Møller et al. (2006) took the proposal for the auxiliary variable  $q(\mathbf{y}'|\theta')$  as  $h(\mathbf{y}'|\theta')/Z(\theta')$  which is the same as the model for the data  $\mathbf{x}$  given the parameter value  $\theta'$  proposed from  $q(\theta'|\theta)$ . The resulting algorithm for the above joint density with this proposal density has Metropolis–Hastings acceptance probability

$$\alpha = \min \left\{ 1, \frac{f(\mathbf{y}'|\theta', \mathbf{x})p(\theta')h(\mathbf{x}|\theta')Z(\theta)h(\mathbf{y}|\theta)Z(\theta')q(\theta|\theta')}{f(\mathbf{y}|\theta, \mathbf{x})p(\theta)h(\mathbf{x}|\theta)Z(\theta')h(\mathbf{y}'|\theta')Z(\theta)q(\theta'|\theta)} \right\}. \quad (2)$$

The resulting Metropolis–Hastings acceptance probability (2) does not contain the normalizing functions because they get cancelled out. Since  $\int_{\mathcal{X}} \pi(\theta, \mathbf{y}|\mathbf{x})d\mathbf{y} = \pi(\theta|\mathbf{x})$ , the marginal  $\theta$  samples follow the original target distribution. We will henceforth use AVM to refer to this auxiliary variable MCMC algorithm. We summarize the algorithm in Figure 1. We should note that updating  $\mathbf{y}$  requires drawing a sample with exact distribution  $h(\cdot|\theta')/Z(\theta')$ . For some models this can be achieved via perfect sampling (Propp and Wilson 1996), a clever method that uses bounding Markov chains to construct a sampler where the draws are exactly (not just asymptotically) from the target distribution. Although perfect sampling is possible for some models, for instance certain Markov random field (MRF) models, it is not trivial to construct a perfect sampler in general. This is a major practical limitation of the AVM.

*Components to be tuned:* Møller et al. (2006) reported that the choice of  $f(\mathbf{y}|\theta, \mathbf{x})$  impacts the mixing of the chain. Suppose  $Z(\theta)$  is known so that we can choose  $f(\mathbf{y}|\theta, \mathbf{x}) = h(\mathbf{y}|\theta)/Z(\theta)$ . Then it is easily seen that (2) is identical to the acceptance probability of the Metropolis–Hastings algorithm with the stationary density  $\pi(\theta|\mathbf{x})$ , which implies that this chain has the same

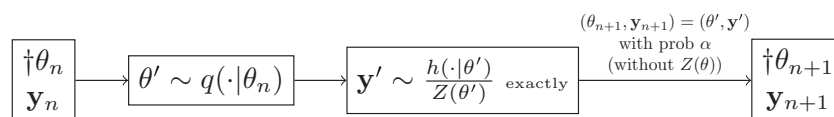


Figure 1. Illustration of the  $(n + 1)$  step update for auxiliary variable MCMC. The dagger symbols indicate the parameter of interest.

**Algorithm 1** Auxiliary variable MCMC (AVM)

Given  $\{\theta_n, \mathbf{y}_n\} \in \Theta \times \mathcal{X}$  at  $n$ th iteration.

1. Propose  $\theta' \sim q(\cdot|\theta_n)$ .
2. Propose the auxiliary variable exactly from probability model at  $\theta'$ :  $\mathbf{y}' \sim \frac{h(\cdot|\theta')}{Z(\theta')}$  using perfect sampling.
3. Accept  $\{\theta_{n+1}, \mathbf{y}_{n+1}\} = \{\theta', \mathbf{y}'\}$  with probability  $\alpha = \min \left\{ 1, \frac{f(\mathbf{y}'|\theta', \mathbf{x})p(\theta')h(\mathbf{x}|\theta')h(\mathbf{y}_n|\theta_n)q(\theta_n|\theta')}{f(\mathbf{y}_n|\theta_n, \mathbf{x})p(\theta_n)h(\mathbf{x}|\theta_n)h(\mathbf{y}'|\theta')q(\theta'|\theta_n)} \right\}$ , else reject (set  $\{\theta_{n+1}, \mathbf{y}_{n+1}\} = \{\theta_n, \mathbf{y}_n\}$ ).

convergence properties as the Markov chain where the normalizing function is known and the same proposal  $q(\theta'|\theta)$  is used for  $\theta$ . Therefore, a reasonable choice of conditional density for the auxiliary variable is  $f(\mathbf{y}|\theta, \mathbf{x})$  that approximates  $h(\mathbf{y}|\theta)/Z(\theta)$ . A simple choice is  $f(\mathbf{y}|\theta, \mathbf{x}) = h(\mathbf{y}|\hat{\theta})/Z(\hat{\theta})$ , where  $\hat{\theta}$  may be an approximation to the MLE.  $\hat{\theta}$  should be predetermined before implementing the AVM algorithm. For example, the maximum pseudolikelihood estimate (MPLE) proposed by Besag (1974) may be an option though for some problems the MPLE may be a poor approximation to the MLE. Then  $f(\mathbf{y}'|\theta', \mathbf{x})/f(\mathbf{y}|\theta, \mathbf{x}) = h(\mathbf{y}'|\hat{\theta})/h(\mathbf{y}|\hat{\theta})$ , which makes it possible to calculate (2) because there are no intractable terms. Another possible choice of  $f(\mathbf{y}|\theta, \mathbf{x})$  is using a normalizable density without  $Z(\theta)$ . AVM mixes better as  $f(\mathbf{y}|\theta, \mathbf{x})$  more closely resembles  $h(\mathbf{y}|\theta)/Z(\theta)$ .

*Theoretical justification:* The Markov chain satisfies the detailed balance condition if  $\pi(\theta|\mathbf{x})T(\theta'|\theta) = \pi(\theta'|\mathbf{x})T(\theta|\theta')$ , where  $T(\theta'|\theta)$  is the transition kernel of the Markov chain with the target density  $\pi(\theta|\mathbf{x})$ . Since the Markov chain satisfies detailed balance with respect to the augmented target distribution, the marginal distribution of which is the posterior distribution of  $\theta$ , AVM is an asymptotically exact MCMC algorithm. However, to achieve the detailed balance condition, we need to sample  $\mathbf{y}$  exactly from the likelihood function  $h(\mathbf{y}|\theta')/Z(\theta')$ .

**2.2. The Exchange Algorithm**

Appearing almost simultaneously with Møller et al. (2006), the exchange algorithm (Murray, Ghahramani, and MacKay 2006) also constructs an augmented target distribution and updates the augmented state via the Metropolis–Hastings algorithm. Consider the auxiliary variable  $\mathbf{y}$  which follows  $h(\mathbf{y}|\theta')/Z(\theta')$  and conditional density of  $\theta'$  for given  $\theta$  as  $q(\theta'|\theta)$  where  $\theta$  is the parameter setting of data  $\mathbf{x}$ . Then the augmented joint density is

$$\begin{aligned} \pi(\theta, \theta', \mathbf{y}|\mathbf{x}) &\propto p(\theta)L(\theta|\mathbf{x})q(\theta'|\theta)L(\theta'|\mathbf{y}) \\ &= p(\theta)\frac{h(\mathbf{x}|\theta)}{Z(\theta)}q(\theta'|\theta)\frac{h(\mathbf{y}|\theta')}{Z(\theta')}. \end{aligned} \quad (3)$$

For this augmented density,  $\{\theta', \mathbf{y}\}$  is updated through block-Gibbs samplers;  $\theta'$  is generated from the proposal  $q(\cdot|\theta)$ ; and

the auxiliary variable  $\mathbf{y}$  is generated from  $h(\cdot|\theta')/Z(\theta')$ . The first two arrows in Figure 2 correspond to the update of  $\{\theta', \mathbf{y}\}$ . Then  $\theta$  is updated through exchanging parameter settings. Let  $\{\theta, \theta'\}$  be the current parameter settings for  $\{\mathbf{x}, \mathbf{y}\}$ . Consider a swapping proposal  $s(\{\theta^*, \theta'^*\}|\{\theta, \theta'\}) = \delta(\theta^* - \theta')\delta(\theta'^* - \theta)$ , where  $\delta$  denotes the Dirac delta function. After swapping is proposed, data  $\mathbf{x}$  follow  $\theta'$  instead of  $\theta$  and the auxiliary variable  $\mathbf{y}$  follows  $\theta$  instead of  $\theta'$ . The symmetric swapping proposal results in the acceptance probability ( $\alpha$ ),

$$\begin{aligned} &\min \left\{ 1, \frac{s(\{\theta, \theta'\}|\{\theta^*, \theta'^*\})\pi(\theta', \theta, \mathbf{y}|\mathbf{x})}{s(\{\theta^*, \theta'^*\}|\{\theta, \theta'\})\pi(\theta, \theta', \mathbf{y}|\mathbf{x})} \right\} \\ &= \min \left\{ 1, \frac{p(\theta')h(\mathbf{x}|\theta')Z(\theta)h(\mathbf{y}|\theta)Z(\theta')q(\theta|\theta')}{p(\theta)h(\mathbf{x}|\theta)Z(\theta')h(\mathbf{y}|\theta')Z(\theta)q(\theta'|\theta)} \right\}, \end{aligned} \quad (4)$$

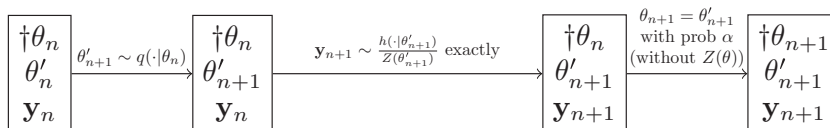
with intractable terms cancelled. Since  $\int_{\mathcal{X}} \int_{\Theta} \pi(\theta, \theta', \mathbf{y}|\mathbf{x})d\theta'd\mathbf{y} = \pi(\theta|\mathbf{x})$ , the marginal  $\theta$  samples follow the original target distribution, as shown in Figure 2. This idea of swapping parameter settings is connected to the parallel tempering algorithm (Geyer 1991). Parallel tempering swaps the states while preserving joint distributions in the product space. The exchange algorithm also swaps parameter settings between  $\mathbf{x}$  and  $\mathbf{y}$ , while preserving the augmented distribution. Since we do not need to keep  $\{\theta', \mathbf{y}\}$  in the actual implementation of the algorithm, the exchange algorithm may be simply written as Algorithm 2.

**Algorithm 2** Exchange algorithm

Given  $\theta_n \in \Theta$  at  $n$ th iteration.

1. Propose  $\theta' \sim q(\cdot|\theta_n)$ .
2. Generate the auxiliary variable exactly from probability model at  $\theta'$ :  $\mathbf{y} \sim \frac{h(\cdot|\theta')}{Z(\theta')}$  using perfect sampling.
3. Accept  $\theta_{n+1} = \theta'$  with probability  $\alpha = \min \left\{ 1, \frac{p(\theta')h(\mathbf{x}|\theta')h(\mathbf{y}|\theta_n)q(\theta_n|\theta')}{p(\theta_n)h(\mathbf{x}|\theta_n)h(\mathbf{y}|\theta')q(\theta'|\theta_n)} \right\}$ , else reject (set  $\{\theta_{n+1}, \mathbf{y}_{n+1}\} = \{\theta_n, \mathbf{y}_n\}$ ).

The exchange algorithm can also be extended through “bridging” (Murray, Ghahramani, and MacKay 2006), Algorithm 3. The idea may be summarized as follows. In the acceptance probability of Algorithm 2, one may think of the ratio  $h(\mathbf{x}|\theta')/h(\mathbf{x}|\theta)$  as measuring whether the proposed  $\theta'$  explains the data  $\mathbf{x}$  better than current  $\theta$  or not. Also,  $h(\mathbf{y}|\theta)/h(\mathbf{y}|\theta')$  represents how well the auxiliary variable  $\mathbf{y}$  is supported under  $\theta$  versus  $\theta'$ . Therefore, even if  $\theta'$  is a much better candidate than  $\theta$  under the data  $\mathbf{x}$ , swapping can be rejected if the auxiliary variable  $\mathbf{y}$  is improbable under  $\theta$ , which can lead to slow mixing of the chain. To improve this, annealed importance sampling (AIS; Neal 1996, 2001) can be combined with the exchange algorithm. Instead of sampling a single auxiliary variable, a series of auxiliary variables  $\{\mathbf{y}_0, \mathbf{y}_1, \dots, \mathbf{y}_K\}$  are sampled from the intermediate densities between  $h(\mathbf{y}|\theta')$  and  $h(\mathbf{y}|\theta)$ . Intermediate densities



**Figure 2.** Illustration for the exchange algorithm. The dagger symbols indicate the parameter of interest.

can be constructed as

$$f_k(\mathbf{y}|\theta, \theta') = h(\mathbf{y}|\theta')^{1-\beta_k} h(\mathbf{y}|\theta)^{\beta_k},$$

$$\beta_k = \frac{k}{K+1}, \quad k = 1, \dots, K. \quad (5)$$

After proposing  $\theta'$ , generate  $\mathbf{y}_0$  from  $h(\mathbf{y}|\theta')/Z(\theta')$ . Then generate each  $\mathbf{y}_k$  from  $T_k(\mathbf{y}_k|\mathbf{y}_{k-1})$ , the Metropolis–Hastings (MH) transition kernel whose stationary density is  $f_k(\mathbf{y}|\theta, \theta')$ . This bridging step helps to generate  $\mathbf{y}_K$  that is more probable under  $\theta$ .

### Algorithm 3 Exchange algorithm with bridging

Given  $\theta_n \in \Theta$  at  $n$ th iteration.  
 1. Propose  $\theta' \sim q(\cdot|\theta_n)$ .  
 2. Generate the series of auxiliary variables  $\{\mathbf{y}_0, \mathbf{y}_1, \dots, \mathbf{y}_K\}$  from (5):  $\mathbf{y}_0 \sim \frac{h(\cdot|\theta')}{Z(\theta')}$  using perfect sampling,  $\mathbf{y}_k \sim T_k(\mathbf{y}_k|\mathbf{y}_{k-1})$  via 1-step MH update for  $k = 1, \dots, K$ .  
 3. Accept  $\theta_{n+1} = \theta'$  with bridging probability  
 $\alpha = \min \left\{ 1, \frac{p(\theta')h(\mathbf{x}|\theta')q(\theta_n|\theta')}{p(\theta_n)h(\mathbf{x}|\theta_n)q(\theta'|\theta_n)} \prod_{k=0}^K \frac{f_{k+1}(\mathbf{y}_k|\theta_n, \theta')}{f_k(\mathbf{y}_k|\mathbf{y}_{k-1}, \theta')} \right\}$ ,  
 else reject (set  $\{\theta_{n+1}, \mathbf{y}_{n+1}\} = \{\theta_n, \mathbf{y}_n\}$ ).

*Components to be tuned:* The basic exchange algorithm does not require any tuning besides the usual tuning of the proposal for  $\theta$ . For the extended version, there are many options for bridging but even here (5) provides an automated schedule. Larger  $K$  in (5) can lead to better mixing of the chain at the expense of computing time. Effective sample size (ESS) (Kass et al. 1998) provides a rough diagnostic for how well the chain is mixing. Therefore, effective sample size per unit time (ESS/T) can be used to determine  $K$  that provides a compromise between mixing and computational cost (Murray 2007).

*Theoretical justification:* Both the exchange algorithm and the extended exchange algorithm with bridging are asymptotically exact MCMC in that constructed Markov chains satisfy detailed balance condition with respect to the target distribution. However just like AVM, the exchange algorithm is of limited applicability because it requires perfect sampling to generate the auxiliary variable.

### 2.3. The Double Metropolis–Hastings Sampler

Both AVM and the exchange algorithm require perfect sampling which can be very expensive or impossible for complicated probability models. Liang (2010) replaced perfect sampling for

the auxiliary variable with a standard Metropolis–Hastings algorithm; the last state of the resulting Markov chain is then treated like a draw from the perfect sampler. This is called double Metropolis–Hastings (DMH) because a Metropolis–Hastings algorithm is used within another Metropolis–Hastings algorithm. As in Figure 3, one is the “outer sampler” to generate  $\theta$  draws while the other is the “inner sampler” used to generate the auxiliary variable  $\mathbf{y}$ .

Define  $T_{\theta'}^m(\mathbf{y}|\mathbf{x})$  as  $m$ -MH updates from  $\mathbf{x}$  to  $\mathbf{y}$  under  $\theta'$  whose stationary density is  $h(\mathbf{y}|\theta')/Z(\theta')$ ; this satisfies the detailed balanced condition

$$h(\mathbf{x}|\theta')T_{\theta'}^m(\mathbf{y}|\mathbf{x}) = h(\mathbf{y}|\theta')T_{\theta'}^m(\mathbf{x}|\mathbf{y}). \quad (6)$$

Plugging this result into (4) yields

$$\alpha = \min \left\{ 1, \frac{p(\theta')h(\mathbf{y}|\theta')T_{\theta'}^m(\mathbf{x}|\mathbf{y})q(\theta|\theta')}{p(\theta)h(\mathbf{x}|\theta)T_{\theta'}^m(\mathbf{y}|\mathbf{x})q(\theta'|\theta)} \right\}, \quad (7)$$

which is the acceptance probability of the DMH algorithm. DMH approximates perfect sampling through  $m$  steps of an (inner) Metropolis–Hastings algorithm. A similar approximate approach which replaces perfect sampling with MH updates for the social network models is discussed in Caimo and Friel (2011). As we show later, DMH is simpler to implement and computationally more efficient relative to all other algorithms. However, its computational efficiency directly depends on the efficiency of the inner sampler. If the inner sampler update  $T_{\theta'}^m(\mathbf{y}|\mathbf{x})$  is expensive or needs large  $m$  to get probable auxiliary variable  $\mathbf{y}$ , then DMH becomes computationally expensive. Since  $m$  is usually proportional to the dimension of the data  $\mathbf{x}$ , the inner sampler is the main computational bottleneck for large data.

### Algorithm 4 Double Metropolis–Hastings algorithm

Given  $\theta_n \in \Theta$  at  $n$ th iteration.  
 1. Propose  $\theta' \sim q(\cdot|\theta_n)$ .  
 2. Generate the auxiliary variable approximately from probability model at  $\theta'$ :  $\mathbf{y} \sim T_{\theta'}^m(\cdot|\mathbf{x})$  using  $m$ -MH updates.  
 3. Accept  $\theta_{n+1} = \theta'$  with probability  
 $\alpha = \min \left\{ 1, \frac{p(\theta')h(\mathbf{x}|\theta')h(\mathbf{y}|\theta_n)q(\theta_n|\theta')}{p(\theta_n)h(\mathbf{x}|\theta_n)h(\mathbf{y}|\theta')q(\theta'|\theta_n)} \right\} =$   
 $\min \left\{ 1, \frac{p(\theta')h(\mathbf{y}|\theta_n)T_{\theta'}^m(\mathbf{x}|\mathbf{y})q(\theta_n|\theta')}{p(\theta_n)h(\mathbf{x}|\theta_n)T_{\theta'}^m(\mathbf{y}|\mathbf{x})q(\theta'|\theta_n)} \right\}$ ,  
 else reject (set  $\{\theta_{n+1}, \mathbf{y}_{n+1}\} = \{\theta_n, \mathbf{y}_n\}$ ).

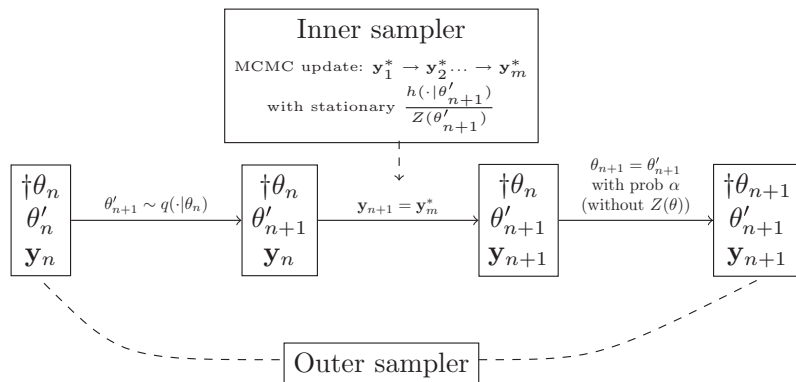


Figure 3. Illustration of the double Metropolis–Hastings algorithm. The dagger symbols indicate the parameter of interest.



*Components to be tuned:* Users need to decide the number of MH updates  $m$  to generate  $\mathbf{y}$ . There are numerous stopping rules for MCMC if computing expectations is the goal (see Flegal, Haran, and Jones 2008; Gong and Flegal 2015). However, here determining convergence to the stationary distribution is of interest (Rosenthal 1995; Jones and Hobert 2001), which is a very challenging problem in general. A simple heuristic for determining  $m$  is to set  $m$  to be proportional to the size of the data, for instance  $m = 5n$ , where  $n$  is the size of the data. We recommend that DMH should be run for larger  $m$  values, for instance  $m = 10n$  to confirm that the results do not change much as  $m$  is increased.

*Theoretical justification:* DMH is an asymptotically inexact algorithm because the detailed balance condition does not hold for the outer sampler unless the inner sampler length ( $m$ ) approaches infinity; in practice the inner sampler length is of course finite.

## 2.4. Adaptive Exchange Algorithm

To address the asymptotic inexactness of DMH, Liang et al. (2016) proposed the adaptive exchange algorithm (AEX). AEX runs two chains simultaneously: the auxiliary chain and the target chain (see Figure 4). The auxiliary chain generates a sample from a family of distributions,  $\{h(\mathbf{x}|\theta^{(1)})/Z(\theta^{(1)}), \dots, h(\mathbf{x}|\theta^{(d)})/Z(\theta^{(d)})\}$ , where  $\{\theta^{(1)}, \dots, \theta^{(d)}\}$  are prespecified “particles” covering a parameter space. The generated sample is stored at each iteration. Then the target chain generates a posterior sample from the  $\pi(\theta|\mathbf{x})$  via the exchange algorithm. The difference is that the auxiliary variable  $\mathbf{y}$  is sampled from the cumulative samples in the auxiliary chain until current iteration (resampling). With increasing iterations, cumulative samples from the auxiliary chain grow, so that the resampling of  $\mathbf{y}$  in the target chain becomes identical to exact sampling of  $\mathbf{y}$ . Then similar to the exchange algorithm, the marginal density of  $\theta$  converges to the target  $\pi(\theta|\mathbf{x})$ . AEX is therefore an asymptotically exact algorithm that does not require perfect sampling. We begin with some notation for the  $n$ th iteration of AEX.

- Particles:  $\{\theta^{(1)}, \dots, \theta^{(d)}\}$ , each  $\theta^{(i)} \in \Theta$ . These particles are fixed through the algorithm.
- Particle index:  $I_n \in \{1, \dots, d\}$  returns index of a chosen particle at  $n$ th iteration.

- Auxiliary data:  $\mathbf{x}_n \in \mathcal{X}$  is an approximate sample from the probability model at selected particle  $\theta^{(I_n)}$ . That is,  $\mathbf{x}_n \sim h(\cdot|\theta^{(I_n)})/Z(\theta^{(I_n)})$ .
- Normalizing function approximation at each particle:  $\mathbf{w}_n = \{w_n^{(1)}, \dots, w_n^{(d)}\} \in W$ . For  $i = 1, \dots, d$ , as  $n$  gets large  $w_n^{(i)}$  converges to  $Z(\theta^{(i)})$  via the stochastic approximation algorithm (Liang, Liu, and Carroll 2007).
- Cumulative information:  $H_n = \cup_{j=1}^n \{I_j, \mathbf{x}_j, \mathbf{w}_j\}$  is necessary for constructing the algorithm. Therefore,  $\{I_n, \mathbf{x}_n, \mathbf{w}_n\}$  should be stored at each iteration.
- Posterior sample:  $\theta_n \in \Theta$

AEX updates a non-Markovian stochastic process  $\{I_{n+1}, \mathbf{x}_{n+1}, \mathbf{w}_{n+1}, \theta_{n+1}\} \in \{1, \dots, d\} \times \mathcal{X} \times R^d \times \Theta$  in the  $(n+1)$ st iteration. For each iteration, the auxiliary chain generates a sample  $\mathbf{x}_{n+1}$  from the mixture density,  $(1/d) \sum_{i=1}^d h(\mathbf{x}|\theta^{(i)})/Z(\theta^{(i)})$  through stochastic approximation Monte Carlo (SAMC; Liang, Liu, and Carroll 2007). Particles are  $d$ -number of points in the grid over a parameter space  $\Theta$ . Particles are chosen so that they can cover the important region of the parameter space and the sample spaces of neighboring distributions  $\{h(\mathbf{x}|\theta^{(i)})/Z(\theta^{(i)})\}$  overlap each other. A clever choice of particles can improve mixing of the auxiliary chain and is important for AEX to be a practical algorithm. We refer the reader to the supplementary material for strategies for choosing particles. For the  $(n+1)$ st update of AEX, cumulative information  $H_{n+1}$  is constructed through the auxiliary chain. For  $(n+1)$ st iteration of the algorithm,  $\{I_{n+1}, \mathbf{x}_{n+1}, \mathbf{w}_{n+1}\}$  is updated for given  $\{I_n, \mathbf{x}_n, \mathbf{w}_n\}$  using SAMC and then added to  $H_n$  to construct  $H_{n+1}$ . Hence, the accumulated information ( $H_{n+1}$ ) becomes larger with each iteration of AEX. See Algorithm 5 for details.

Then the target chain generates  $\theta_{n+1}$  from the posterior using the exchange algorithm. For proposed  $\theta'$ , the auxiliary variable  $\mathbf{y}$  is sampled from  $\{\mathbf{x}_1, \dots, \mathbf{x}_{n+1}\}$  through a dynamic importance sampling procedure. The resampling probability is

$$P(\mathbf{y} = \mathbf{x}_l | \theta') \propto \sum_{j=1}^{n+1} w_j^{(I_j)} \frac{h(\mathbf{x}_j | \theta')}{h(\mathbf{x}_j | \theta^{(I_j)})} 1\{\mathbf{x}_j = \mathbf{x}_l\},$$

$$l = 1, \dots, n+1. \quad (8)$$

In (8),  $h(\mathbf{x}_j | \theta^{(I_j)})/w_j^{(I_j)}$  is the importance function at  $j$ th iteration. Since the importance function changes as  $I_j$  varies, it is called a dynamic importance function. Here,  $\mathbf{x}_l$ 's more

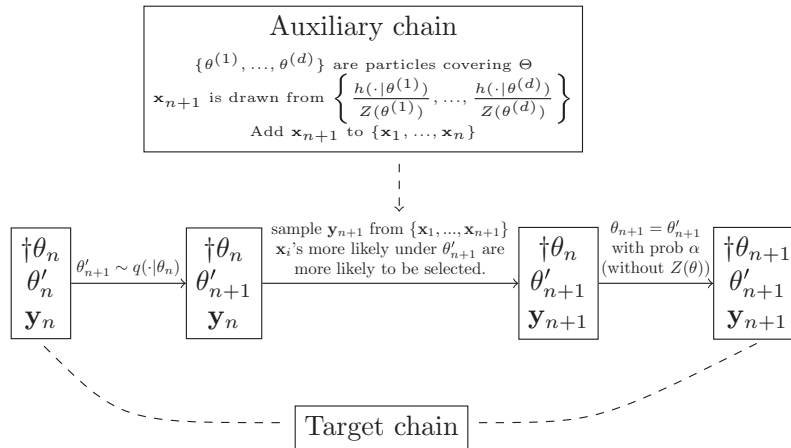


Figure 4. Illustration for the adaptive exchange algorithm. The dagger symbols indicate the parameter of interest.

likely under  $\theta'$  are more likely to be sampled. Liang et al. (2016) showed that the distribution of resampled  $\mathbf{y}$  converges to  $h(\mathbf{y}|\theta')/Z(\theta)$  as  $n$  increases. Since the proposal for  $\mathbf{y}$  in (8) changes with each iteration, AEX appears at first to be an adaptive MCMC algorithm. However, unlike basic adaptive MCMC algorithms the target distribution also varies as  $H_{n+1}$  grows.

In Step 2(a) of Algorithm 5,  $\{a_n\}$  is called the gain factor, which is the step size of update for the abundance factor  $\mathbf{w}_n$ . A large gain factor can lead to faster movement around all  $\{h(\mathbf{x}|\theta^{(i)})/Z(\theta^{(i)})\}$ . Since  $a_n = n_0/\max(n_0, n)$  becomes smaller as  $n$  increases, movement between particles becomes slower with increasing iterations. Liang, Liu, and Carroll (2007) suggested larger  $n_0$  for complicated problems so that makes it fast to move around all state space. Step 2(b) illustrates varying truncation (Jin and Liang 2013), which can help the convergence of  $\mathbf{w}_n$  regardless of starting point. For  $\mathbf{w}_n \in W$ , consider the sequence of set  $\{K_s\} \in W$ , which satisfies  $\cup_{s \geq 0} K_s = W$  and  $K_s \subset K_{s+1}^o$ , where  $K_{s+1}^o$  denotes the interior of set  $K_{s+1}$ . A truncation function  $T$  maps a point in  $\mathcal{X} \times W$  to a random point in  $\mathcal{X}_0 \times K_0$  for  $\mathcal{X}_0 \subset \mathcal{X}$ , and  $\sigma_n$  is the number of truncations that occurred by the  $n$ th iteration ( $\sigma_0 = 0$ ). From such truncation, when an updated  $\mathbf{w}_{n+1/2}$  is outside the interesting target region ( $K_s$ ), it is reinitialized with a smaller gain factor and expanded target region ( $K_{s+1} \supset K_s$ ). This varying truncation can help SAMC to find an appropriate step size  $\{a_n\}$  and starting point of abundance factor automatically (Jin and Liang 2013). The reader may understand such varying truncation method as a safeguard for when the abundance factor becomes degenerate.

---

**Algorithm 5** Adaptive Exchange algorithm (**Part 1**): Auxiliary chain Constructing cumulative information  $H_{n+1}$

---

Initialize  $\{I_0, \mathbf{x}_0, \mathbf{w}_0\}$ , for example,  $\mathbf{x}_0 = \text{data } \mathbf{x}$ ,  $I_0 \sim \{1, \dots, d\}$  uniformly,  $\mathbf{w}_0 = \mathbf{1}$ . Set  $H_0 = \{I_0, \mathbf{x}_0, w_0^{(I_0)}\}$ .

For  $(n+1)$ st update, given  $\{I_n, \mathbf{x}_n, \mathbf{w}_n, \theta_n\} \in \{1, \dots, d\} \times \mathcal{X} \times W \times \Theta$ .

1. Update  $I_{n+1}$  or  $\mathbf{x}_{n+1}$  with equal probability.

(a) With probability 0.5, update  $I_{n+1}$  at  $\mathbf{x}_n$ . Propose  $\theta^{(I')}$  from the  $k$ -nearest particles (say in Euclidean distance) of  $\theta^{(I_n)}$  with equal probability. Accept  $\{I_{n+1}, \mathbf{x}_{n+1}\} = \{I', \mathbf{x}_n\}$  with  $\alpha_1 = \min\{1, \frac{w_n^{(I_n)} h(\mathbf{x}_n|\theta^{(I')})}{w_n^{(I')} h(\mathbf{x}_n|\theta^{(I_n)})}\}$ .

(b) With probability 0.5, update  $\mathbf{x}_{n+1}$  at  $\theta^{(I_n)}$ . Propose  $\mathbf{x}'$  through  $m$ -MH updates:  $\mathbf{x}' \sim T_{I_n}^m(\cdot|\mathbf{x}_n)$  whose target is  $\frac{h(\mathbf{x}|\theta^{(I_n)})}{Z(\theta^{(I_n)})}$ . Accept  $\{I_{n+1}, \mathbf{x}_{n+1}\} = \{I_n, \mathbf{x}'\}$  with  $\alpha_2 = \min\{1, \frac{h(\mathbf{x}'|\theta^{(I_n)})T(\mathbf{x}_n|\mathbf{x}')}{h(\mathbf{x}_n|\theta^{(I_n)})T(\mathbf{x}'|\mathbf{x}_n)}\}$ .

2. Update approximation of  $Z(\theta^{(i)})$  at  $i = 1, \dots, d$ :

Following stochastic approximation (Liang, Liu, and Carroll 2007)

(a)  $\log(w_{n+1/2}^{(i)}) = \log(w_n^{(i)}) + a_{n+1}(1\{I_{n+1} = i\} - 1/d)$ ,  $a_n = n_0/\max(n_0, n)$ .

(b) Stochastic truncation (T) is implemented if  $\mathbf{w}_{n+1/2}$  is outside of the target region ( $K_{\sigma_n}$ ), and records the number of truncation ( $\sigma_n$ ):  $\{\mathbf{w}_{n+1}, \mathbf{x}_{n+1}\} = \begin{cases} T(\{\mathbf{w}_n, \mathbf{x}_n\}), \sigma_{n+1} = \sigma_n + 1 & \mathbf{w}_{n+1/2} \notin K_{\sigma_n} \\ \{\mathbf{w}_{n+1/2}, \mathbf{x}_{n+1}\}, \sigma_{n+1} = \sigma_n & \text{o.w.} \end{cases}$

3. Cumulative information is updated:  $H_{n+1} = H_n \cup \{I_{n+1}, \mathbf{x}_{n+1}, w_{n+1}^{(I_{n+1})}\}$ .

---

**Algorithm 5** Adaptive Exchange algorithm (**Part 2**): Target chain Obtain  $\theta_{n+1}$  approximately from  $\pi(\theta|\mathbf{x})$  by exchange algorithm, using resampled  $\mathbf{y}$

---

Continued from Part 1: for  $(n+1)$ st update,

4. Propose  $\theta'$ :  $\theta' \sim q(\cdot|\theta_n)$ .

5. Sample the auxiliary variable from collected dataset in auxiliary chain:  $\mathbf{y} \sim \{\mathbf{x}_1, \dots, \mathbf{x}_{n+1}\}$ , with resampling probabilities as  $P(\mathbf{y} = \mathbf{x}_l|\theta') \propto \sum_{j=1}^{|H_{n+1}|} w_j^{(I_j)} \frac{h(\mathbf{x}_j|\theta')}{h(\mathbf{x}_j|\theta^{(I_j)})} 1\{\mathbf{x}_j = \mathbf{x}_l\}$ ,  $l = 1, \dots, |H_{n+1}|$ .

6. Accept  $\theta_{n+1} = \theta'$  with probability

$\alpha_3 = \min\{1, \frac{p(\theta')h(\mathbf{x}|\theta')h(\mathbf{y}|\theta_n)q(\theta_n|\theta')}{p(\theta_n)h(\mathbf{x}|\theta_n)h(\mathbf{y}|\theta')q(\theta'|\theta_n)}\}$ , else reject (set  $\{\theta_{n+1}, \mathbf{y}_{n+1}\} = \{\theta_n, \mathbf{y}_n\}$ ).

---

In practice, Liang et al. (2016) suggested that only the auxiliary chain (Algorithm 5 Part 1) is implemented  $N_1$  iterations at first to construct the database, and then the auxiliary and target chain can be run simultaneously. In detail, through  $N_1$  preliminary iterations of the auxiliary chain, information:  $\cup_{j=1}^{N_1} \{I_j, \mathbf{x}_j, \mathbf{w}_j\}$  is constructed. Using this information as initial  $H_0 = \cup_{j=1}^{N_1} \{I_j, \mathbf{x}_j, \mathbf{w}_j\}$ , we can implement the entire Algorithm 5. In this case, the size of cumulative information becomes  $|H_0| = N_1$ , and  $|H_{n+1}| = N_1 + n + 1$  (without preliminary step,  $|H_0| = 0$ , and  $|H_{n+1}| = n + 1$ ). The reason of conducting preliminary step is because the early performance of the target chain can be affected by initial starting values of the auxiliary chain. This preliminary step can construct an initial database  $H_0$  and improves the mixing of the target chain.

Although we classify AEX as an auxiliary variable approach, AEX also shares characteristics with likelihood approximation approaches. The abundance factor  $\{w_j^{(i)}\}$  converges to  $\{Z(\theta^{(i)})\}$  for each particle, which is guaranteed by SAMC. This implies that likelihood approximation approaches are applied to estimate  $Z(\theta^{(i)})$  directly. Then the auxiliary variable is resampled to cancel out  $Z(\theta)$  in the target chain. A strength of AEX is its broad applicability, while still remaining asymptotically exact. It is reasonably efficient computationally. However because of its adaptive structure, there can be serious memory issues. When there are low-dimensional sufficient statistics, as is typically the case with exponential family models, only the sufficient statistics of  $\mathbf{x}_{n+1}$  need to be stored in Step 3 of the Algorithm 5. However, in the absence of such sufficient statistics,  $\mathbf{x}_{n+1}$  itself need to be stored at each step. Furthermore, without sufficient statistics, resampling probability calculations in Step 5 become expensive because  $h(\mathbf{x}_j|\theta')$  should be recalculated; when there are sufficient statistics, one can simply take the product of  $\theta'$  and the sufficient statistic.

*Components to be tuned:* We provide details about choosing particles (number of particles, strategies for choosing particles) in the supplementary material. In the auxiliary chain, the number of neighbors for updating  $\theta^{(I_n)}$ , the number of preliminary runs for the auxiliary chain ( $N_1$ ), the number of MH updates ( $m$ ), gain factor component ( $n_0$ ), and the target region ( $K_s$ ) of abundance factor should be tuned. For particles, neighbor is defined using the distance from each other. We can define the closest 20 points as neighboring particles, and this choice appears to work well both in simulated examples in this manuscript as well as the examples in Liang et al. (2016). For

updating  $\mathbf{x}_n$ ,  $m = 1$  cycle of MH updates are enough, because we do not need to generate independent samples. For complicated problems, Liang et al. (2016) recommend larger values for  $N_1$  and  $n_0$ . For the simulation examples, Liang et al. (2016) used around  $N_1 = 10,000d$ ,  $n_0 = 25,000$  where  $d$  is the number of particles. We used similar values in our simulated examples. If the preliminary auxiliary chain does not appear to converge, larger  $n_0$  and  $N_1$  should be used. The convergence of the auxiliary chain can be checked by comparing sampling frequencies  $v(i)/N_1$  with the target probabilities  $1/d$ , where  $v(i)$  is the number of visitations at each particle. If  $v(i)/N_1 \approx 1/d$  for each particle, then the preliminary run of auxiliary chain can be diagnosed as having converged. For the choice of compact subset of  $W$ , Liang et al. (2016) set  $K_s = [0, 10^{100+10s}]^d$  and  $\mathcal{X}_0 = \mathcal{X}$  in their examples. Detailed simulation settings about the social network model also can be checked in Jin, Yuan, and Liang (2013).

*Theoretical justification:* The AEX is an asymptotically exact algorithm and the ergodic average satisfies the Weak Law of Large Numbers. Since AEX is adaptive in both proposal and target, conventional theory for adaptive MCMC does not directly apply. Liang et al. (2016) extended a result in (Roberts and Rosenthal 2007, Theorem 1), which shows the ergodicity of the adaptive chain with the same target distribution. Although the original proof (Liang et al. 2016) assumes compactness of the parameter space  $\Theta$ , Jin, Yuan, and Liang (2013) extended the results without a compactness assumption for  $\Theta$ . The details of the assumptions for the proof(s) are provided across several results in (Jin, Yuan, and Liang 2013, Lemma 3.1, 3.2 and Theorem 3.1). To make it easy for readers to understand the assumptions from a practical point of view, we have attempted to distill the key assumptions as follows: (1) both  $\mathcal{X}$  and  $W$  are compact, (2)  $h(\mathbf{x}|\theta)$  is bounded away from 0 and  $\infty$ , (3)  $\lim_{n \rightarrow \infty} a_n = 0$ ,  $\sum_{n=1}^{\infty} a_n = \infty$ ,  $\sum_{n=1}^{\infty} a_n^\eta < \infty$  for some  $\eta \in (1, 2]$ , (4) Doeblin condition holds for kernel of auxiliary chain. Some of these conditions are easily verified in practice. The sample space  $\mathcal{X}$  for data is compact for finite lattice problem and point process with finite number of points on bounded domain. Sample space  $W$  for abundance factor is also compact, if we set  $K_s$  as suggested in the components to be tuned above. The second assumption is also satisfied for realistic parameter settings. The third is a technical assumption for SAMC, and is satisfied if we construct  $a_n = n_0 / \max(n_0, n)$ . The sufficient condition for the last assumption is local positiveness of  $m$ -MH updates in Step 1(b) of Algorithm 5, which means  $\exists \delta > 0, \epsilon > 0$  such that,  $\forall \mathbf{x} \in \mathcal{X}$ ,  $|\mathbf{x} - \mathbf{y}| \leq \delta$  implies  $T(\mathbf{y}|\mathbf{x}) \geq \epsilon$ . This is satisfied for a simple Gibbs sampler, birth-death sampler in point process, and tie-no-tie sampler in social network models. Therefore, broadly speaking, the assumptions hold in most problems. We point the readers to the manuscript to find the complete set of assumptions for these results.

### 3. Likelihood Approximation Methods

While the auxiliary variable approaches we have discussed so far avoid approximating  $Z(\theta)$ , the approaches based on likelihood approximations we discuss in this section directly approximate  $Z(\theta)$  through Monte Carlo and substitute the approximation into the Metropolis–Hastings acceptance probability.

#### 3.1. Atchade, Lartillot, and Robert's Adaptive MCMC

Atchade, Lartillot, and Robert (2008) provided an adaptive MCMC algorithm, henceforth the ALR algorithm, which approximates  $Z(\theta)$  through importance sampling in the acceptance probability. ALR constructs a non-Markovian stochastic process and approximates  $Z(\theta)$  at every iteration using the entire sample path of the process. ALR is an asymptotically exact algorithm that does not require perfect sampling.

The ALR algorithm uses importance sampling ideas developed in MCMC-MLE (Geyer and Thompson 1992). Consider  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ , samples from  $h(\mathbf{x}|\theta^{(0)})/Z(\theta^{(0)})$  for some  $\theta^{(0)}$ . Then  $Z(\theta)/Z(\theta^{(0)})$  is an expectation that can be approximated by importance sampling. However the approximation might be poor if  $\theta$  is far from  $\theta^{(0)}$ . Atchade, Lartillot, and Robert (2008) introduced multiple particles,  $\{\theta^{(1)}, \dots, \theta^{(d)}\}$ , and a linear combination of importance sampling estimates for  $\{Z(\theta)/Z(\theta^{(1)}), \dots, Z(\theta)/Z(\theta^{(d)})\}$  approximates  $Z(\theta)$ . Larger weights are assigned to the estimate of  $Z(\theta)/Z(\theta^{(i)})$  when  $\theta$  and  $\theta^{(i)}$  are closer to each other. This idea of sampling from a mixture is related to umbrella sampling (Torrie and Valleau 1977) which provides a more robust approximation than a single point importance sampling estimate, as long as one of the  $\theta^{(i)}$ 's is close to  $\theta$ . We begin with some notation for the  $n$ th iteration of ALR.

- Particles:  $\{\theta^{(1)}, \dots, \theta^{(d)}\}$ , each  $\theta^{(i)} \in \Theta$ . These particles are fixed through the algorithm.
- Particle index:  $I_n \in \{1, \dots, d\}$  returns the index of the selected particle at  $n$ th iteration.
- Auxiliary data:  $\mathbf{x}_n \in \mathcal{X}$  is an approximate sample from the probability model at a selected particle  $\theta^{(I_n)}$ , that is,  $\mathbf{x}_n \sim h(\cdot|\theta^{(I_n)})/Z(\theta^{(I_n)})$ .
- Normalizing function approximation at each particle:  $\mathbf{c}_n = \{c_n^{(1)}, \dots, c_n^{(d)}\} \in \mathbb{R}^d$ . For  $i = 1, \dots, d$ , as  $n$  gets large  $c_n^{(i)}$  converges to  $\log Z(\theta^{(i)})$  by stochastic approximation (Wang and Landau 2001; Atchade and Liu 2004).
- Cumulative information:  $H_n = \cup_{j=1}^n \{I_j, \mathbf{x}_j\}$  is necessary for constructing the algorithm. Therefore,  $\{I_n, \mathbf{x}_n\}$  should be stored at each iteration.
- Posterior sample:  $\theta_n \in \Theta$ .

By extending a result in Wang and Landau (2001) and Atchade and Liu (2004), the ALR algorithm constructs a non-Markovian stochastic process  $\{\mathbf{x}_{n+1}, I_{n+1}, \mathbf{c}_{n+1}, \theta_{n+1}\} \in \mathcal{X} \times \{1, \dots, d\} \times \mathbb{R}^d \times \Theta$ .  $Z(\theta)$  is approximated through multiple importance sampling and plugged into the acceptance probability to generate samples from  $\pi(\theta|\mathbf{x})$ . A stochastic approximation for estimating marginal densities, similar to the approximation used in ALR, is also described in Liang (2007). The decomposition of  $Z(\theta)$  is

$$\begin{aligned} Z(\theta) &= \int_{\mathcal{X}} h(\mathbf{x}|\theta) d\mathbf{x} \\ &= \sum_{i=1}^d k(\theta, \theta^{(i)}) \int_{\mathcal{X}} \frac{h(\mathbf{x}|\theta)}{h(\mathbf{x}|\theta^{(i)})} h(\mathbf{x}|\theta^{(i)}) d\mathbf{x} \\ &= \sum_{i=1}^d k(\theta, \theta^{(i)}) Z(\theta^{(i)}) \int_{\mathcal{X}} \frac{h(\mathbf{x}|\theta)}{h(\mathbf{x}|\theta^{(i)})} \frac{h(\mathbf{x}|\theta^{(i)})}{Z(\theta^{(i)})} d\mathbf{x}. \end{aligned} \quad (9)$$



This suggests the following Monte Carlo approximation of  $Z(\theta)$ ,

$$\hat{Z}_{n+1}(\theta) = \sum_{i=1}^d k(\theta, \theta^{(i)}) \exp(c_{n+1}^{(i)}) \frac{\sum_{k=1}^{n+1} \frac{h(\mathbf{x}_k|\theta)}{h(\mathbf{x}_k|\theta^{(i)})} 1\{I_k = i\}}{\sum_{k=1}^{n+1} 1\{I_k = i\}}, \quad (10)$$

a weighted sum of the importance sampling estimate  $Z(\theta^{(i)})$  at each particle. A “similarity kernel”  $k(\theta, \theta^{(i)})$  measures the distance between  $\theta$  and  $\theta^{(i)}$  and assigns larger weights to particles closer to  $\theta$ . Similar to AEX, particles are prespecified before starting the algorithm and are chosen so that they can cover the parameter space. In the examples in Atchade, Lartillot, and Robert (2008), particles are collected from the stochastic approximation (Younes 1988), though particles could also be chosen as in AEX (see supplementary material for details). Algorithm 6 provides ALR details.

**Algorithm 6** ALR algorithm (Part 1): Find good initial values  $\{\mathbf{x}_\tau, I_\tau, \mathbf{c}_\tau\}$  for use in Part 2

Initialize  $\{\mathbf{x}_0, I_0, \mathbf{c}_0\}$ , for example,  $\mathbf{x}_0 = \text{data } \mathbf{x}$ ,  $I_0 \sim \{1, \dots, d\}$  uniformly,  $\mathbf{c}_0 = \mathbf{0}$ .

For  $(n+1)$ st update, given  $\{\mathbf{x}_n, I_n, \mathbf{c}_n\} \in \mathcal{X} \times \{1, \dots, d\} \times \mathbb{R}^d$ .

**while**  $\gamma_{n+1} > \epsilon_1$  **do** ( $\gamma_n$  is step size used in step 3 below)

Reset  $\mathbf{v} = \mathbf{0} \in \mathbb{R}^d$  ( $v(i)$  counts the number of visitations at each particle  $\theta^{(i)}$ )

**while**  $\max_i |v(i) - \bar{v}| > \epsilon_2 \bar{v}$  **do** (until each particle has been visited equally)

For  $i = 1, \dots, d$ , set  $v(i) = v(i) + 1_i(I_{n+1})$ .

1. Update  $\mathbf{x}_{n+1}$  through  $m$ -MH step for given particle:

$\mathbf{x}_{n+1} \sim T_{I_n}^m(\cdot | \mathbf{x}_n)$  whose stationary density is  $\frac{h(\mathbf{x}|\theta^{(I_n)})}{Z(\theta^{(I_n)})}$ .

2. Decide where to visit in the next iteration:

$I_{n+1} \sim \{1, \dots, d\}$  with probabilities  $h(\mathbf{x}_{n+1}|\theta^{(i)}) \exp(-c_n^{(i)})$  for  $i = 1, \dots, d$ .

3. Update approximation of  $\log Z(\theta^{(i)})$ : Following stochastic approximation (Wang and Landau 2001; Atchade and Liu 2004)

$$c_{n+1}^{(i)} = c_n^{(i)} + \gamma_n \frac{h(\mathbf{x}_{n+1}|\theta^{(i)}) \exp(-c_n^{(i)})}{\sum_{j=1}^d h(\mathbf{x}_{n+1}|\theta^{(j)}) \exp(-c_n^{(j)})} \quad \text{for}$$

$i = 1, \dots, d$ , and  $\gamma_{n+1} = \gamma_n$ .

**end while**

Step size become smaller:  $\gamma_{n+1} = \gamma_n/2$

**end while**

Return  $\{\mathbf{x}_\tau, I_\tau, \mathbf{c}_\tau\}$ , where  $\tau$  is the total number of iterations for Part 1.

We split the ALR algorithms into two parts (see Figure 5). The purpose of Part 1 of the Algorithm 6 is to obtain a reasonable starting value for the normalizing function approximation without the memory and computational costs associated with updating the entire process  $\{\mathbf{x}_n, I_n, \mathbf{c}_n, \theta_n\}$ . The step size  $\{\gamma_n\}$  plays an important role in the convergence of  $\mathbf{c}_n$  (analogous to how the gain factor  $a_n$  is important for the abundance factor  $\mathbf{w}_n$  in AEX). If  $\gamma_n$  is too small, it takes a long time for convergence. If  $\gamma_n$  is too large it can lead to large variance of  $\mathbf{c}_n$ . Therefore in the beginning of the algorithm,  $\gamma_n$  is set to be a large value so that it can

**Algorithm 7** ALR algorithm (Part 2): Update the entire process.

Initialize  $\{\mathbf{x}_0, I_0, \mathbf{c}_0\} = \{\mathbf{x}_\tau, I_\tau, \mathbf{c}_\tau\}$  from Part 1 and set  $H_0 = \{\mathbf{x}_0, I_0\}$ .

For  $(n+1)$ st update, obtain  $\{\mathbf{x}_{n+1}, I_{n+1}, \mathbf{c}_{n+1}\}$  using Step 1-3 in the Part 1 with deterministic step size,  $\gamma_{n+1} = \epsilon_1/(n+1)^{0.7}$ .

Then append dataset:  $H_{n+1} = H_n \cup \{\mathbf{x}_{n+1}, I_{n+1}\}$ .

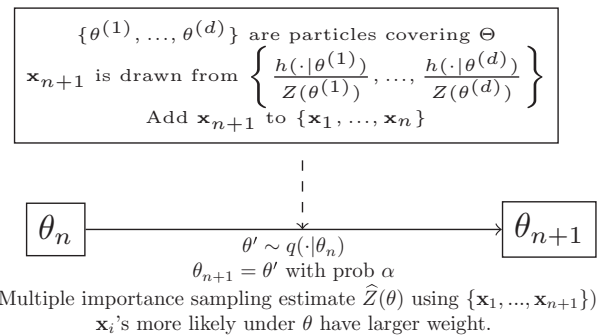
4. Approximate  $Z(\theta)$  adaptively using all previous samples ( $H_{n+1}$ ) and  $\mathbf{c}_{n+1}$ :  $\hat{Z}_{n+1}(\theta) = \sum_{i=1}^d k(\theta, \theta^{(i)}) \exp(c_{n+1}^{(i)}) \frac{\sum_{k=1}^{n+1} \frac{h(\mathbf{x}_k|\theta)}{h(\mathbf{x}_k|\theta^{(i)})} 1\{I_k = i\}}{\sum_{k=1}^{n+1} 1\{I_k = i\}}$ .

5. Propose  $\theta' \sim q(\cdot|\theta)$  and accept  $\theta_{n+1} = \theta'$  with probability  $\alpha = \min \left\{ 1, \frac{p(\theta')h(\mathbf{x}|\theta')\hat{Z}_{n+1}(\theta_n)q(\theta_n|\theta')}{p(\theta_n)h(\mathbf{x}|\theta_n)\hat{Z}_{n+1}(\theta')q(\theta'|\theta_n)} \right\}$ , else reject (set  $\{\theta_{n+1}, \mathbf{y}_{n+1}\} = \{\theta_n, \mathbf{y}_n\}$ ).

move around the space fast; it slowly becomes smaller after moving around the state space. In Algorithm 6 (Part 1), the number of visits to each particle is recorded through  $\mathbf{v}$ . If  $\mathbf{v}$  is close to uniform distribution ( $\{I_n\}$  has visited  $\{1, \dots, d\}$  about equally),  $\gamma_n$  becomes smaller. If step size  $\gamma_n$  is smaller than some convergence criteria  $\epsilon_1$ , we can assume  $\mathbf{c}_n$  has converged to a reasonable value.

The strengths and weakness of ALR are similar to those of AEX: ALR is asymptotically exact without requiring perfect sampling. However, the computational and memory costs can be large. Without low-dimensional sufficient statistics, the entire  $\mathbf{x}_n$  chain needs to be stored with each iteration to calculate  $\hat{Z}_{n+1}(\theta)$  in Step 4 of the Algorithm 6. Furthermore, without sufficient statistics, calculations in Step 3 and Step 4 become expensive, because  $h(\mathbf{x}_k|\theta)$ ,  $h(\mathbf{x}_k|\theta^{(i)})$  need to be recalculated; with sufficient statistics, one can simply take the product of  $\theta$  or  $\theta^{(i)}$  and the sufficient statistic of  $\mathbf{x}_k$ .

*Components to be tuned:* In this algorithm, step size ( $\gamma_n$ ), the convergence checking components ( $\epsilon_1, \epsilon_2$ ), number of MH updates ( $m$ ), number of particles ( $d$ ), and kernel  $k(\theta, \theta^{(i)})$  need to be tuned. Atchade, Lartillot, and Robert (2008) set initial  $\gamma_0$  as 1,  $\epsilon_1 = 0.001$ ,  $\epsilon_2 = 0.2$ . Under these settings, consider the sequence of bounded stopping times  $0 = \tau_0 < \tau_1 < \dots < \tau_{10} = \tau$ , where  $\tau$  is the total number of iteration in Algorithm 6 (Part 1). Until  $\tau_1$ , initial  $\gamma_0$  is used as step size. For this step size,  $\tau_1$  is the stopping time until each particle has been visited equally (i.e.,  $\max_{i=1, \dots, d} |v(i) - \bar{v}| \leq \epsilon_2 \bar{v}$ ). Then  $\gamma_{\tau_1+1}$  becomes  $\gamma_0/2 = 1/2$  and is kept until  $\tau_2$ . This is repeated until  $\tau_{10}$  where  $\gamma_{\tau_{10}+1} = 1/2^{10} < 0.001 = \epsilon_1$ . Likewise, step size is controlled to



**Figure 5.** Illustration for Atchade, Lartillot, and Robert's (ALR) algorithm.

hasten convergence of  $\mathbf{c}_n$ . Once  $\mathbf{c}_n$  appears to have converged, Algorithm 6 (Part 2) is implemented.  $\gamma_{n+1}$  is updated as the deterministic sequence  $0.001/(n+1)^{0.7}$  in Algorithm 6 (Part 2). Similar to AEX,  $m = 1$  cycle of MH updates are enough in practice. For the number of particles,  $d = 100p$  appears to work well in practice to cover the  $p$ -dimensional parameter space. Although various choice of kernel  $k(\theta, \theta^{(i)})$  is possible, uniform kernel with bandwidth  $h$  is used in this manuscript. This kernel gives  $1/h$  weights for the  $h$  closest particles and gives 0 weights for others. Bandwidth should be determined by trials and errors and we used  $h = 20$ .

*Theoretical justification:* The ALR is an asymptotically exact algorithm and the ergodic average from the chain satisfies the Strong Law of Large Numbers. Atchade, Lartillot, and Robert (2008) provided a proof that  $\{\theta_n\}$  from the generated process converges to the target distribution exactly. Three assumptions are required: (1)  $h(\mathbf{x}|\theta)$  is bounded away from 0 and  $\infty$ , (2)  $q^{n_0}(\theta'|\theta) \geq \epsilon$  for all  $\theta, \theta' \in \Theta$  where  $q$  is proposal density and  $n_0 \geq 1$  is an integer, (3)  $\{\gamma_n\}$  is random and adaptively updated based on the previous generated process, which satisfies  $\sum \gamma_n = \infty$  and  $\sum \gamma_n^2 < \infty$  almost surely. The first assumption holds for finite  $\mathcal{X}$  and realistic parameter settings. The second assumption also holds for most symmetric kernels. The third assumption is technical, and taken from Wang and Landau (2001). Because we can typically control the sequence  $\{\gamma_n\}$ , this assumption is also generally easy to satisfy. Therefore, the assumptions appear to hold in most cases in practice.

### 3.2. Pseudo-Marginal MCMC

Suppose  $\hat{L}(\theta|\mathbf{x})$  is a positive and unbiased Monte Carlo approximation of  $L(\theta|\mathbf{x})$ . An MCMC algorithm that uses  $\hat{L}(\theta|\mathbf{x})$  in place of  $L(\theta|\mathbf{x})$  is called pseudo-marginal MCMC (Beaumont 2003; Andrieu and Roberts 2009) and its stationary distribution is equal to the desired target posterior distribution.

The pseudo-marginal framework in our context requires an unbiased approximation for  $1/Z(\theta)$ . An unbiased approximation of  $Z(\theta)$  can be obtained using importance sampling estimate  $\hat{Z}(\theta)$  via MCMC samples from the likelihood. However,  $1/\hat{Z}(\theta)$  is a consistent but biased approximation for  $1/Z(\theta)$ . The Russian Roulette algorithm (Lyne et al. 2015) addresses this bias through a clever geometric series correction (details are in the supplement). This is an asymptotically exact algorithm and the theoretical assumptions are generally satisfied for general forms of the probability model  $h(\mathbf{x}|\theta)$ . However, it requires multiple  $\hat{Z}_i(\theta)$ 's at each iteration and each  $\hat{Z}_i(\theta)$  approximation itself requires multiple MCMC samples from  $h(\mathbf{x}|\theta)$ , making it computationally very expensive.

### 3.3. Noisy MCMC and Hybrids

If a Markov chain with transition kernel  $P$  satisfies detailed balance with respect to the target  $\pi$ , it is asymptotically exact. However, when  $P$  is approximated by  $\hat{P}$ , the samples generated may only approximately follow the target  $\hat{\pi}$ . Such algorithms may be generically referred to as “noisy MCMC.” Alquier et al. (2016) described a broad class of noisy MCMC algorithms and used total variational distance to quantify the distance between the asymptotically exact and inexact chain. Noisy MCMC in

its broader sense is a large class of algorithms, including, for instance, pseudo-marginal MCMC and ALR. Here, we discuss a noisy MCMC algorithm that builds upon the exchange algorithm; hence this algorithm may be thought of as a hybrid between auxiliary and likelihood approximation-based methods.

In the exchange algorithm, a single auxiliary variable  $\mathbf{y}$  is generated from  $h(\mathbf{y}|\theta')/Z(\theta')$ . Therefore, the  $h(\mathbf{y}|\theta)/h(\mathbf{y}|\theta')$  term in the (4) may be thought of as a one-sample unbiased importance sampling estimate of  $Z(\theta)/Z(\theta')$ . Instead of using a single  $\mathbf{y}$ , if multiple auxiliary variables  $\{\mathbf{y}_1, \dots, \mathbf{y}_N\}$  are generated from the  $h(\mathbf{y}|\theta')/Z(\theta')$ , the resulting importance sampling approximation will have smaller variance. This is called the noisy exchange algorithm with acceptance probability

$$\alpha = \min \left\{ 1, \frac{p(\theta')h(\mathbf{x}|\theta')q(\theta|\theta')}{p(\theta)h(\mathbf{x}|\theta)q(\theta'|\theta)} \frac{1}{N} \sum_{i=1}^N \frac{h(\mathbf{y}_i|\theta)}{h(\mathbf{y}_i|\theta')} \right\}. \quad (11)$$

For  $1 < N < \infty$ , the algorithm is asymptotically inexact because the detailed balance condition does not hold. Liang and Jin (2013) also proposed a similar approach called Monte Carlo MH (MCMH). Alquier et al. (2016) reported that the noisy exchange algorithm shows better mixing than the exchange algorithm, and proposed a noisy Metropolis-adjusted Langevin (MALA) exchange algorithm, where the gradient of the intractable distribution is approximated to construct proposals, that further improves upon mixing. Also, the estimate for the ergodic average has smaller bias than the exchange algorithm in empirical studies. Of course, there may be a tradeoff between improving mixing and increasing computational costs per iteration. We skip details but mention the general noisy MCMC approach as it may lead to other useful algorithms. For instance, DMH with multiple auxiliary variables, say “noisy DMH,” would be a simple extension.

*Components to be tuned:* The number of auxiliary variables,  $N$ , is an additional component to be tuned compared to the exchange algorithm or DMH. The choice of  $N$  depends on some tradeoffs: as  $N$  becomes large, the constructed chain can have better mixing and estimates from the chain can have lower variance at the expense of computing time. Parallel computing may be helpful for sampling  $N$   $\mathbf{y}_i$ 's, and evaluating  $h(\mathbf{y}_i|\theta)/h(\mathbf{y}_i|\theta')$  independently.

*Theoretical justification:* Alquier et al. (2016) provided an upper bound for total variation norm distance between  $P$ , transition kernel for the exchange algorithm, and  $\hat{P}$ , transition kernel for the noisy exchange algorithm, based on a result in (Mitrophanov 2005, Corollary 3.1), which requires uniform ergodicity of  $P$ . Two assumptions are used in this derivation: (1) the prior  $p(\theta)$  is bounded away from 0 and  $\infty$ , (2)  $q(\theta'|\theta)$  is bounded away from 0 and  $\infty$ ; both assumptions are typically satisfied.

## 4. Simulated and Real Data Examples

We now study the algorithms in the context of three models that are of general interest: (1) the Ising model, (2) a social network model, and (3) an attraction-repulsion point process model. These models also present different computational challenges. The code for this is implemented in R (Ihaka and Gentleman 1996) and C++, using the Rcpp and RcppArmadillo

packages (Eddelbuettel et al. 2011). We use the examples to compare the efficiency as well as practical implementation challenges of the algorithms. Simulation settings for the Russian roulette algorithm are provided in the supplementary material.

#### 4.1. The Ising Model

Consider an Ising model (Lenz 1920; Ising 1925) on an  $m$  by  $n$  lattice with parameter  $\theta$ . The observed lattice  $\mathbf{x} = \{x_{ij}\}$  has binary values  $x_{i,j} = \{-1, 1\}$ , where  $i, j$  denotes row and column location in the lattice. The model is

$$L(\theta|\mathbf{x}) = \frac{1}{Z(\theta)} \exp \{ \theta S(\mathbf{x}) \},$$

$$S(\mathbf{x}) = \sum_{i=1}^m \sum_{j=1}^{n-1} x_{i,j} x_{i,j+1} + \sum_{i=1}^{m-1} \sum_{j=1}^n x_{i,j} x_{i+1,j}, \quad (12)$$

where spatial dependence is imposed via  $S(\mathbf{x})$ . The larger  $\theta$  becomes, the stronger the interactions between the data on the lattice. Summation over all  $2^{mn}$  possible configurations of this model is required for the calculation of the normalizing function, which is computationally expensive even for lattices of moderate size. The simulations are conducted using perfect sampling (Propp and Wilson 1996) on a  $10 \times 10$  lattice with uniform prior  $[0,1]$  with parameter settings representing moderate dependence, with  $\theta = 0.2$ , and strong dependence, with  $\theta = 0.43$ .

All algorithms are tuned according to the methods described in the previous section.  $h(\mathbf{y}|\hat{\theta})/Z(\hat{\theta})$  is chosen as the conditional density of the auxiliary variable  $\mathbf{y}$  for the AVM algorithm, where  $\hat{\theta}$  is MPLE. The auxiliary variable is generated by 10 cycles of Gibbs updates in DMH. For AEX, 100 particles are selected among 5000 samples from fractional DMH (descriptions are in the supplementary material). Then the preliminary run of the auxiliary chain is implemented for 420,000 iterations with the first 20,000 iterations discarded for burn-in. To expedite computing and reduce memory management issues, the resulting samples are thinned (at equally spaced intervals of 20) to obtain 20,000 samples. In the auxiliary chain, we set  $n_0 = 20,000$ ,

$K_s = [0, 100^{100+10s}]^{100}$ , and  $\mathbf{x}_n$  is updated by a single cycle of Gibbs updates. For ALR, 100 particles are drawn from the uniform prior and  $\mathbf{x}_n$  is updated by a single cycle of Gibbs updates. The kernel giving equal weights for  $h = 10$  nearest particles and 0 for others is used. In noisy DMH, 100 samples are used to produce importance sampling estimate and same inner sampler is constructed as DMH. We note that to make the computations feasible, we used parallel computing to obtain importance sampling estimates for both noisy DMH and the Russian Roulette algorithms. The parallel computing was implemented through OpenMP with the samples generated in parallel across eight processors. We treated a run from the exchange algorithm as our gold standard; it was run for 1,010,000 iterations with 10,000 discarded for burn-in and 10,000 thinned samples are obtained from the remaining 1,000,000. Same simulation settings are used for  $\theta = 0.43$  case. Since perfect sampling takes very long for  $\theta = 0.43$  case, we use the ALR algorithm as the gold standard. All algorithms were run until the Monte Carlo standard errors calculated by batch means (Jones et al. 2006; Flegal, Haran, and Jones 2008) are below 0.01.

Table 1 shows that the estimates from different algorithms are well matched to the gold standard when  $\theta = 0.2$ . Posterior densities in Figure 6 also indicate agreement. Likelihood approximation approaches show larger ESS than auxiliary variable approaches. On the other hand, auxiliary variable approaches have smaller computational costs than likelihood approximation approaches. In particular, DMH shows the shortest computing time. For this parameter setting, inference results from both asymptotically exact and inexact algorithms are accurate.

However when  $\theta = 0.43$ , inference results are biased for asymptotically inexact algorithms. In Table 1 and Figure 6, it is observed that the inference results from DMH and Noisy DMH do not match the gold standard. Due to the strong dependence at this parameter setting, mixing of the inner sampler to generate samples from the  $h(\mathbf{x}|\theta)/Z(\theta)$  is slower than for the  $\theta = 0.2$  case. Therefore, a large number of Gibbs updates are necessary for accurate inference. Although Noisy DMH provides a closer estimate to the true value than DMH because it uses multiple samples for estimating  $Z(\theta)/Z(\theta')$ , it is still biased. The number

**Table 1.** Inference results for an Ising model on a  $10 \times 10$  lattice. 20,000 MCMC samples are generated from each algorithm. The highest posterior density (HPD) is calculated by using coda package in R. The calculation of effective sample size (ESS) follows Kass et al. (1998) and Robert and Casella (2013). "Acc" represents acceptance probability.

$\theta = 0.2$	Mean	95%HPD	ESS	Acc	Time (sec)	ESS/Time
AVM	0.20	(0.08,0.32)	1527.88	0.37	12.46	122.62
Exchange	0.20	(0.08,0.33)	2061.23	0.49	15.29	134.81
DMH	0.21	(0.08,0.34)	2068.29	0.48	1.78	1161.96
AEX	0.21	(0.07,0.35)	1778.15	0.47	23.72	74.96
ALR	0.20	(0.08,0.34)	3647.34	0.59	9.33	390.93
RussianR	0.20	(0.06,0.33)	2650.83	0.50	13849.10	0.19
NoisyDMH	0.20	(0.06,0.33)	3347.76	0.59	71.02	47.14
Gold	0.20	(0.08,0.33)	9845.89			
$\theta = 0.43$	Mean	95%HPD	ESS	Acc	Time (sec)	ESS/Time
AVM	0.44	(0.35,0.57)	1020.12	0.32	10395.91	0.10
Exchange	0.43	(0.33,0.54)	2146.05	0.40	4889.32	0.44
DMH	0.47	(0.34,0.61)	2029.24	0.48	1.72	1179.79
AEX	0.43	(0.32,0.53)	2048.01	0.41	24.96	82.05
ALR	0.43	(0.32,0.53)	4125.59	0.51	9.88	417.57
RussianR	0.44	(0.33,0.54)	2708.55	0.45	30247.70	0.09
NoisyDMH	0.47	(0.32,0.61)	3486.99	0.62	75.17	46.39
Gold	0.43	(0.33,0.54)	10000.00			

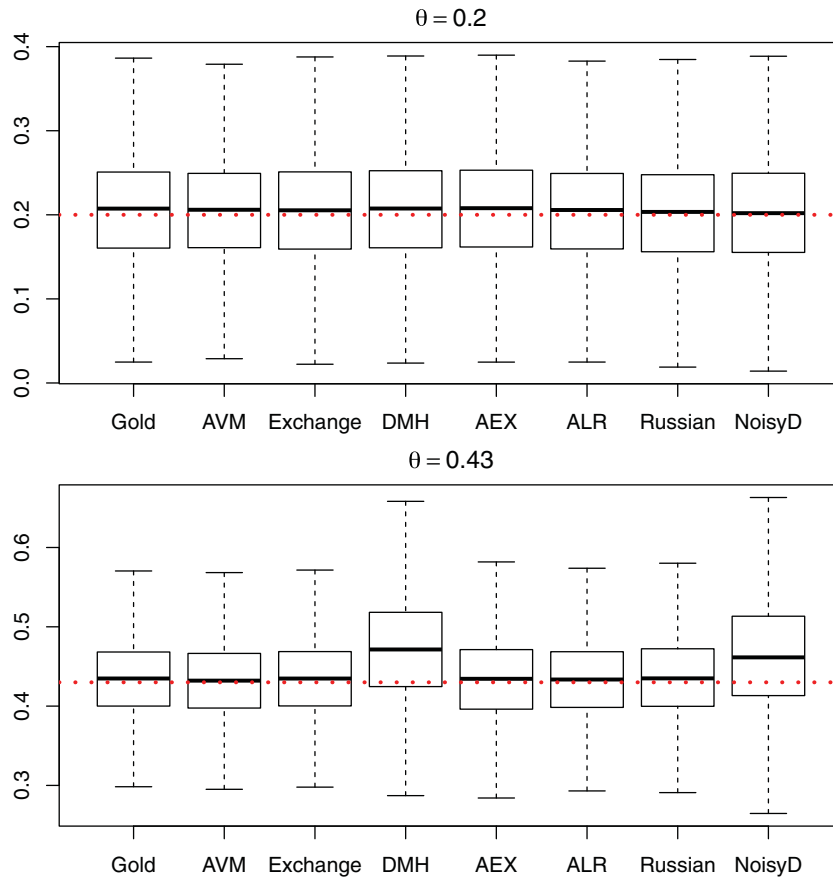


Figure 6. Posterior densities for different  $\theta$ . Dotted line indicates true value.

of iterations for the inner sampler in the Russian Roulette algorithm also has to be increased due to the slow mixing (see supplementary material for details). We can also observe that AVM and the exchange algorithm take several hours, because perfect sampling takes longer to achieve coalescence. On the other hand, there are no extra costs for AEX and ALR in the presence of strong dependence. For both algorithms, still a single cycle of Gibbs update is enough for updating  $\mathbf{x}_n$ . This is because both algorithms require updating not sampling  $\mathbf{x}_n$  from the probability model.

**Summary:** In the case where the Ising model has only moderate dependence, all the algorithms work reasonably well. Here, DMH may be preferable because it is the easiest to construct and the fastest to run. On the other hand, in the context of Ising models that exhibit strong dependence, the computational problem becomes more challenging. For the asymptotically inexact algorithms like DMH, to obtain accurate results in the presence of strong dependence in the Ising model, the number of iterations for the inner sampler needs to be very large. This results in a much larger computational burden. AVM and the exchange algorithms also become impractical because perfect sampling takes too long. AEX and ALR can be useful for spatial autologistic models because both algorithms can attain accurate estimates with moderate computing speed even in the presence of strong dependence.

## 4.2. Social Network Models

Exponential random graph models (ERGM; Robins et al. 2007; Hunter et al. 2008) provide an approach for modeling

relationships among nodes of a network. Consider the undirected ERGM with  $n$  nodes. For all  $i \neq j$ ,  $x_{i,j} = 1$  if the  $i$ th node and  $j$ th node are connected, otherwise  $x_{i,j} = 0$  and  $x_{i,i}$  is defined as 0. This forms an  $n$  by  $n$  adjacency matrix,  $\mathbf{x}$ . Calculation of the normalizing function requires summation over all  $2^{n(n-1)/2}$  configuration, which is infeasible. Here, we investigate two ERGM examples, the first with four parameters and the second with nine parameters; the latter example also involves a more complicated summary statistic which poses an additional computational challenge. We have not implemented AVM and the exchange algorithms because perfect samplers are possible only for a few special cases (Butts 2012).

### 4.2.1. A Florentine Business Network

In this example, we study the Florentine business dataset (Padgett 1994), which describes the business networks among 16 Florentine families. Consider the undirected ERGM, where the probability model is

$$L(\theta|\mathbf{x}) = \frac{1}{Z(\theta)} \exp \{ \theta_1 S_1(\mathbf{x}) + \theta_2 S_2(\mathbf{x}) + \theta_3 S_3(\mathbf{x}) + \theta_4 S_4(\mathbf{x}) \} \quad (13)$$

$$S_l(\mathbf{x}) = \sum_{i=1}^n \binom{x_{i+}}{l}, l = 1, 2, 3; \quad S_4(\mathbf{x}) = \sum_{i < j < k} x_{i,j} x_{j,k} x_{k,i}.$$

Sufficient statistics  $S(\mathbf{x}) = \{S_1(\mathbf{x}), S_2(\mathbf{x}), S_3(\mathbf{x}), S_4(\mathbf{x})\}$  represent the number of edges, two-stars, three-stars, and triangles, respectively, where  $x_{i+}$  indicates row sum of adjacency matrix. Triangle represents the number of cyclic relationships



and  $k$ -star indicates the number of nodes which have exactly  $k$  relationships.

Each of the algorithms is tuned according to the guidelines in the previous sections. The auxiliary variable is generated by 10 cycles of Gibbs updates in DMH. In AEX, 200 particles are selected among 5000 samples from fractional DMH, as described in the supplementary material. Then the preliminary run of the auxiliary chain is implemented for 630,000 iterations with the first 30,000 iterations discarded for burn-in. The resulting samples are thinned (at equally spaced intervals of 20) to obtain 30,000 samples. In the auxiliary chain, we set  $n_0 = 20,000$ ,  $K_s = [0, 100^{100+10s}]^{200}$ , and  $\mathbf{x}_n$  is updated by a single cycle of Gibbs updates. For ALR, 400 particles are chosen from the short run of DMH, and a single cycle of Gibbs updates are used to update  $\mathbf{x}_n$ . The same kernel with  $h = 20$  is used as previous example. One hundred samples are used for the importance sampling estimate in noisy DMH. Importance sampling estimates in Russian Roulette and noisy DMH algorithms are obtained in parallel as in the previous example. To obtain a gold standard for comparisons, we run the AEX algorithm 10 times independently. Each run consists of 101,000 iterations with 1000 iterations discarded as burn-in. Then 10,000 samples are obtained by thinning from the remaining 10 sets of 100,000 samples. All algorithms were run until the Monte Carlo standard error is at or below 0.02.

Here, we only provide the inference results regarding  $\theta_2$  because similar results are observed for the other parameters. Table 2 indicates that the estimates from the different algorithms are similar to those of the gold standard. This is because 10 cycles of Gibbs updates are enough to generate auxiliary samples from the likelihood in this example. However, for fewer iterations, say 1 to 2 cycles, the asymptotically inexact algorithms are biased. As in the Ising model example, less correlated samples can be generated from likelihood approximation approaches, at additional computational expense. Both ALR and AEX are computationally efficient compare to the Russian Roulette algorithm. This is because we can effectively store the previous samples using sufficient statistics in the likelihood. However, ALR is relatively expensive here (unlike in the Ising model), because it takes longer to visit the entire state space equally with increasing dimensions. The performance of AEX is relatively robust in the multidimensional case and is the fastest among asymptotically exact algorithms if the particles are carefully chosen.

**Summary:** AEX is the most reliable approach because it is asymptotically exact while at the same time retaining some computational efficiency because this model has low-dimensional

statistics. On the other hand, DMH is simple and computationally efficient and as long as the length of the inner sampler is reasonable, it also provides accurate inference.

#### 4.2.2. An Emergent Multi-Organizational Network (EMON)

In this example, we study the Mount St. Helens emergent multi-organizational network (EMON) dataset (Drabek et al. 1981), which describes communication networks for the search and rescue operations among 27 organizations. There are three node attributes in the model: (1) sponsorship level (city, county, federal, private, and state), (2) command rank score (higher score indicates higher rank), (3) whether headquarters are sited locally to the impact area or not (L/NL). Consider the undirected ERGM, where the probability model is

$$L(\theta|\mathbf{x}) = \frac{1}{Z(\theta)} \exp \left\{ \sum_{l=1}^9 \theta_l S_l(\mathbf{x}) \right\}, \quad (14)$$

$$S_1(\mathbf{x}) = \sum_{i=1}^n \binom{x_{i+}}{1}, \quad S_2(\mathbf{x}) = \sum_{i < j < k} x_{i,j} x_{j,k} x_{k,i}$$

$$S_3(\mathbf{x}) = \sum_{i < j} x_{i,j} (1\{\text{sponsor}_i = \text{county}\} + 1\{\text{sponsor}_j = \text{county}\})$$

$$S_4(\mathbf{x}) = \sum_{i < j} x_{i,j} (1\{\text{sponsor}_i = \text{federal}\} + 1\{\text{sponsor}_j = \text{federal}\})$$

$$S_5(\mathbf{x}) = \sum_{i < j} x_{i,j} (1\{\text{sponsor}_i = \text{private}\} + 1\{\text{sponsor}_j = \text{private}\})$$

$$S_6(\mathbf{x}) = \sum_{i < j} x_{i,j} (1\{\text{sponsor}_i = \text{state}\} + 1\{\text{sponsor}_j = \text{state}\})$$

$$S_7(\mathbf{x}) = \sum_{i < j} x_{i,j} (\text{command}_i + \text{command}_j)$$

$$S_8(\mathbf{x}) = \sum_{x_{i,j} \in L} GD_4, \quad S_9(\mathbf{x}) = \sum_{x_{i,j} \in NL} GD_4.$$

The sufficient statistics are  $S_1(\mathbf{x})$  (the number of edges),  $S_2(\mathbf{x})$  (triangles),  $S_3(\mathbf{x}) - S_6(\mathbf{x})$  (node factor for sponsorship level),  $S_7(\mathbf{x})$  (node covariance for command rank score), and  $S_8(\mathbf{x}) - S_9(\mathbf{x})$  (graphlet orbit factor for location of headquarters). Graphlet statistics are small, connected subgraphs which represent the certain topological structure of a network (Pržulj, Corneil, and Jurisica 2004; Pržulj 2007). Here, we used Graphlet 4 with automorphism orbit 7 for describing brokerage roles between the 27 organizations (Yaveroglu et al. 2014).

In this model, there are practical implementation issues for both particle-based algorithms (ALR and AEX): (1) Both algorithms require increasing number of particles to cover the high-dimensional parameter space. (2) The algorithms used for generating particles can exhibit slow mixing. For example, the default implementation of DMH mixes slowly, which implies that it takes longer chains to generate the necessary particles for the ALR or AEX algorithms. (3) Even after the particles are generated, the stochastic approximation algorithm can still be slow when the number of parameters is large. Even with large gain factor (AEX) or step size (ALR), visiting each state space equally is infeasible when the parameter dimension is high (as in the nine-dimensional case above). This issue is perhaps even more problematic than issues (1) and (2) above.

The Russian Roulette algorithm is also computationally expensive because of the complexity of summary statistics in this

**Table 2.** Inference results for 2-star in ERGM for Florentine business dataset. 30,000 MCMC samples are generated from each algorithm. The highest posterior density (HPD) is calculated by using coda package in R. The calculation of effective sample size (ESS) follows Kass et al. (1998) and Robert and Casella (2013). "Acc" represents acceptance probability.

$\theta_2$	Mean	95%HPD	ESS	Acc	Time (sec)	ESS/Time
DMH	1.24	(0.02,2.57)	1026.67	0.24	10.45	98.25
AEX	1.24	(0.17,2.58)	991.87	0.20	126.46	7.84
ALR	1.25	(0.16,2.52)	1456.57	0.33	2500.37	0.58
RussianR	1.27	(0.03,2.68)	1433.60	0.31	33534.96	0.04
NoisyDMH	1.28	(0.03,2.59)	1600.90	0.32	297.35	5.38
Gold	1.27	(0.08,2.50)	9655.90			

**Table 3.** Inference results for graphlet orbit factor for “NL” in EMON dataset. 40,000 MCMC samples are generated from each algorithm. The highest posterior density (HPD) is calculated by using the coda package in R. The calculation of effective sample size (ESS) follows Kass et al. (1998) and Robert and Casella (2013). “Acc” represents acceptance probability.

$\theta_9 \times 10^2$	Mean	95%HPD	ESS	Acc	Time (min)	ESS/Time
DMH	1.93	(0.61,2.98)	583.08	0.05	16.10	36.21
NoisyDMH	1.93	(0.08,3.45)	796.01	0.08	1026.15	0.78
Gold	1.89	(0.58,2.98)	5494.31			

example (wall time of roughly 1 month). Therefore, we study only the DMH and noisy DMH algorithms for this example. The details of our implementation are as follows. We generate auxiliary variables via 10 cycles of Gibbs updates, and 300 samples are used to construct an importance sampling estimate in noisy DMH. Importance sampling approximations in noisy DMH are evaluated through parallel methods, as in previous examples. We use the DMH algorithm with 20 cycles as the gold standard; it was run for 101,000 iterations with 1000 samples discarded for burn-in and 10,000 thinned samples are obtained from the remaining samples. All algorithms are run until the Monte Carlo standard error is no larger than 0.03.

Here, we provide inference results for  $\theta_9$  because similar results are observed for the other parameters. Table 3 shows that posterior means from the different algorithms are well matched to the gold standard. It is observed that the acceptance rate is much lower than in the simpler Florentine business data example. The highest posterior density (HPD) interval obtained from Noisy DMH is slightly wider than that of the gold standard. As in the previous examples, a Noisy DMH algorithm can generate less correlated samples with higher acceptance rates compared to the DMH algorithm.

**Summary:** Particle-based algorithms (AEX and ALR) are infeasible for the nine-dimensional parameter case even though there are summary statistics to be stored. Because of computationally expensive graphlet statistics, Russian Roulette is infeasible. With the choice of an appropriate length for the inner sampler, DMH can provide accurate inference, with the highest ESS/T. This fact suggests that for these challenging cases, DMH may still be a practical approach.

### 4.3. Spatial Interaction Point Process

A spatial point process  $\mathbf{x} = \{x_1, \dots, x_n\}$  is a realization of random points in a bounded plane  $S \subset R^2$ . By introducing an interaction function  $\phi(D_{ij})$  which is the function of distance between the coordinates of  $x_i$  and  $x_j$ , a probability model may be used to describe spatial patterns among the point. Extending the Strauss process (Strauss 1975) which explains repulsion patterns among point, Goldstein et al. (2015) developed a point process model to explain both attraction and repulsion patterns of the cells infected with human respiratory syncytial virus (RSV). The interaction function is

$$\phi(D) = \begin{cases} 0 & 0 \leq D \leq R \\ \theta_1 - \left( \frac{\sqrt{\theta_1}}{\theta_2 - R} (D - \theta_2) \right)^2 & R < D \leq D_1 \\ 1 + \frac{1}{(\theta_3(D - D_2))^2} & D > D_1 \end{cases} \quad (15)$$

and the probability model is

$$L(\theta|\mathbf{x}) = \frac{\lambda^n \left[ \prod_{i=1}^n \exp \left\{ \min \left( \sum_{i \neq j} \log(\phi(D_{i,j})), 1.2 \right) \right\} \right]}{Z(\theta)}, \quad (16)$$

$$\theta = \{\lambda, \theta_1, \theta_2, \theta_3\}.$$

There are four parameters  $\{\lambda, \theta_1, \theta_2, \theta_3\}$  in the model:  $\lambda$  controls intensity of the point process and  $\{\theta_1, \theta_2, \theta_3\}$  are the parameters controlling the interaction function.  $\theta_1$  is the peak value of  $\phi$ ,  $\theta_2$  is value of  $D$  at the peak of  $\phi$ , and  $\theta_3$  represents descent rate after the peak. When the points are too close to each other,  $\phi$  value in (15) is less than 1 which means points have a tendency to remain apart. However as the distance between points is larger,  $\phi$  value becomes increased which means that points clump together. Attraction patterns become smaller as the distance between the points becomes larger. Likewise, this model can capture attraction repulsion spatial association among infected cells. Calculation of the normalizing function requires integration over the continuous domain  $S$ , which is intractable.

Goldstein et al. (2015) implemented DMH for three independent replicates of 3200 points in a well with radius 1350 pixels, which is about 10,000 points. Even with code written in C, inference took roughly 12 hr. This is because the number of points  $n$  is large, and for each iteration of the DMH, thousands of birth-death MCMC steps are required to generate the auxiliary variable. To allow for a comparison with other algorithms, which are computationally more expensive than DMH, we work with a smaller point pattern; however, this pattern is still computationally challenging enough to serve as a good testbed for the various algorithms we consider. Simulations are conducted on a well with 337.5 radius pixels and number of points  $n \approx 200$  without replicates. Point process  $\mathbf{x}$  is generated through a long run of birth-death MCMC. We follow RSV-B simulation settings in Goldstein et al. (2015), where the true parameter is  $\{\lambda \times 10^4, \theta_1, \theta_2, \theta_3\} = \{4, 1.2, 15, 0.3\}$ . We use uniform priors on  $\theta$ . Since the number of data points is small, descent rate parameter  $\theta_3$  is not recovered well. Therefore,  $\theta_3$  is fixed at 0.3 and we infer the other three parameters.

Implementing AVM or the exchange algorithm is infeasible because perfect sampling is impossible for this example. Although both AEX and ALR can be implemented theoretically, it is not practical because there are no summary statistics. Without summary statistics, we need to store the distance matrix of cumulative point process samples with varying dimension around 200 by 200, which is a burden on memory. Therefore, we study DMH, noisy DMH, and the Russian Roulette algorithms for this example. Each algorithm is tuned according to the previous sections. For all the algorithms, samples are generated from the likelihood through 2000 iterations of birth-death MCMC. For noisy DMH, 100 samples are used to construct an importance sampling approximation with each iteration. Importance sampling approximations in noisy DMH and Russian Roulette are evaluated through parallel methods, as in the past. We use the Russian Roulette algorithm as the gold standard; it was run for 101,000 iterations with 1000 samples discarded for burn-in and 10,000 thinned samples are obtained from the remaining samples. All algorithms were run until the Monte Carlo standard error is at or below 0.01.

**Table 4.** Inference results for outputs for  $\theta_1$  in attraction repulsion point process model ( $n \approx 200$ ). 40,000 MCMC samples are generated from each algorithm. The highest posterior density (HPD) is calculated by using coda package in R. The calculation of effective sample size (ESS) follows Kass et al. (1998) and Robert and Casella (2013). "Acc" represents acceptance probability.

$\theta_1$	Mean	95%HPD	ESS	Acc	Time (min)	ESS/Time
DMH	1.20	(1.03,1.33)	1343.11	0.27	3.93	341.35
RussianR	1.19	(1.04,1.35)	1867.54	0.32	3299.76	0.57
NoisyDMH	1.20	(1.04,1.34)	2676.49	0.42	77.73	34.44
Gold	1.19	(1.03,1.34)	4397.39			

Here, we only provide results for  $\theta_1$  because similar results are observed for the other parameters. Table 4 shows that inference results from the different algorithms are well matched to the gold standard and the highest posterior density (HPD) intervals cover true values. As in the previous examples, likelihood approximation approaches can generate less correlated samples with higher acceptance rates.

**Summary:** DMH has significant advantages in terms of computational efficiency. Compared to the Russian Roulette algorithm which takes about 55 hr, DMH only takes several minutes to be implemented. This is because both calculation of the  $h(\mathbf{x}|\theta)$  and the inner sampler to generate auxiliary variables are expensive compared to previous examples. Furthermore, considering ESS/T, DMH can generate effective sample within the shortest time, while simulated samples are close to the gold standard. This fact demonstrates that for computationally expensive problem, especially the model without low-dimensional sufficient statistics, we recommend DMH as a practical approach, though asymptotically exact inference is not guaranteed theoretically. Considering that the example is 16 times smaller than the original, it is infeasible to implement other algorithms for the original problem.

## 5. Computational Complexity

Here, we investigate the computational complexity of the algorithms, that is, how the algorithms scale as we increase the number of data points. We believe this is particularly relevant in modern statistics since a common and important question to ask is how well an algorithm works as data sizes get large. Algorithms requiring perfect sampling will not be considered, because perfect sampling can only be constructed for limited cases. Let  $L(\cdot)$  be the complexity of evaluating  $h(\mathbf{x}|\theta)$ ,  $G(\cdot)$  be the complexity of inner sampler, and  $n$  be the number of points in data  $\mathbf{x}$ .

We begin with a few caveats. It is challenging to calculate complexity because some algorithms contain random quantities (stopping times), or have many components that need to be tuned specifically for different problems. Also, the mixing of the Markov chains and the quality of asymptotically exact algorithms are influenced by the data itself (for instance, dependence) and not just the size of the data. We also note that to simplify calculations, we assume the dimensions of the data  $\mathbf{x}$  and the auxiliary variable  $\mathbf{y}$  are assumed to be same ( $n$ ). This is not strictly accurate for the point process example where the auxiliary variable is generated through birth–death MCMC with varying dimensions. However, the dimensions are approximately similar to original data's dimension  $n$ . In the point

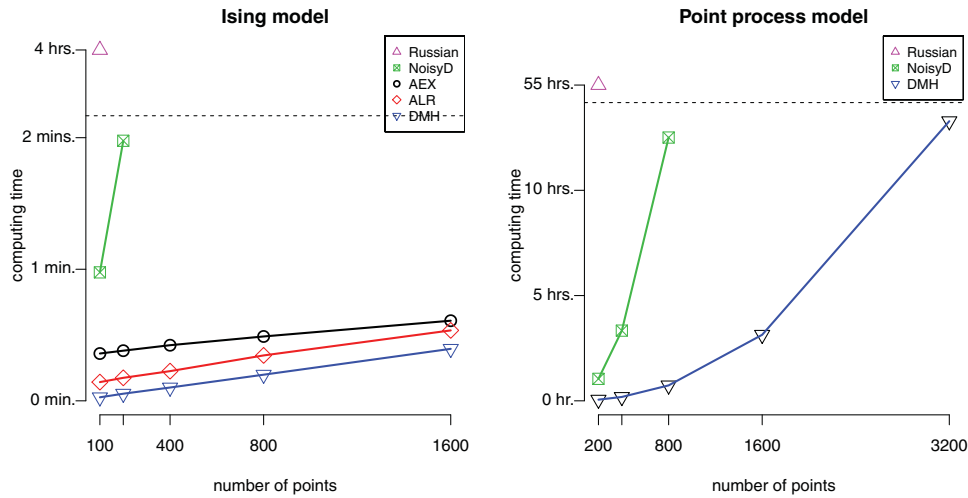
process example, we can store data  $\mathbf{x}$  itself or distance matrix of  $\mathbf{x}$  in the adaptive algorithms (AEX, ALR) to evaluate  $h(\mathbf{x}|\theta)$ . The latter can avoid recomputations at the expense of memory. Here, we assume that distance matrix of  $\mathbf{x}$  will be stored. We assume that the length of the inner sampler is proportional to  $n$ . It is true when the inner sampler is used for update (AEX, ALR) which requires a single cycle of update (proportional to  $n$ ). However, it may not be correct if the purpose of the inner sampler is sampling (DMH, Russian Roulette, Noisy DMH) from the probability model. The mixing of the inner sampler can be different depending on the parameter settings. However, several cycles of update (proportional to  $n$ ) appears to work in practice.

All the algorithms require sampling (DMH, Noisy DMH, Russian Roulette) or updating (AEX, ALR) from the probability model and evaluating  $h(\mathbf{x}|\theta)$ . In exponential family models such as Ising model or ERGM, complexity of  $L(\cdot)$  is of order  $n$ , because the number of calculations for  $S(\mathbf{x})$  is proportional to  $n$ . For the inner sampler, when a single point is proposed to be changed, only neighboring points (fixed number regardless of  $n$ ) are affected due to the Markovian assumption. Because the length of the inner sampler is assumed to increase with  $n$ ,  $G(\cdot)$  is also of order  $n$ . On the other hand, complexity of  $L(\cdot)$  for point process model is of order  $n^2$ , because evaluation of the  $h(\mathbf{x}|\theta)$  requires calculating an  $n$  by  $n$  distance matrix for  $n$  data points, and evaluating interaction function  $\phi(\cdot)$  on the corresponding distance matrix. For the inner sampler, when a single point is proposed to be added (birth) or deleted (death), the distance of a proposed point from other  $n$  points should be calculated, and then  $\phi(\cdot)$  should be evaluated at each point. Since the length of the inner sampler is assumed to be proportional to  $n$ ,  $G(\cdot)$  is of order  $n^2$ .

There is a big difference in calculating  $h(\mathbf{x}|\theta)$  for exponential family and point process model. For the exponential family model, once we evaluate  $S(\mathbf{x})$ , we can simply take the product of  $\theta$  and  $S(\mathbf{x})$  for evaluation of  $h(\mathbf{x}|\theta)$  in different  $\theta$ . On the other hand, for point process model,  $h(\mathbf{x}|\theta)$  should be recalculated;  $\phi(\cdot)$  should be evaluated at the distance matrix of  $\mathbf{x}$  with different parameters. Here, we provide our main observations (see supplement for details): (1) Complexity of exponential family is  $\mathcal{O}(n)$  and  $\mathcal{O}(n^2)$  in point process for all the algorithms. (2) Although algorithms have the same complexity, there are major differences in calculations per iteration. The amount of calculations per iteration of AEX, ALR, and DMH are similar for exponential family models. Considering complexity and memory costs, DMH is the cheapest algorithm for both models. (3) Adaptive algorithms (AEX and ALR) require more memory. With increasing iterations, both algorithms become slower because algorithms use cumulative samples  $H_{n+1}$  in calculations per iteration (AEX: Step 5 of the Algorithm 5, ALR: Step 4 of the Algorithm 6). Memory may not be an issue for the exponential family model. However, for models without low-dimensional summary statistics, memory costs can become prohibitively expensive.

Figure 7 is the observed computing time for several algorithms with different scales in both models. We only include parts of results for Noisy DMH and the Russian Roulette algorithm. This is because both algorithms are too expensive to compare with other algorithms for large scales. For different scales, mixing of each algorithm is determined to be similar based on effective sample sizes, once we use appropriate step





**Figure 7.** Illustration of the observed computing time for algorithms. For Ising model, time is measured for  $\theta = 0.2$  with 20,000 iterations. For point process model, time is measured for RSV-B simulation settings in Goldstein et al. (2015) with 40,000 iterations.

size (covariance) for proposal. Step size can be tuned to achieve similar acceptance rate for different scales; the larger the data size becomes, the smaller step is used. We can also adaptively update step size (Atchadé 2006; Atchade, Lartillot, and Robert 2008). Figure 7 supports our calculations about  $\mathcal{O}(n)$  complexity for exponential family and  $\mathcal{O}(n^2)$  complexity for point process model. Also in the exponential family model, slopes of AEX, ALR, and DMH are similar to each other. However, DMH is the fastest because of memory requirements of both adaptive algorithms. These facts are consistent with our calculations.

## 6. Connections and Summary of Results

Here, we point out connections between the algorithms. Also, based on our study, we provide some conclusions about advantages and disadvantages of each algorithm.

### 6.1. Connections and Observations

All the algorithms require sampling from the probability model either approximately (MCMC or particle methods) or exactly (perfect sampling). These samples, in turn, are used in some form of an importance sampling estimate. This is clear in the likelihood approximation approaches; for the auxiliary variable approaches one may think of the acceptance probability of the Metropolis–Hastings algorithm as containing a single sample importance sampling approximation of  $Z(\theta)/Z(\theta')$ . Let the conditional density of the auxiliary variable in AVM be  $f(\mathbf{y}|\theta, \mathbf{x}) = h(\mathbf{y}|\hat{\theta})/Z(\hat{\theta})$ , where  $\hat{\theta}$  is the MPLE. The part of the acceptance probability only related to the auxiliary variable is

$$\frac{f(\mathbf{y}'|\theta', \mathbf{x})h(\mathbf{y}|\theta)}{h(\mathbf{y}'|\theta')f(\mathbf{y}|\theta, \mathbf{x})} = \frac{h(\mathbf{y}'|\hat{\theta})h(\mathbf{y}|\theta)}{h(\mathbf{y}'|\theta')h(\mathbf{y}|\hat{\theta})}. \quad (17)$$

Murray, Ghahramani, and MacKay (2006) pointed out that since  $\mathbf{y} \sim h(\cdot|\theta)/Z(\theta)$  and  $\mathbf{y}' \sim h(\cdot|\theta')/Z(\theta')$ ,  $h(\mathbf{y}'|\hat{\theta})/h(\mathbf{y}'|\theta')$  and  $h(\mathbf{y}|\hat{\theta})/h(\mathbf{y}|\theta)$  may be thought of as one-sample unbiased importance sampling approximations of  $Z(\hat{\theta})/Z(\theta')$  and  $Z(\hat{\theta})/Z(\theta)$ , respectively. Therefore, (17) is the ratio of two unbiased estimates which is a biased estimate of  $Z(\theta)/Z(\theta')$ .

Murray, Ghahramani, and MacKay (2006) explained that compared to AVM, the exchange algorithm is more direct because  $h(\mathbf{y}|\theta)/h(\mathbf{y}|\theta')$  in (4) is one-sample unbiased importance sampling estimate for  $Z(\theta)/Z(\theta')$  where  $\mathbf{y} \sim h(\cdot|\theta')/Z(\theta')$ . Though sampling schemes for the auxiliary variable are different, the same logic is applied to DMH and AEX. The only difference is that DMH generates  $\mathbf{y}$  via MCMC, and AEX generates  $\mathbf{y}$  via resampling (dynamic importance sampling from the mixture distribution). Both classes of algorithms are clearly connected through importance sampling. We summarize these connections in Figure 8.

AEX lies at the intersection of both classes of algorithms. As in the likelihood approximation approach, in the auxiliary chain  $w_n^{(i)}$  approximates (up to a constant)  $Z(\theta^{(i)})$  for each  $\theta^{(i)}$ . As in the auxiliary variable methods, the intractable functions are cancelled in the target chain at each iteration. Furthermore, AEX is closely connected to ALR. Both algorithms approximate  $Z(\theta^{(i)})$  at each particle— $w_n^{(i)}$  in AEX and  $c_n^{(i)}$  in ALR. Using their respective approximations, both collect a sample from a family of distributions  $\{h(\mathbf{x}|\theta^{(1)})/Z(\theta^{(1)}), \dots, h(\mathbf{x}|\theta^{(d)})/Z(\theta^{(d)})\}$  via (different) stochastic approximations and save a sample with each iteration. As the number of iterations grows, the accumulated dataset  $H_{n+1}$  grows, which makes both algorithm asymptotically exact. In AEX, the resampling distribution of  $\mathbf{y}$  becomes close to  $h(\mathbf{y}|\theta')/Z(\theta')$  and in ALR the approximation  $\hat{Z}_{n+1}(\theta)$  converges to the truth. The difference is that AEX cancels out  $Z(\theta)$  while ALR plugs-in  $\hat{Z}_{n+1}(\theta)$  into the acceptance probability.

In general, likelihood approximation approaches have better mixing than auxiliary variable approaches. However, auxiliary variable approaches are less expensive per iteration resulting in higher effective sample size per second. Hence, auxiliary variable approaches tend to often be faster than likelihood approximation approaches. However, the efficiency of algorithms can change depending on the model, parameter settings, and the data, as well as decisions about tuning components in each algorithm. In particular, strong dependence in observations tends to be slow sample generation, both with approximate and exact samplers. Perfect sampling takes longer to achieve coalescence, which can slow the AVM and exchange algorithm. Also, MCMC sampling requires longer chains, increasing computing time for



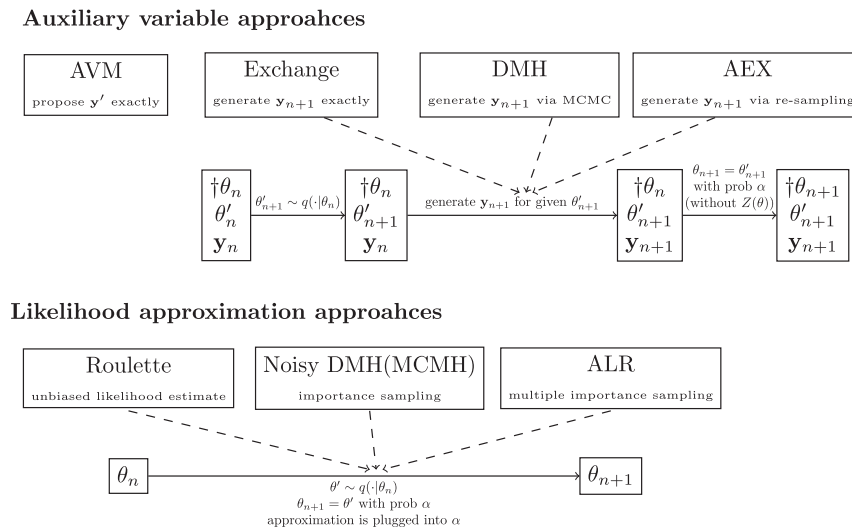


Figure 8. Connections between the algorithms. The dagger symbols indicate the parameter of interest in auxiliary variable approaches.

the DMH and the Russian Roulette algorithms as well. In this case, AEX and ALR show relatively robust performance because both algorithms require updating  $\mathbf{x}_n$ , not sampling  $\mathbf{x}_n$  from the probability model. For both algorithms, the choice of particles is crucial in determining both statistical and computational efficiency. If the parameter has higher dimensions, fewer particles are likely to be in the higher density regions of the posterior distribution of  $\theta$ . Then the approximation of  $Z(\theta)$  might be inaccurate (ALR) or resampled auxiliary variable might be improbable (AEX), which can lead to poor mixing. Requiring an increase in the number of particles also slows computing. Therefore, careful choice of particles is important (see the supplementary material), especially for multidimensional parameters.

## 6.2. Guidelines and Recommendations

Comparing asymptotically exact and inexact algorithms is challenging in general because we have to simultaneously account for two different aspects of the algorithms: (1) mixing of the Markov chain, which affects the rate at which sample-based approximations converge to the true values, and (2) the quality of the target approximation in the case of asymptotically inexact algorithms. To provide reasonable comparisons, we always check the final results in each case using a variety of diagnostics, for example, plots of marginal distributions as the Monte Carlo sample size increases, to convince ourselves that the approximation we obtain finally is close enough to the truth so that the above issues are avoided. Once this is the case, we can compare the efficiency of the algorithms via effective sample size (ESS) and effective samples per unit time (ESS/T). We have used ESS as a practical criteria of comparing mixing of the different algorithms; ESS/T can measure the computational efficiency as well.

Table 5 summarizes algorithms based on the following criteria:

- Class: which class the each algorithm falls in. A,L indicates auxiliary variable approaches and likelihood approximation approaches, respectively, and A/L indicates both.
- Exact or inexact: whether the stationary distribution of the Markov chain is exactly equal to the target posterior in the long run or not.

Table 5. Comparison between algorithms.

Method	Class	Exact or inexact	Adaptation	Requirement	Ease of use
AVM	A	Exact	No	Perfect sampling	3
Exchange	A	Exact	No	Perfect sampling	2
DMH	A	Inexact	No	Nothing	1
AEX	A/L	Exact	Yes	Sufficient statistics	4
ALR	L	Exact	Yes	Sufficient statistics	4
Russian Roulette	L	Exact	No	Nothing	4
Noisy DMH	L	Inexact	No	Nothing	2

- Adaptation: whether the algorithm is adaptive (stationary distribution changes) or not.
- Requirement: key requirement(s) for the algorithm to be practical.
- Ease of use (rank): how much users need to make decision to run the algorithm. The smallest number indicates the simplest one.

AVM and the exchange algorithm propose the novel idea of cancelling out  $Z(\theta)$  and thus have inspired many other algorithms. However, it is difficult to implement both algorithms in practice because the algorithms require perfect sampling. Even if we can implement perfect sampling it still has limited applicability because perfect sampling is very slow, either in the presence of strong dependence or for large-scale data.

AEX, ALR, and Russian Roulette are asymptotically exact without perfect sampling. AEX is the fastest, especially when there are low-dimensional sufficient statistics for the model. AEX is faster than ALR for multidimensional  $\theta$  cases, because AEX does not have random stopping time ( $\tau$ ). For ALR, it is observed that  $\tau$  in Algorithm 6 can be large for multidimensional  $\theta$  problems, because it becomes difficult to move around the state-space equally. This makes ALR slower than AEX for the complex parameter space. Therefore, we recommend AEX for the model with low-dimensional summary statistics such as spatial autologistic model and ERGM. However, as we describe in Section 4.2.2, both algorithms become less practical as the parameter space increases.

Without low-dimensional summary statistics, AEX poses memory burdens. In principle, the Russian Roulette

algorithm can be applicable for the general form of the likelihood without sufficient statistics. The Russian Roulette algorithm does not require summary statistics or the assumption that  $h(\mathbf{x}|\theta)$  should be bounded. However, it may not be practical because constructing unbiased likelihood estimate requires  $\tau$  number of unbiased estimates  $\hat{Z}_i(\theta)$  per iteration. If  $\hat{Z}_i(\theta)$  requires  $N$  samples,  $\tau N$  samples need to be generated from the likelihood using inner sampler per iteration. For this reason, the Russian Roulette algorithm is computationally expensive compared to other methods.

We recommend DMH for computationally expensive likelihoods without low-dimensional summary statistics. Since only a single sample from the likelihood is required per iteration, DMH is computationally cheaper than any other algorithms. Furthermore, the algorithm is widely applicable without requiring any assumptions and relatively robust for high-dimensional parameter space. Once, we choose the appropriate number of the inner sampler updates, we can obtain plausible inference results in practice as well as the effective samples within the shortest time (ESS/T). However, exactness is not guaranteed in DMH. Finally, if we can afford some additional computing expenses, noisy DMH may be useful because it can result in fast mixing chains.

## 7. Discussion

We have discussed MCMC algorithms for likelihoods with intractable normalizing functions, explaining some theoretical underpinnings as well as examining practical implementation issues and efficiencies. We find that auxiliary variable approaches tend to be more efficient than likelihood approximation approaches, though efficiencies vary quite a bit depending on the dimension of the problem, and on the availability of low-dimensional sufficient statistics for models (especially for AEX and ALR algorithms). The specifics of the dataset also matter, for instance datasets with strong dependence can slow down algorithms for inferring parameters of spatial point processes. There is some overlap among the central ideas for most algorithms and all have a clear connection to importance sampling.

For users, practical implementation challenges are as important as theoretical justifications. Chopin and Ridgway (2017) also pointed out ease of use (fewer components to be tuned) and generalizability of code as important criteria for comparing Bayesian computation algorithms. From this perspective, AEX and ALR may be difficult choices for practitioners though both algorithms are asymptotically exact and moderately fast. This is because users have to manually tune lots of components, and do so differently for different problems. AEX seems to be a good choice for social network examples, at least in exponential random graph models where models have low-dimensional sufficient statistics. However, even here, Jin, Yuan, and Liang (2013) discussed the need to carefully tune AEX to address some well-known degeneracy problems (Handcock 2003; Schweinberger 2011).

Overall, we suggest starting with the DMH algorithm for its ease of implementation and the fact that it is typically computationally efficient. If instabilities in results are observed, for instance when the inner sampler length is increased, then AEX may be a more robust approach, even though it is much harder to implement. A preliminary run of DMH can be helpful for

both tuning (e.g., selecting particles) and debugging of the AEX algorithm. Sophisticated users of MCMC and particle-based algorithms could also consider ALR as an alternative, and noisy MCMC methods may also be useful if there are severe mixing issues.

Finally, we hope that we have convinced readers that there are a number of very ingenious approaches and ideas to consider for tackling inference with intractable normalizing function problems. There are opportunities for extending or even combining some of these algorithms—it is not hard to believe that with some adjustments to the algorithms and clever uses of, say, parallel computing, any of these algorithms, for example, the Russian roulette algorithm, is likely to become more efficient and hence more practical. Many of the ideas here relate to the even broader modern research problem of inference with models where likelihood functions are entirely unavailable or intractable. In the intractable normalizing function case, it is apparent that there is still a dearth of generic, automated approaches for constructing efficient MCMC algorithms. For many models inference is infeasible even for moderately high-dimensional problems (see Goldstein et al. 2015). Furthermore, careful approaches for comparing the efficiency and diagnosing convergence of asymptotically inexact algorithms is still an open problem. We hope that this manuscript will encourage further research on all of these challenging topics.

## Supplementary Materials

The supplementary material provides details about how to select particles in the adaptive exchange (AEX) algorithm and a description of the Russian roulette algorithm. It also provides details about how computational complexity calculations for various algorithm was calculated.

## Acknowledgment

The authors are grateful to Anne-Marie Lyne, Ick Hoon Jin, and Yves Atchade for providing useful sample code and advice, and to Faming Liang, Galin Jones and John Hughes for helpful comments.

## Funding

MH was partially supported by the National Science Foundation through NSF-DMS-1418090.

## References

- Alquier, P., Friel, N., Everitt, R., and Boland, A. (2016), “Noisy Monte Carlo: Convergence of Markov Chains With Approximate Transition Kernels,” *Statistics and Computing*, 26, 29–47. [1380]
- Andrieu, C., and Roberts, G. O. (2009), “The Pseudo-Marginal Approach for Efficient Monte Carlo Computations,” *The Annals of Statistics*, 37, 697–725. [1380]
- Atchade, Y., Lartillot, N., and Robert, C. P. (2008), “Bayesian Computation for Statistical Models With Intractable Normalizing Constants,” *Brazilian Journal of Probability and Statistics*, 27, 416–436. [1378, 1379, 1380, 1386]
- Atchadé, Y. F. (2006), “An Adaptive Version for the Metropolis Adjusted Langevin Algorithm With a Truncated Drift,” *Methodology and Computing in Applied Probability*, 8, 235–254. [1386]
- Atchade, Y. F., and Liu, J. S. (2004), “The Wang-Landau Algorithm for Monte Carlo Computation in General State Spaces,” *Statistica Sinica*, 20, 209–233. [1378]

- Beaumont, M. A. (2003), "Estimation of Population Growth or Decline in Genetically Monitored Populations," *Genetics*, 164, 1139–1160. [1380]
- Besag, J. (1974), "Spatial Interaction and the Statistical Analysis of Lattice Systems," *Journal of the Royal Statistical Society, Series B*, 36, 192–236. [1372,1374]
- Brooks, S., Gelman, A., Jones, G., and Meng, X.-L. (2011), *Handbook of Markov Chain Monte Carlo*, Boca Raton, FL: CRC Press. [1372]
- Butts, C. T. (2012), *A Perfect Sampling Method for Exponential Random Graph Models*, Irvine, CA: University of California. [1382]
- Caimo, A., and Friel, N. (2011), "Bayesian Inference for Exponential Random Graph Models," *Social Networks*, 33, 41–55. [1375]
- Chopin, N., and Ridgway, J. (2017), "Leave Pima Indians Alone: Binary Regression as a Benchmark for Bayesian Computation," *Statistical Science*, 32, 64–87. [1388]
- Drabek, T. E., Tamminga, H. L., Kilijaneck, T. S., and Adams, C. R. (1981), *Managing Multiorganizational Emergency Responses: Emergent Search and Rescue Networks in Natural Disaster and Remote Area Settings*, Boulder, CO: Inst. of Behavioral Science. [1383]
- Eddelbuettel, D., François, R., Allaire, J., Chambers, J., Bates, D., and Ushey, K. (2011), "Rcpp: Seamless R and C++ Integration," *Journal of Statistical Software*, 40, 1–18. [1381]
- Flegal, J. M., Haran, M., and Jones, G. L. (2008), "Markov Chain Monte Carlo: Can We Trust the Third Significant Figure?" *Statistical Science*, 23, 250–260. [1376,1381]
- Geyer, C. J. (1991), "Markov Chain Monte Carlo Maximum Likelihood," in *Computing Science and Statistics: Proceedings of the 23rd Symposium on the Interface*, pp. 156–163. [1374]
- (2011), "Introduction to Markov Chain Monte Carlo," in *Handbook of Markov Chain Monte Carlo*, eds. S. Brooks, A. Gelman, X.-L. Meng, and G. L. Jones, Boca Raton, FL: Chapman & Hall, pp. 295–311. [1372]
- Geyer, C. J., and Thompson, E. A. (1992), "Constrained Monte Carlo Maximum Likelihood for Dependent Data," *Journal of the Royal Statistical Society, Series B*, 54, 657–699. [1372,1378]
- Goldstein, J., Haran, M., Simeonov, I., Fricks, J., and Chiaromonte, F. (2015), "An Attraction-Repulsion Point Process Model for Respiratory Syncytial Virus Infections," *Biometrics*, 71, 376–385. [1373,1384,1388]
- Gong, L., and Flegal, J. M. (2015), "A Practical Sequential Stopping Rule for High-Dimensional Markov Chain Monte Carlo," *Journal of Computational and Graphical Statistics*, 25, 684–700. [1376]
- Handcock, M. S. (2003), *Statistical Models for Social Networks: Inference and Degeneracy* (Vol. 126), Washington, DC: National Academies Press. [1373,1388]
- Hastings, W. K. (1970), "Monte Carlo Sampling Methods Using Markov Chains and Their Applications," *Biometrika*, 57, 97–109. [1372]
- Hughes, J., Haran, M., and Caragea, P. (2011), "Autologistic Models for Binary Data on a Lattice," *Environmetrics*, 22, 857–871. [1372]
- Hunter, D. R., and Handcock, M. S. (2012), "Inference in Curved Exponential Family Models for Networks," *Journal of Computational and Graphical Statistics*, 15, 565–583. [1372]
- Hunter, D. R., Handcock, M. S., Butts, C. T., Goodreau, S. M., and Morris, M. (2008), "ergm: A Package to Fit, Simulate and Diagnose Exponential-Family Models for Networks," *Journal of Statistical Software*, 24, 1–29. [1382]
- Hunter, D. R., Krivitsky, P. N., and Schweinberger, M. (2012), "Computational Statistical Methods for Social Network Models," *Journal of Computational and Graphical Statistics*, 21, 856–882. [1373]
- Ibáñez, M. V., and Simó, A. (2003), "Parameter Estimation in Markov Random Field Image Modeling With Imperfect Observations. A Comparative Study," *Pattern Recognition Letters*, 24, 2377–2389. [1372]
- Ihaka, R., and Gentleman, R. (1996), "R: A Language for Data Analysis and Graphics," *Journal of Computational and Graphical Statistics*, 5, 299–314. [1380]
- Ising, E. (1925), "Beitrag zur Theorie des Ferromagnetismus," *Zeitschrift für Physik A Hadrons and Nuclei*, 31, 253–258. [1372,1381]
- Jin, I. H., and Liang, F. (2013), "Fitting Social Network Models Using Varying Truncation Stochastic Approximation MCMC Algorithm," *Journal of Computational and Graphical Statistics*, 22, 927–952. [1377]
- Jin, I. H., Yuan, Y., and Liang, F. (2013), "Bayesian Analysis for Exponential Random Graph Models Using the Adaptive Exchange Sampler," *Statistics and Its Interface*, 6, 559–576. [1378,1388]
- 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, 1537–1547. [1381]
- Jones, G. L., and Hobert, J. P. (2001), "Honest Exploration of Intractable Probability Distributions via Markov Chain Monte Carlo," *Statistical Science*, 16, 312–334. [1376]
- Kass, R. E., Carlin, B. P., Gelman, A., and Neal, R. M. (1998), "Markov Chain Monte Carlo in Practice: A Roundtable Discussion," *The American Statistician*, 52, 93–100. [1375,1381,1383,1384,1385]
- Lenz, W. (1920), "Beitrag zum Verständnis der Magnetischen Erscheinungen in Festen Körpern," *Physikalische Zeitschrift*, 21, 613–615. [1372,1381]
- Liang, F. (2007), "Continuous Contour Monte Carlo for Marginal Density Estimation With an Application to a Spatial Statistical Model," *Journal of Computational and Graphical Statistics*, 16, 608–632. [1378]
- Liang, F. (2010), "A Double Metropolis–Hastings Sampler for Spatial Models With Intractable Normalizing Constants," *Journal of Statistical Computation and Simulation*, 80, 1007–1022. [1375]
- Liang, F., and Jin, I.-H. (2013), "A Monte Carlo Metropolis–Hastings Algorithm for Sampling From Distributions With Intractable Normalizing Constants," *Neural Computation*, 25, 2199–2234. [1380]
- Liang, F., Jin, I. H., Song, Q., and Liu, J. S. (2016), "An Adaptive Exchange Algorithm for Sampling From Distributions With Intractable Normalizing Constants," *Journal of the American Statistical Association*, 111, 377–393. [1373,1376,1377,1378]
- Liang, F., Liu, C., and Carroll, R. J. (2007), "Stochastic Approximation in Monte Carlo Computation," *Journal of the American Statistical Association*, 102, 305–320. [1377]
- Lyne, A.-M., Girolami, M., Atchade, Y., Strathmann, H., and Simpson, D. (2015), "On Russian Roulette Estimates for Bayesian Inference With Doubly-Intractable Likelihoods," *Statistical Science*, 30, 443–467. [1380]
- Metropolis, N., Rosenbluth, A. W., Rosenbluth, M. N., Teller, A. H., and Teller, E. (1953), "Equation of State Calculations by Fast Computing Machines," *The Journal of Chemical Physics*, 21, 1087–1092. [1372]
- Mitrophanov, A. Y. (2005), "Sensitivity and Convergence of Uniformly Ergodic Markov Chains," *Journal of Applied Probability*, 42, 1003–1014. [1380]
- Møller, J., Pettitt, A. N., Reeves, R., and Berthelsen, K. K. (2006), "An Efficient Markov Chain Monte Carlo Method for Distributions With Intractable Normalising Constants," *Biometrika*, 93, 451–458. [1373,1374]
- Murray, I., Ghahramani, Z., and MacKay, D. J. C. (2006), "MCMC for Doubly-Intractable Distributions," in *Proceedings of the 22nd Annual Conference on Uncertainty in Artificial Intelligence (UAI-06)*, AUAI Press, pp. 359–366. [1374,1386]
- Murray, I. A. (2007), *Advances in Markov Chain Monte Carlo Methods*, London: University of London. [1375]
- Neal, R. M. (1996), "Sampling From Multimodal Distributions Using Tempered Transitions," *Statistics and Computing*, 6, 353–366. [1374]
- (2001), "Annealed Importance Sampling," *Statistics and Computing*, 11, 125–139. [1374]
- Padgett, J. F. (1994), "Marriage and Elite Structure in Renaissance Florence, 1282–1500," *Paper Delivered to the Social Science History Association*. [1382]
- Propp, J. G., and Wilson, D. B. (1996), "Exact Sampling With Coupled Markov Chains and Applications to Statistical Mechanics," *Random Structures and Algorithms*, 9, 223–252. [1373,1381]
- Pržulj, N. (2007), "Biological Network Comparison Using Graphlet Degree Distribution," *Bioinformatics*, 23, e177–e183. [1383]
- Pržulj, N., Corneil, D. G., and Jurisica, I. (2004), "Modeling Interactome: Scale-Free or Geometric?" *Bioinformatics*, 20, 3508–3515. [1383]
- Robert, C., and Casella, G. (2013), *Monte Carlo Statistical Methods*, Berlin: Springer Science & Business Media. [1381,1383,1384,1385]
- Roberts, G. O., and Rosenthal, J. S. (2007), "Coupling and Ergodicity of Adaptive Markov Chain Monte Carlo Algorithms," *Journal of Applied Probability*, 44, 458–475. [1378]
- Robins, G., Pattison, P., Kalish, Y., and Lusher, D. (2007), "An Introduction to Exponential Random Graph ( $p^*$ ) Models for Social Networks," *Social Networks*, 29, 173–191. [1372,1382]

- Rosenthal, J. S. (1995), “Minorization Conditions and Convergence Rates for Markov Chain Monte Carlo,” *Journal of the American Statistical Association*, 90, 558–566. [1376]
- Schweinberger, M. (2011), “Instability, Sensitivity, and Degeneracy of Discrete Exponential Families,” *Journal of the American Statistical Association*, 106, 1361–1370. [1388]
- Strauss, D. J. (1975), “A Model for Clustering,” *Biometrika*, 62, 467–475. [1384]
- Torrie, G. M., and Valleau, J. P. (1977), “Nonphysical Sampling Distributions in Monte Carlo Free-Energy Estimation: Umbrella Sampling,” *Journal of Computational Physics*, 23, 187–199. [1372,1378]
- Wang, F., and Landau, D. (2001), “Efficient, Multiple-Range Random Walk Algorithm to Calculate the Density of States,” *Physical Review Letters*, 86, 2050–2053. [1380]
- Yaveroglu, O. N., Fitzhugh, S. M., Kurant, M., Markopoulou, A., Butts, C. T., and Przulj, N. (2014), “ergm. Graphlets: A Package for ERG Modeling Based on Graphlet Statistics,” *Journal of Statistical Software*, 65, 1–29. [1383]
- Younes, L. (1988), “Estimation and Annealing for Gibbsian Fields,” *Annales de l’IHP Probabilités et Statistiques*, 24, 269–294. [1372,1379]