
Parallel Markov Chain Monte Carlo via Spectral Clustering

Guillaume Basse
Harvard University

Natesh Pillai
Harvard University

Aaron Smith
University of Ottawa

Abstract

As it has become common to use many computer cores in routine applications, finding good ways to parallelize popular algorithms has become increasingly important. In this paper, we present a parallelization scheme for Markov chain Monte Carlo (MCMC) methods based on spectral clustering of the underlying state space, generalizing earlier work on parallelization of MCMC methods by state space partitioning. We show empirically that this approach speeds up MCMC sampling for multimodal distributions and that it can be usefully applied in greater generality than several related algorithms. Our algorithm converges under reasonable conditions to an ‘optimal’ MCMC algorithm. We also show that our approach can be asymptotically far more efficient than naive parallelization, even in situations such as completely flat target distributions where no unique optimal algorithm exists. Finally, we combine theoretical and empirical bounds to provide practical guidance on the choice of tuning parameters.

1 INTRODUCTION

Markov chain Monte Carlo (MCMC) is a powerful and popular method for sampling from target distributions. As a sampling method, it is inherently parallel: simply run independent copies of the Markov chain on every available core. However, as MCMC has been used for a wider variety of problems, it has become clear that this ‘naive’ parallelization can often be improved upon. A major problem in the field is to develop new parallelization methods and find conditions under which they are better than the naive parallelization.

Appearing in Proceedings of the 19th International Conference on Artificial Intelligence and Statistics (AISTATS) 2016, Cadiz, Spain. JMLR: W&CP volume 41. Copyright 2016 by the authors.

Past approaches to parallelizing MCMC algorithms can be broadly divided into four categories: (i) *single-step* speedups, (ii) *exploration* speedups, (iii) learning or partitioning *data* and (iv) learning or partitioning the *state-space*. The first method uses several cores to increase the speed at which individual steps of a Metropolis-Hastings algorithm are taken (Wilkinson (2005); Calderhead (2014); Brockwell (2006); Angelino et al. (2014); Feng et al. (2003)), while the second method uses several cores to run slightly different chains, increasing the speed at which the algorithm mixes (see *e.g.*, parallel-tempering in Altekar et al. (2004) and the elliptical slice sampler in Nishihara et al. (2014)). Contrary to the first two approaches, the third approach applies mainly in the context of sampling from posterior distributions, and involves partitioning the data into batches and later recombining the MCMC samples (Scott et al. (2013); Wang and Dunson (2013); Neiswanger et al. (2013); Huang and Gelman (2005)). This approach is especially useful for large data sets as they are often stored across different machines. The fourth method often involves finding good partitions of the state space and running a different MCMC chain in each part of the partition (Hallgren and Koski (2014); VanDerwerken and Schmidler (2013)), though there exist other methods in this category (Craiu et al. (2009)). Although we describe four categories, these methods can generally all be applied at the same time. In addition, several data-augmentation chains incorporate the underlying dataset into the state space of the Markov chain (see, *e.g.*, Maclaurin and Adams (2014)), meaning that statespace-partitioning schemes, including the approach described in this paper, can be used as a first step in data-partitioning schemes.

1.1 Our Contributions

We propose a novel collection of methods for parallelizing MCMC by partitioning the underlying state space. Our key idea is to find partitions based on spectral clustering (see Von Luxburg (2007)). The main intuition behind all state space partitioning methods

is to replace a single Markov chain targeting a highly multimodal distribution with several Markov chains, each targeting distinct unimodal distributions. Since Markov chains tend to mix much more quickly on unimodal distributions than on multimodal distributions, this should improve computational efficiency. Our approach differs from existing space-partitioning approaches in that it includes a general way to find a partition (unlike Hallgren and Koski (2014)) and that the allowable partitions are extremely general and can in particular include non convex sets (this is in contrast to the more limited family of Voronoi partitions suggested in VanDerwerken and Schmidler (2013)). Perhaps surprisingly, we find that the additional flexibility of our family of partitions generally does not greatly increase the cost of finding a ‘good’ partition, *as long as a ‘good’ partition exists*.

We provide empirical evidence that our approach works well for benchmark problems and provide examples for which our method outperforms its competitors. We also provide a theoretical grounding for our approach. This includes some guarantees that the partitions we find converge to a ‘good’ partition and that ‘good’ partitions result in efficient MCMC chains. We use fully worked-out examples to illustrate the gains that our method can provide under optimal circumstances, as well as the fact that state space partitioning can give an advantage over naive parallelization even when the target distribution does not have strong clusters and when the partitions used are neither stable nor close to optimal. Finally, we discuss heuristics for the amount of computational effort that should be spent on finding a partition.

2 Intuition Behind Partitioning

Before discussing how our algorithm chooses a partition, we give notation and explain why state space partitioning methods work well once a good partition has been found. Let $\{Y_t\}_{t \in \mathbb{N}}$ be a Markov chain on a state space Ω with stationary distribution π and transition kernel K . Throughout this paper, we assume that the kernel K is a Metropolis-Hastings kernel associated with a proposal kernel Q , though our approach can be applied in other settings. For any $T \in \mathbb{N}$ and π -measurable function h , the usual MCMC estimate of $\mu \equiv \pi(h)$ is

$$\hat{\mu} = \frac{1}{T} \sum_{t=1}^T h(Y_t). \quad (1)$$

The computation of this estimate can be naively parallelized by running n independent chains

$\{Y_t^{(i)}\}_{t \in \mathbb{N}, 1 \leq i \leq n}$ and writing

$$\hat{\mu}_{\text{naive}} = \frac{1}{nT} \sum_{t=1}^T \sum_{i=1}^n h(Y_t^{(i)}). \quad (2)$$

Alternatively, fix a partition of Ω into disjoint subsets $\{\Omega_i\}_{i=1}^n$. Define the weights $w_i = \pi(\Omega_i)$ and distributions

$$\tilde{\pi}_i(A) = \pi(A \cap \Omega_i) \text{ and } \pi_i(A) = \frac{1}{w_i} \tilde{\pi}_i(A). \quad (3)$$

For $1 \leq i \leq n$, define K_i to be the Metropolis-Hastings kernel with proposal kernel Q and target distribution π_i and let $\{X_t^{(i)}\}_{t \in \mathbb{N}}$ be a Markov chain evolving according to K_i . For each fixed i , the estimate of $\mu_i \equiv \pi_i(h)$ that is analogous to (1) is given by

$$\hat{\mu}_{i,\text{par}} = \frac{1}{T} \sum_{t=1}^T h(X_t^{(i)}), \quad (4)$$

and since $\mu = \sum_i w_i \mu_i$, we use

$$\hat{\mu}_{\text{par}} = \sum_{i=1}^n w_i \hat{\mu}_{i,\text{par}} \quad (5)$$

as the full estimate of μ .

We now derive conditions under which the estimator in (5) has smaller variance than the estimator in (3). Let $\lambda = \sup\{|\lambda^*| : (\lambda^*I - K)^{-1} \text{ is not a bounded linear operator on } L^2(\pi), \lambda^* \neq 1\}$ and denote by $(1 - \lambda)$ the spectral gap of reversible kernel K , and similarly by $(1 - \lambda_i)$ the spectral gap of K_i . The normalized variance of the estimate (2) can be bounded (see e.g. Prop 4.29 of Aldous and Fill (2002)) by

$$T \text{Var}[\hat{\mu}_{\text{naive}}](1 - O(T^{-1})) \leq \nu_{\text{naive}} \equiv \frac{2\|h\|_{2,\pi}^2}{n(1-\lambda)}, \quad (6)$$

and generically there exists a function h for which this is close to equality for large T . Similarly,

$$T \text{Var}[\hat{\mu}_{\text{par}}](1 - O(T^{-1})) \leq \nu_{\text{par}} \equiv 2 \sum_i \frac{w_i^2 \|h\|_{2,\pi_i}^2}{1 - \lambda_i}, \quad (7)$$

and again this is generically close to equality for large T and worst-case h . Denote by Φ the set of all measurable n -partitions of Ω and by $\mathcal{P} = (\Omega_1, \dots, \Omega_n)$ an element of Φ . Equations (6) and (7) suggest that the estimate (5) obtained from a partitioned state space should be more efficient than the naive estimate (2) when

$$\sum_i \frac{nw_i^2(1-\lambda)}{1-\lambda_i} \frac{\|h\|_{2,\pi_i}^2}{\|h\|_{2,\pi}^2} < 1. \quad (8)$$

Since $\|h\|_{2,\pi_i}^2 = \int_x h(x)^2 \pi_i(dx) \leq \frac{1}{w_i} \|h\|_{2,\pi}^2$ for all i , this suggests choosing the partition

$$\mathcal{P} = \operatorname{argmin}_{\mathcal{P} \in \Phi} \left(\sum_i \frac{w_i}{1 - \lambda_i} \right) \quad (9)$$

to find an estimator that is efficient for generic functions h . Finding this partition is computationally difficult, and so we settle for finding a partition that makes a good proxy for $\sum_i \frac{w_i}{1 - \lambda_i}$ small. Define the *conductance* of K_i by

$$\begin{aligned} \phi_i(S) &= \frac{\int_{x \in S} K_i(x, S^c) d\pi_i(x)}{\pi_i(S) \pi_i(S^c)} \\ \phi_i &= \inf_{S: 0 < \pi_i(S) < \frac{1}{2}} \phi_i(S), \end{aligned} \quad (10)$$

with an analogous definition for the conductances associated with K . By Lawler and Sokal (1988),

$$\frac{\phi_i^2}{2} \leq (1 - \lambda_i) \leq 2\phi_i, \quad (11)$$

and so

$$\sum_i \frac{w_i}{1 - \lambda_i} \leq 2 \sum_i \frac{w_i}{\phi_i^2}. \quad (12)$$

Thus, we approximately minimize the objective function (9) by making the upper bound (12) small. It is known (see Meila and Shi (2000); Kannan et al. (2004)) that spectral clustering approximately finds

$$\mathcal{P} = \operatorname{argmin}_{\mathcal{P} \in \Phi} \sum_{i=1}^n w_i \phi(\mathcal{P}_i). \quad (13)$$

Although it is not obvious, this choice of partition also approximately minimizes the right-hand side of Equation (12) (see Lee et al. (2014)), and so throughout this paper we will generally choose our partitions via spectral clustering. By inequalities (6), (7) and (12), the condition

$$\sum_i \frac{nw_i^2(1 - \lambda)}{\phi_i^2} \frac{\|h\|_{2,\pi_i}^2}{\|h\|_{2,\pi}^2} \ll 1 \quad (14)$$

implies that $\nu_{\text{par}} \ll \nu_{\text{naive}}$.

Inequality (14) gives a sufficient condition under which partitioning results in a more efficient sampler than naive parallelization. We give some examples showing that the ‘optimal’ partition defined by the heuristic (13) can satisfy this condition, and also that even very poor approximations of this partition can vastly improve sampling efficiency. The following illustrates the enormous improvement that partitioning can achieve when each mode of a strongly multimodal target density is in a separate part of the partition:

Example 2.1 (Mixture of Gaussians). *Fix constants $0 < \tau \ll \sigma \ll 1$ and consider the Random Walk Metropolis-Hastings chain K on \mathbb{R} with proposal kernel $Q(x, \cdot) = \frac{1}{2\tau} \text{Unif}[x - \tau, x + \tau]$ and target distribution $\pi = \frac{1}{2} \text{N}(-1, \sigma) + \frac{1}{2} \text{N}(1, \sigma)$. By considering the set $S = (-\infty, 0]$, we can calculate from Equation (10) that this chain has conductance (and thus spectral gap by Equation (11)) at most $O(\tau^{-2} e^{-c\sigma^{-2}})$ for some fixed $0 < c < \infty$. We consider speeding up simulation from the target distribution by partitioning the state space \mathbb{R} into $n = 2$ parts. Although spectral clustering attempts to minimize the objective function (13) rather than the ‘correct’ objective function (9), Figure 1 suggests that both have the same minimizer: the partition $\mathcal{P} = \{(-\infty, 0], (0, \infty)\}$.*

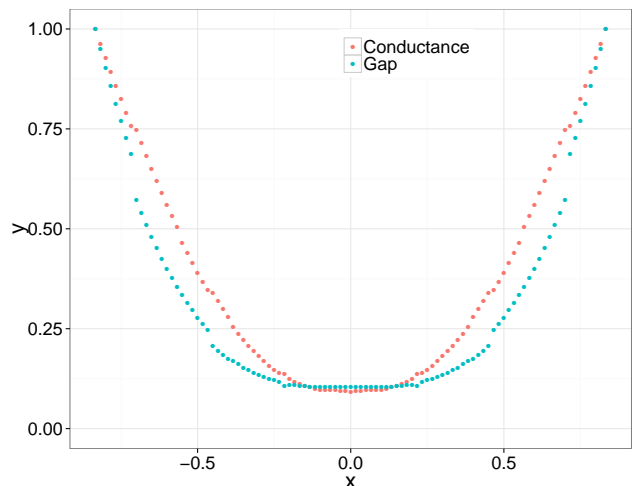


Figure 1: Comparing Two Objective Functions ($\sigma = 0.4$, $\tau = 0.2$) (Red is objective function for spectral clustering; blue is objective function for MCMC).

The Metropolis-Hastings chains K_1 and K_2 with proposal distribution Q and target distributions $\pi_1(x) \propto \pi(x) \mathbf{1}_{x \leq 0}$ and $\pi_2(x) \propto \pi(x) \mathbf{1}_{x > 0}$ have conductances that are at least on the order of $\frac{\tau}{\sigma}$ in the same regime (see Theorem 4.3.3 of Woodard (2007)), and thus spectral gaps that are at least on the order of $\frac{\tau^2}{\sigma^2}$ (again, see Equation (11)). Thus, the ratio of efficiencies (14) is at most $O(\frac{1}{\sigma^2} e^{-c\sigma^{-2}}) \ll 1$. Figure 2 plots the the inverse of the true value of this ratio, showing that it is enormous for reasonable values of σ and illustrating the gains of state space partitioning over naive parallelization. Note - this plot is on a logarithmic scale!

Figure 1 suggests that spectral clustering is essentially optimizing the right objective function, and Figure 2 shows that the estimator $\hat{\mu}_{\text{par}}$ associated with the best partition can be vastly more efficient than $\hat{\mu}_{\text{naive}}$. Unfortunately, in realistic examples, we will not have ac-

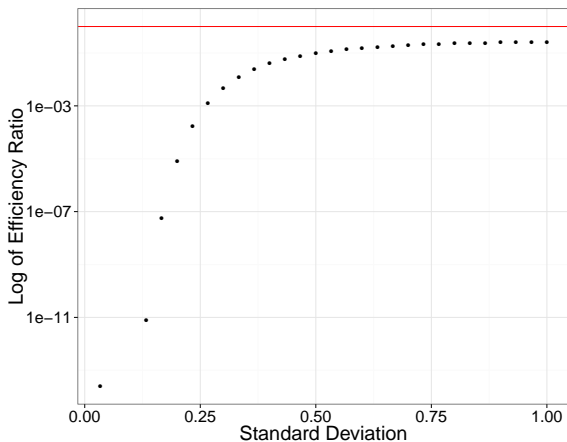


Figure 2: Log of Efficiency Ratio ($\log(\text{Var}(\hat{\mu}_{\text{par}})/\text{Var}(\hat{\mu}_{\text{naive}}))$ for $\tau = 0.2$).

cess to an optimal partition. Thus, it is natural to ask if similar gains can be obtained for partitions that are closer to those that might be seen in practice. Figure 3 shows the relative efficiency of $\hat{\mu}_{\text{par}}$ as the partition $\mathcal{P} = \{(-\infty, R], (R, \infty)\}$ ranges over various values of $R \geq 0$. It shows that $\hat{\mu}_{\text{par}}$ can be much more efficient than $\hat{\mu}_{\text{naive}}$ even when the partition used is quite far from optimal.

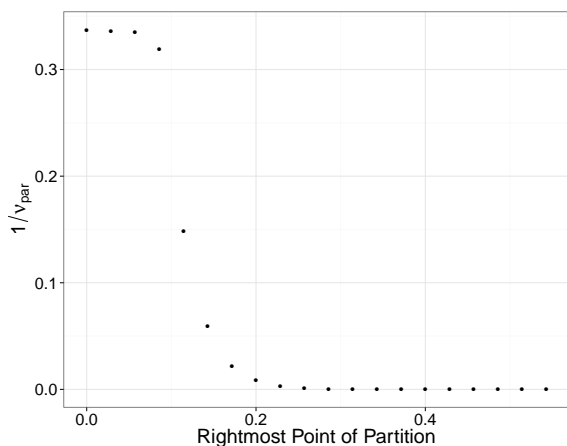


Figure 3: Inverse of the bound of the variance of $\hat{\mu}_{\text{par}}$ as R changes ($\tau = 0.2$, $\sigma = 0.15$).

Together, the three Figures 1, 2 and 3 justify our approach by showing that the partition returned by spectral clustering is closely related to the best partition, that the best partition can substantially increase efficiency, and finally that there is a fairly wide range of partitions that give rise to estimators with nearly-optimal efficiency.

A central point of this paper is that our parallelization

scheme can offer large improvements even when the partition used is far from optimal, and even in the absence of multimodality. This is most starkly illustrated by:

Example 2.2 (Simple Random Walk on the Cycle). Fix m and define the graph (V_m, E_m) with $V_m = \{1, 2, \dots, m\}$ and $E_m = \{(x, y) \in V_m : |x - y| \leq 1\} \cup \{(1, m)\}$. Recall that the simple random walk on the circle (here $\Omega = V_m$), which has transition kernel $\mathbb{P}[X_{t+1} = y | X_t = x] = \frac{1}{3} \mathbf{1}_{(x,y) \in E_m}$, has stationary distribution $\pi = \text{Unif}(\Omega)$ and spectral gap of order $\frac{1}{m^2}$ (see Section 12.3 of Levin et al. (2009)). Again, we fix a function h with $\|h\|_{2,\pi}^2 = 1$ and compare the efficiency of $\hat{\mu}_{\text{par}}$ to the efficiency of $\hat{\mu}_{\text{naive}}$. By Equation (6), the variance of $\hat{\mu}_{\text{naive}}$ is on the order of $\frac{m^2}{nT}$. If we partition Ω into n connected components $\{\mathcal{P}_j\}_{j=1}^n$, the associated kernels K_j have spectral gaps on the order of $\frac{1}{|\mathcal{P}_j|^2}$ (again, see Levin et al. (2009)). Thus, by Equation (7), the associated estimator $\hat{\mu}_{\text{par}}$ has variance on the order of $\frac{1}{m^2T} \sum_{j=1}^n |\mathcal{P}_j|^4$. It is easy to check that this is minimized by choosing $|\mathcal{P}_j| = \frac{m}{n}$, in which case $\hat{\mu}$ has an asymptotic variance of only $\frac{m^2}{n^3T}$. This is much smaller than the variance $\frac{m^2}{nT}$ of $\hat{\mu}_{\text{naive}}$, despite the fact that the target distribution is completely flat. We note again that suboptimal partitions can yield substantial improvements. For example, any choice of partition for which $\sup_j |\mathcal{P}_j| = o\left(\frac{m}{\sqrt{n}}\right)$ will lead to an estimator $\hat{\mu}_{\text{par}}$ with a smaller asymptotic variance than the estimator $\hat{\mu}_{\text{naive}}$. We also note that, due to the rotational symmetry of the simple random walk, there is not a unique optimal partition in this example. This non-uniqueness does not prevent us from using the algorithm presented in this paper. In most clustering applications, the resulting partition is of interest in and of itself, and is meaningful only if a unique good partition exists. In contrast, we only use the partitioning method as a way to reduce the variance of our estimator, and don't attach any particular meaning to it: we only need the existence of at least one good partition.

The remainder of this paper is concerned with describing an algorithm that produces a good partition efficiently, as well as describing its performance.

3 Methods

In this section, we lay out our approach and give some useful variations. Our main approach, summarized in Algorithm 1, has four steps: an initial exploration step, followed by repeated partitioning, sampling and weighting. In the first step, we explore the state space and try to capture as many modes as possible; this is necessary if we hope to have a reasonable first par-

tion. In the second step, we use spectral clustering and the history of the algorithm to find a ‘good’ partition $\{\Omega_i\}_{i=1}^n$ of the state space - that is, one for which the restricted chains K_i all have large conductance. In the third step, we run the chains $\{X_t^{(i)}\}_{t \in \mathbb{N}}$ in each component Ω_i of the partition. In the last step, we estimate the weights w_i of each element of the partition. The algorithm requires as input the number n of cores to be used (corresponding to the number of disjoint sets in the partitions), the proposal kernel Q , the target distribution π , the number N_0 of samples to be obtained by the initial exploration stage, the number ℓ of times that repartitioning occurs, the number $\{N_i\}_{1 \leq i \leq \ell}$ of samples used in each repartitioning step (first step of Algorithm 2), and the number $\{T_i\}_{1 \leq i \leq \ell}$ of steps to run the Markov chain between repartitionings. In practice, it is often helpful to have ℓ , T and especially N depend on the previous history of the chain rather than fixing them in advance.

Algorithm 1: Our Method

input : $Q, \pi, n, N_0, \ell, \{N_i\}_{1 \leq i \leq \ell}, \{T_i\}_{1 \leq i \leq \ell}$,
output: $\hat{\mu}$
begin
 Initialize X as in Section 3.1;
 for $i = 1$ **to** ℓ **do**
 $(\Omega_1, \dots, \Omega_n), (V_1, \dots, V_n) \leftarrow$
 DoSpectralClustering (X, n, N_i, Q, π) ;
 Compute $(\tilde{\pi}_1, \dots, \tilde{\pi}_n)$ as in Equation (3) ;
 $X \leftarrow X \cup \text{RunParallelChains}$
 $(Q, \tilde{\pi}_1, \tilde{\pi}_2, \dots, \tilde{\pi}_n, T_i)$;
 Compute $(\hat{w}_1, \dots, \hat{w}_n)$ as in Equations (17) and
 (18) ;
 end
 Estimate $\hat{\mu}$ as in Equations (20);
end

3.1 Explore State Space

In this step, we create a sample $X = (X_1, X_2, \dots, X_{N_0})$ of points from the state space Ω ; these points will be used to create an initial partition. Ideally, these points should cover every mode of the target π . In simulations, we have found that generating these points by running the parallel tempering algorithm (Swendsen and Wang (1986); Geyer (1991); Earl and Deem (2005)) from several well-dispersed initial points works well.

3.2 Partition State Space

In the *partitioning step* we obtain a partition of the state space, given the collection of points X that we have seen so far. We summarize our approach in Al-

gorithm 2.

Algorithm 2: DoSpectralClustering

input : X, n, N, Q, π
output: $(\Omega_1, \dots, \Omega_n), (V^1, \dots, V^n)$
begin
 Subsample N points uniformly and without
 replacement $X_1, \dots, X_N \sim \text{Unif}(X)$;
 Define the matrix $\hat{Q}_{ij} = Q(X_i, X_j)$ for
 $i, j \in \{1, \dots, N\}$;
 Define the diagonal matrix $D_{ii} = \sum_j \hat{Q}_{ij}$ for
 $i \in \{1, \dots, N\}$;
 Let $L = D^{-1/2} \hat{Q} D^{-1/2}$;
 Let V^1, \dots, V^n be the n normalized leading
 eigenvectors of L ;
 For $i=1, \dots, N$ let
 $Z_i = (V^1[i], \dots, V^n[i]) / \|(V^1[i], \dots, V^n[i])\|$;
 Define the map
 $\sigma : \{Z_1, Z_2, \dots, Z_N\} \rightarrow \{1, 2, \dots, n\}$ by **kmeans**
 $(Z_1, \dots, Z_N; n)$;
 Let $\{C_1, \dots, C_n\}$ be the n centers obtained by
 kmeans;
 Define the partition $\Omega = \sqcup_{i=1}^n \Omega_i$ by Equation
 (16);
end

The first step, in which we subsample N points from the original data X , is used to keep the computational burden manageable, as the dataset can be very large. The second-last step of this algorithm, **kmeans**, is the popular k -means clustering algorithm (see *e.g.*, Chapter 13 of Hastie et al. (2009)). The last step is to extend the partition σ of the set $\{Z_1, \dots, Z_N\}$ to a partition of the entire state space Ω . Let $\lambda_1, \dots, \lambda_n$ be the eigenvalues associated with V^1, \dots, V^n . Following Equations (8) through (12) of Bengio et al. (2003), for $x \in \Omega$ and $1 \leq i \leq n$, we define

$$Z_i(x) = \frac{\sqrt{N}}{\lambda_i} \sum_{j=1}^N V^i[j] Q(x, x_j). \quad (15)$$

Set $Z(x) = (Z_1(x), \dots, Z_n(x))$ and for $1 \leq i \leq n$ define

$$\Omega_i = \{x \in \Omega : \operatorname{argmin}_{j \in \{1, \dots, n\}} (\|C_j - Z(x)\|) = i\}. \quad (16)$$

3.3 Run Chains

We define the method **RunParallelChains**. For each $i \in \{1, 2, \dots, n\}$, we let $\{X_t^{(i)}\}_{t=1}^T$ be a Metropolis-Hastings chain with proposal kernel Q and target distribution $\tilde{\pi}_i$. The method then returns $\{X_t^{(i)}\}_{1 \leq t \leq T, 1 \leq i \leq n}$. We do not specify the initial points $X_1^{(i)}$, but have found in practice that the point

in $X \cap \Omega_i$ corresponding to the centroid given by the k -means algorithm is a good choice.

3.4 Estimate Weights

The final step is to estimate the weights $w_i = \pi(\Omega_i)$. There are a few options here, but we find that bridge sampling (Meng and Wong (1996), Gelman and Meng (1998)) works well. Recall that bridge sampling requires, for each $1 \leq i \leq n$, a proposal distribution p_i and bridge function α_i . In this paper, we generally choose p_i to be a normal or student distribution whose first two moments match the empirical moments of $\text{Unif}(X_i)$ and use the geometric bridge $\alpha_i = (p_i \tilde{\pi}_i)^{-1/2}$, where $\tilde{\pi}_i(x) = w_i \pi_i(x)$ is the unnormalized version of π_i .

For fixed p_i and α_i , let $X_i = X \cap \Omega_i$, $n_i = |X_i|$, and let $\{\theta_j\}_{1 \leq j \leq n_i}$ be i.i.d. draws from p_i . Then, define

$$\hat{c}_1^{(i)} = \frac{\frac{1}{n_i} \sum_{j=1}^{n_i} \tilde{\pi}_i(\theta_j) \alpha(\theta_j)}{\frac{1}{n_i} \sum_{\theta \in X_i} p_i(\theta) \alpha(\theta)}. \quad (17)$$

Since the estimates $\hat{c}_1^{(1)}, \dots, \hat{c}_1^{(n)}$ do not generally add up to 1, we renormalize:

$$\hat{w}_i = \frac{\hat{c}_1^{(i)}}{\sum_{j=1}^n \hat{c}_1^{(j)}}. \quad (18)$$

3.5 Estimating $\hat{\mu}$

We conclude by defining an estimator for $\hat{\mu}$. We start by writing:

$$\mu = \mathbb{E}_\pi(h(x)) = \sum_{i=1}^n \mu_i w_i \quad (19)$$

where $\mu_i = \pi(\Omega_i) \mathbb{E}_\pi(h(x) 1_{x \in \Omega_i}) = \mathbb{E}_{\pi_i}(h(x))$. We estimate μ_i by $\hat{\mu}_i$ as in Equation (4), and using the estimated weights $\hat{w}_1, \dots, \hat{w}_n$ we have the following estimator for μ :

$$\hat{\mu} = \sum_{i=1}^n \hat{\mu}_i \hat{w}_i \quad (20)$$

4 Applications

We describe the parameters and proposals used for the simulations in the Supplementary Material.

4.1 Example 1: Mixture of Gaussians in 2 dimensions

We start with an example used by VanDerwerken and Schmidler (2013), where the target distribution here is a mixture of bivariate normals. We compare the performance of three methods: parallel tempering, naive

parallelization, and our method. For our method, we obtained $N_0 = 8000$ samples using parallel tempering as our exploration phase (see Figure 4). We then ran our algorithm with $\ell = 1$ round of partitioning, $n = 4$ clusters, and $N_1 = 700$, thus obtaining estimates $\hat{w}_i, i = 1, \dots, 4$ of the weights. Considering those weights as fixed, we ran parallel constrained chains for an additional $T_1 = 4000$ iterations, and obtained an estimate $\hat{\mu}$ as in (19).

To evaluate our method we repeated the last step (the last 4000 iterations) 500 times, computed the Euclidean distance between our estimate and the true expectation for each replication, and computed the average squared error, and the sample standard deviation of the squared error. In total, for one estimate, we ran a total $8000 + 4000 + 4 \times 4000 = 28000$ iterations. But if we assume that each iteration (whether it be for parallel tempering or constrained metropolis) takes t seconds, and consider the fact that 4×4000 iterations are run in parallel on 4 cores, each estimate is obtained in $16000 \times t$ seconds. For parallel tempering, we obtained one estimate of the mean by running the algorithm for 16000 iterations. We again repeated this 500 times to obtain an average squared error and a sample standard deviation. For the naive method, we ran in parallel 4 independent chains initialized randomly from the target distribution, for 16000 iterations. In the end, if we consider that an iteration of parallel tempering takes the same time than an iteration of metropolis hastings, then all the methods take the $16000 \times t$ seconds (in fact, one iteration of parallel tempering takes a bit longer; thus, as measured by clock time, our method would look even better).

method	mean	se
ours	0.008	0.009
parallel tempering	0.21	0.28
naive parallel	14.08	13.5

Table 1: Square distance from the true mean in the 2D mixture of gaussians example

We see in Table 1 that our method dramatically reduces the mean squared error compared to parallel tempering, and improves on naive parallelization even more dramatically.

4.2 Example 2: Why Spectral Clustering?

In this section, we illustrate the flexibility of spectral clustering, and show that it can work in situations where Voronoi clustering fails. Define the two sets $\mathcal{S}_1 = \{(r \cos(\theta), r \sin(\theta)) : \theta \in [\frac{2\pi}{6}, \frac{10\pi}{6}], r \in [1, 1.1]\}$ and $\mathcal{S}_2 = \{(r \cos(\theta), (r-1) \sin(\theta)) : \theta \in [-\frac{4\pi}{6}, \frac{4\pi}{6}], r \in [1, 1.1]\}$, and consider the target distribution π that is

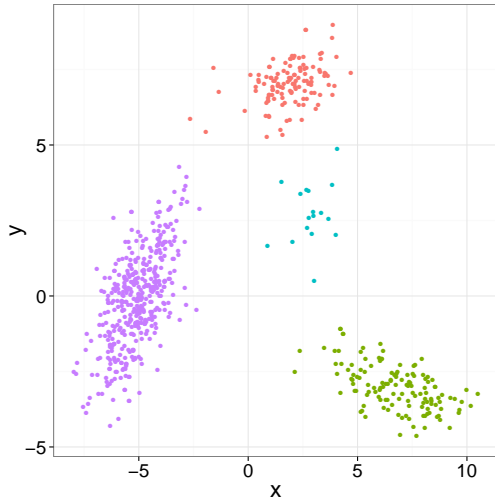


Figure 4: Spectral clustering of sample space for 2D mixture of Gaussians

uniform on $\mathcal{S} = \mathcal{S}_1 \cup \mathcal{S}_2$ (see Figure 5). We consider the simplified scenario where the exploration is carried by two chains of length 5000, initialized in the two sets, which results in a reasonably good picture of the overall density. We then perform spectral clustering on the one hand, and k-means (an instance of Voronoi clustering) on the other hand, using the $N_0 = 10000$ samples obtained in the exploration step. Figure 5 shows the results of the clustering phase: we see that because k-means can only find convex partitions, it doesn't capture the shapes adequately. To illustrate the impact of the choice of clustering method in terms of Monte-Carlo error, we carried the rest of our algorithm (estimating of the weights, running restricted parallel chains, and estimating the mean of the distribution) with parameters ($\ell = 1, N_1 = 700, T_1 = 10000$) for each clustering method. In particular, we simulated the last two steps of algorithm 1 (running restricted chains and estimating the mean) 200 times for each method, and computed the squared Euclidean distance of the estimated mean to the true mean. Using spectral clustering, the average squared distance is 0.04, with standard error 0.05, while with k-means, the average distance is 0.13, with standard error 0.03.

5 Convergence and Optimality

All proofs are in the Supplementary Material.

5.1 Consistency

We do not assume that the sequence of partitions we obtain at each stage $1 \leq i \leq \ell$ of Algorithm 1 is in any sense optimal, or that it converges in any sense to a

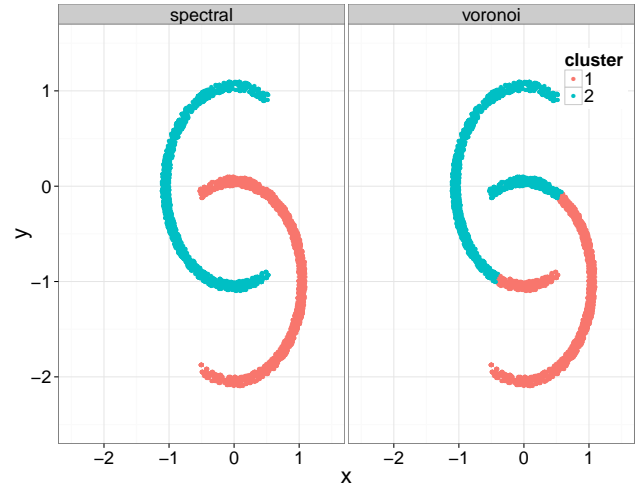


Figure 5: Spectral (left) and k-means (right) clusterings of the sample space

good partition. Indeed, as illustrated by example 2.2, our clustering algorithm can greatly increase computational efficiency even when there is not a unique optimal partition and when the partition used is far from any optimal partition. Even when there is not convergence to a unique optimal partition, the estimator $\hat{\mu}$ of μ returned by Algorithm 1 is generally consistent. We consider a simple setting:

Assumptions 5.1. Assume that the state space $\Omega \subset \mathbb{R}^d$ is bounded and that the target distribution π and the proposal distributions $\{Q(x, \cdot)\}_{x \in \Omega}$ have densities $\rho(\cdot)$ and $q(x, \cdot)$ that satisfy

$$c < \rho(y), q(x, y) < C \tag{21}$$

for some $0 < c < C < \infty$ and all $x, y \in \Omega$.

Theorem 5.2. Let Assumptions 5.1 hold, and assume further that h is bounded. Fix ℓ and $\{T_i\}_{i=1}^{\ell-1}$, and let $\hat{\mu}$ be the estimate returned by Algorithm 1. Then

$$\mathbb{P}[\lim_{T_\ell \rightarrow \infty} \hat{\mu} = \mu] = 1.$$

Although our method doesn't require that our partitions converge to an optimal partition, this convergence is desirable and does occur under reasonable conditions. For a partition $\mathcal{P} = \{\Omega_i\}_{i=1}^n$ of Ω , define an equivalence relation on Ω by writing $x \sim_{\mathcal{P}} y$ if and only if there exists some $1 \leq i \leq n$ such that $x, y \in \Omega_i$. Then define the distance d_π between pairs of partitions $\mathcal{P} = \{\Omega_i\}_{i=1}^n, \mathcal{P}' = \{\Omega'_i\}_{i=1}^n$ by

$$d_\pi(\mathcal{P}, \mathcal{P}') = \mathbb{P}[(X \sim_{\mathcal{P}} Y) \oplus (X \sim_{\mathcal{P}'} Y)],$$

where X, Y are drawn independently from π and \oplus denotes the logical operator *XOR*. For any kernel Q and

distribution π on Ω , Von Luxburg et al. (2008) defines an associated limiting Laplacian $\mathcal{L} = \mathcal{L}(Q, \pi)$. For any limiting Laplacian $\mathcal{L}(Q, \pi)$, Ben-David et al. (2006) defines the notion of a class of partitions $\mathcal{C} = \mathcal{C}(Q, \pi)$ associated with \mathcal{L} . We do not give precise definitions of these objects in this paper; the only heuristic needed is that $\mathcal{C}(Q, \pi)$ generally has exactly one element, unless Q, π have symmetries. We can then state the following corollary to Theorem 16 of Ben-David et al. (2006):

Theorem 5.3 (Convergence of Partitions). *Let Assumptions 5.1 hold. Fix a partition $\mathcal{P} = \{\Omega_j\}_{j=1}^n$ with associated measures $\{\pi_i\}_{i=1}^n$, so that $\mathcal{C}(Q, \frac{1}{n} \sum_{j=1}^n \pi_j)$ has a unique element \mathcal{P} . Fix $\gamma > 0$ and two sequences $\{N(k)\}_{k \in \mathbb{N}}, \{T(k)\}_{k \in \mathbb{N}}$ satisfying*

$$\begin{aligned} \lim_{k \rightarrow \infty} T(k) &= \infty & (22) \\ \lim_{k \rightarrow \infty} \frac{N(k)^{2+2\gamma}}{T(k)} &= 0. \end{aligned}$$

For $k \in \mathbb{N}$, let $X = \{X_t^{(i)}\}_{0 \leq t \leq T(k), 1 \leq i \leq n}$ be the output of the method **RunParallelChains**($Q, \pi_1, \dots, \pi_n, T(k)$) in Algorithm 1. Let $\mathcal{P}_{k,X}$ be the partition returned by **DoSpectralClustering**($X, n, N(k), Q, \pi$). Then for $\epsilon > 0$,

$$\lim_{k \rightarrow \infty} \mathbb{P}[d_\pi(\mathcal{P}_{k,X}, \mathcal{P}) > \epsilon] = 0. \quad (23)$$

This result implies that, if the partition $\{\Omega_i\}_{i=1}^n$ at stage $1 \leq q < \ell$ of Algorithm 1 is close to the optimal partition as measured by the metric d_π on partitions, the metric at stage $q+1$ can be made arbitrarily close as well by choosing N_{q+1}, T_{q+1} large.

5.2 Sample Size Heuristics

In this section, we discuss the choice of the sample size N used to compute each partition in Algorithm 1. We emphasize two facts:

1. Increasing N has essentially no impact on the mixing properties of K_i after a certain point N_{\max} .
2. If there exists an optimal partition $\{\Omega_i\}_{i=1}^n$, and $d_{i,j} = \mathbb{E}_{x \sim \pi_i, y \sim \pi_j} [|x - y|]$ represents the distance between parts of the partition while $1 - \lambda_i$ is the spectral gap of kernel K_i , we often have

$$N_{\max} \approx \left(\frac{\max_{1 \leq i < j \leq n} d_{i,j}}{\min_{1 \leq i \leq n} \pi(\Omega_i)(1 - \lambda_i)} \right)^2. \quad (24)$$

Together, these tell us that for the problems where our methods are most useful (*i.e.* where $\max_{1 \leq i \leq n} (1 - \lambda_i)$ is largest), the amount of effort that spent on finding the partitions should be small.

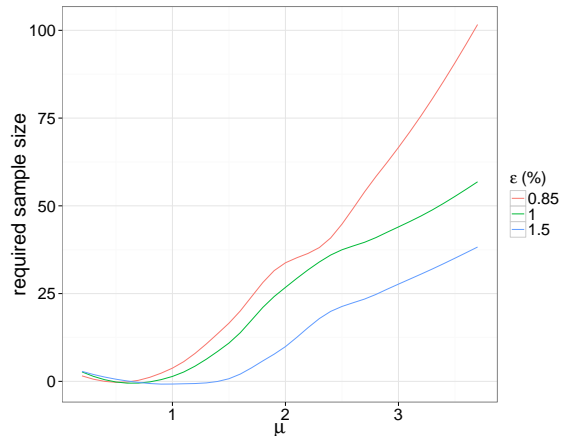


Figure 6: Number of points needed to create a partition within a certain distance of the sample optimal partition

The first heuristic follows from the fact that the partitions returned by **DoSpectralClustering** converge to the optimal partition under moderate conditions (see Theorem 5.3) and that the mixing time and spectral gap are continuous functions of the underlying transition kernel (see *e.g.* the main result of Mitrophanov (2005)).

Our justification for the second heuristic is empirical. Let $Q(x, x^*) = \frac{I(|x-x^*| < \tau)}{2\tau}$, and for $0 < \mu < \infty$ let $\pi_\mu = \frac{1}{2}\mathcal{N}(-\mu, 1) + \frac{1}{2}\mathcal{N}(\mu, 1)$. For any partition $\mathcal{P} = \{\Omega_i\}_{i=1}^n$, let $\lambda(\mathcal{P})$ be the smallest spectral gap of the associated kernels $\{K_i\}_{i=1}^n$ and let $\mathcal{P}_0 = \{(-\infty, 0], [0, \infty)\}$. Finally, for $0 < \epsilon < 1$, define

$$N_{\max}(\epsilon, \mu) = \min\{N : \lambda(\mathcal{C}_{N,\mu}) \geq (1 - \epsilon)\lambda(\mathcal{C}_0)\} \quad (25)$$

the number of points needed to create a partition that is within a factor of $1 - \epsilon$ of the optimal partition.

For each $\mu \in \{0.2, 0.3, \dots, 3.7\}$ we generated i.i.d. samples $X_\mu = \{X_i\} \sim \pi_\mu$ and generated partitions $\mathcal{P}_{N,\mu}$ for $N = \{1, 2, \dots\}$, according to **DoSpectralClustering**($X_\mu, 2, N, Q, \pi_\mu$). Figure 6 presents smoothed versions of the averages of the curves $\{N_{\max}(\epsilon, \mu)\}_{\mu \in \{0.2, 0.3, \dots, 3.7\}}$ for $\epsilon \in \{0.015, 0.01, 0.0085\}$ over 20 runs. This plot agrees fairly well with the heuristic (24), as do other generated plots. The most important property of our heuristic is that one need not spend an unlimited amount of computational resources to learn a ‘good enough’ partitioning of the state space, and that the computational resources required can be very modest if there exists a very good partition. Similarly to Daniely et al. (2012), for problems where we expect the method in this paper to work very well, we find the repartitioning step to be computationally inexpensive.

References

- David Aldous and James Allen Fill. Reversible Markov chains and random walks on graphs, 2002. Unfinished monograph, recompiled 2014, available online.
- Gautam Altekar, Sandhya Dwarkadas, John P Huelsenbeck, and Fredrik Ronquist. Parallel Metropolis coupled Markov chain Monte Carlo for Bayesian phylogenetic inference. *Bioinformatics*, 20(3):407–415, 2004.
- E. Angelino, E. Kohler, A. Waterland, M. Seltzer, and R. P. Adams. Accelerating MCMC via parallel predictive prefetching. *arXiv preprint arXiv:1403.7265*, March 2014.
- Shai Ben-David, Ulrike Von Luxburg, and Dávid Pál. A sober look at clustering stability. In *Learning theory*, pages 5–19. Springer, 2006.
- Yoshua Bengio, Jean-Francois Paiement, and Pascal Vincent. Out-of-sample extensions for lle, isomap, mds, eigenmaps, and spectral clustering. In *In Advances in Neural Information Processing Systems*, pages 177–184. MIT Press, 2003.
- AE Brockwell. Parallel Markov chain Monte Carlo simulation by pre-fetching. *Journal of Computational and Graphical Statistics*, 15(1):246–261, 2006.
- Ben Calderhead. A general construction for parallelizing Metropolis-Hastings algorithms. *Proceedings of the National Academy of Sciences*, 111(49):17408–17413, 2014.
- Radu V Craiu, Jeffrey Rosenthal, and Chao Yang. Learn from thy neighbor: Parallel-chain and regional adaptive mcmc. *Journal of the American Statistical Association*, 104(488):1454–1466, 2009.
- A. Daniely, N. Linial, and M. Saks. Clustering is difficult only when it does not matter. *arXiv preprint arXiv:1205.4891*, May 2012.
- David J Earl and Michael W Deem. Parallel tempering: Theory, applications, and new perspectives. *Physical Chemistry Chemical Physics*, 7(23):3910–3916, 2005.
- Xizhou Feng, Duncan A Buell, John R Rose, and Peter J Waddell. Parallel algorithms for Bayesian phylogenetic inference. *Journal of Parallel and Distributed Computing*, 63(7):707–718, 2003.
- Andrew Gelman and Xiao-Li Meng. Simulating normalizing constants: From importance sampling to bridge sampling to path sampling. *Statistical science*, pages 163–185, 1998.
- Charles J Geyer. Markov chain Monte Carlo maximum likelihood. 1991.
- J. Hallgren and T. Koski. Decomposition sampling applied to parallelization of Metropolis-Hastings. *arXiv preprint arXiv:1402.2828*, February 2014.
- Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The elements of statistical learning*, volume 2. Springer, 2009.
- Zaiying Huang and Andrew Gelman. Sampling for Bayesian computation with large datasets. unpublished, April 2005.
- Ravi Kannan, Santosh Vempala, and Adrian Vetta. On clusterings: Good, bad and spectral. *Journal of the ACM (JACM)*, 51(3):497–515, 2004.
- Gregory F Lawler and Alan D Sokal. Bounds on the l^2 spectrum for Markov chains and Markov processes: a generalization of Cheeger’s inequality. *Transactions of the American mathematical society*, 309(2):557–580, 1988.
- James R Lee, Shayan Oveis Gharan, and Luca Trevisan. Multiway spectral partitioning and higher-order cheeger inequalities. *Journal of the ACM (JACM)*, 61(6):37, 2014.
- David Asher Levin, Yuval Peres, and Elizabeth Lee Wilmer. *Markov chains and mixing times*. American Mathematical Soc., 2009.
- Dougal Maclaurin and Ryan P Adams. Firefly monte carlo: Exact mcmc with subsets of data. *arXiv preprint arXiv:1403.5693*, 2014.
- Marina Meila and Jianbo Shi. Learning segmentation by random walks. In *In Advances in Neural Information Processing*, pages 470–477. MIT Press, 2000.
- Xiao-li Meng and Wing Hung Wong. Simulating ratios of normalizing constants via a simple identity: A theoretical exploration. *Statistica Sinica*, pages 831–860, 1996.
- A Yu Mitrophanov. Sensitivity and convergence of uniformly ergodic markov chains. *Journal of Applied Probability*, pages 1003–1014, 2005.
- W. Neiswanger, C. Wang, and E. Xing. Asymptotically exact, embarrassingly parallel MCMC. *arXiv preprint arXiv:1311.4780*, November 2013.
- Robert Nishihara, Iain Murray, and Ryan P. Adams. Parallel MCMC with generalized elliptical slice sampling. *Journal of Machine Learning Research*, 15:2087–2112, 2014.
- Steven L. Scott, Alexander W. Blocker, and Fernando V. Bonassi. Bayes and big data: The consensus Monte Carlo algorithm. In *Bayes 250*, 2013.
- Robert H Swendsen and Jian-Sheng Wang. Replica Monte Carlo simulation of spin-glasses. *Physical Review Letters*, 57(21):2607, 1986.
- D. N. VanDerwerken and S. C. Schmidler. Parallel Markov chain Monte Carlo. *arXiv preprint arXiv:1312.7479*, December 2013.

Ulrike Von Luxburg. A tutorial on spectral clustering. *Statistics and computing*, 17(4):395–416, 2007.

Ulrike Von Luxburg, Mikhail Belkin, and Olivier Bousquet. Consistency of spectral clustering. *The Annals of Statistics*, pages 555–586, 2008.

X. Wang and D. B. Dunson. Parallelizing MCMC via Weierstrass Sampler. *arXiv preprint arXiv:1312.4605*, December 2013.

Darren J. Wilkinson. Parallel Bayesian computation. In Erricos J. Kontoghiorghes, editor, *Handbook of Parallel Computing and Statistics*, Statistics: A Series of Textbooks and Monographs. Chapman and Hall/CRC, 2005.

Dawn B. Woodard. *Conditions for Rapid and Torpid Mixing of Parallel and Simulated Tempering on Multimodal Distributions*. PhD thesis, Duke University, 2007.