

7-2014

Asymptotically Exact, Embarrassingly Parallel MCMC

Willie Neiswanger
Carnegie Mellon University

Chong Wang
Princeton University

Eric P. Xing
Carnegie Mellon University, epxing@cs.cmu.edu

Follow this and additional works at: http://repository.cmu.edu/machine_learning

 Part of the [Theory and Algorithms Commons](#)

Published In

Proceedings of the 30th International Conference on Conference on Uncertainty in Artificial Intelligence (UAI 2014).

This Conference Proceeding is brought to you for free and open access by the School of Computer Science at Research Showcase @ CMU. It has been accepted for inclusion in Machine Learning Department by an authorized administrator of Research Showcase @ CMU. For more information, please contact research-showcase@andrew.cmu.edu.

Asymptotically Exact, Embarrassingly Parallel MCMC

Willie Neiswanger

Machine Learning Department
Carnegie Mellon University
willie@cs.cmu.edu

Chong Wang

chongw@cs.princeton.edu

Eric P. Xing

School of Computer Science
Carnegie Mellon University
epxing@cs.cmu.edu

Abstract

Communication costs, resulting from synchronization requirements during learning, can greatly slow down many parallel machine learning algorithms. In this paper, we present a parallel Markov chain Monte Carlo (MCMC) algorithm in which subsets of data are processed independently, with very little communication. First, we arbitrarily partition data onto multiple machines. Then, on each machine, any classical MCMC method (e.g., Gibbs sampling) may be used to draw samples from a posterior distribution given the data subset. Finally, the samples from each machine are combined to form samples from the full posterior. This embarrassingly parallel algorithm allows each machine to act independently on a subset of the data (without communication) until the final combination stage. We prove that our algorithm generates asymptotically exact samples and empirically demonstrate its ability to parallelize burn-in and sampling in several models.

1 Introduction

Markov chain Monte Carlo (MCMC) methods are popular tools for performing approximate Bayesian inference via posterior sampling. One major benefit of these techniques is that they guarantee asymptotically exact recovery of the posterior distribution as the number of posterior samples grows. However, MCMC methods may take a prohibitively long time, since for N data points, most methods must perform $O(N)$ operations to draw a sample. Furthermore, MCMC methods might require a large number of “burn-in” steps before beginning to generate representative samples. Further complicating matters is the issue that, for many big data applications, it is necessary to store and process

data on multiple machines, and so MCMC must be adapted to run in these data-distributed settings.

Researchers currently tackle these problems independently, in two primary ways. To speed up sampling, multiple independent chains of MCMC can be run in parallel [20, 11, 13]; however, each chain is still run on the entire dataset, and there is no speed-up of the burn-in process (as each chain must still complete the full burn-in before generating samples). To run MCMC when data is partitioned among multiple machines, each machine can perform computation that involves a subset of the data and exchange information at each iteration to draw a sample [10, 14, 18]; however, this requires a significant amount of communication between machines, which can greatly increase computation time when machines wait for external information [1, 7].

We aim to develop a procedure to tackle both problems simultaneously, to allow for quicker burn-in and sampling in settings where data are partitioned among machines. To accomplish this, we propose the following: on each machine, run MCMC on only a subset of the data (independently, without communication between machines), and then combine the samples from each machine to algorithmically construct samples from the full-data posterior distribution. We’d like our procedure to satisfy the following four criteria:

1. Each machine only has access to a portion of the data.
2. Each machine performs MCMC independently, without communicating (i.e. “embarrassingly parallel”).
3. Each machine can use any type of MCMC to generate samples.
4. The combination procedure yields provably asymptotically exact samples from the full-data posterior.

The third criterion allows existing MCMC algorithms or software packages to be run directly on subsets of the data—the combination procedure then acts as a post-processing step to transform the samples to the correct distribution. Note that this procedure is particularly

suitable for use in a MapReduce [4] framework. Also note that, unlike current strategies, this procedure does not involve multiple “duplicate” chains (as each chain uses a different portion of the data and samples from a different posterior distribution), nor does it involve parallelizing a single chain (as there are multiple chains operating independently). We will show how this allows our method to, in fact, parallelize and greatly reduce the time required for burn-in.

In this paper we will (1) introduce and define the *subposterior* density—a modified posterior given a subset of the data—which will be used heavily, (2) present methods for the embarrassingly parallel MCMC and combination procedure, (3) prove theoretical guarantees about the samples generated from our algorithm, (4) describe the current scope of the presented method (i.e. where and when it can be applied), and (5) show empirical results demonstrating that we can achieve speed-ups for burn-in and sampling while meeting the above four criteria.

2 Embarrassingly Parallel MCMC

The basic idea behind our method is to partition a set of N i.i.d. data points $x^N = \{x_1, \dots, x_N\}$ into M subsets, sample from the *subposterior*—the posterior given a data subset with an underweighted prior—in parallel, and then combine the resulting samples to form samples from the full-data posterior $p(\theta|x^N)$, where $\theta \in \mathbb{R}^d$ and $p(\theta|x^N) \propto p(\theta)p(x^N|\theta) = p(\theta)\prod_{i=1}^N p(x_i|\theta)$.

More formally, given data x^N partitioned into M subsets $\{x^{n_1}, \dots, x^{n_M}\}$, the procedure is:

1. For $m = 1, \dots, M$ (in parallel):
Sample from the subposterior p_m , where

$$p_m(\theta) \propto p(\theta)^{\frac{1}{M}} p(x^{n_m}|\theta). \quad (1)$$

2. Combine the subposterior samples to produce samples from an estimate of the subposterior density product $p_1 \cdots p_M$, which is proportional to the full-data posterior, i.e. $p_1 \cdots p_M(\theta) \propto p(\theta|x^N)$.

We want to emphasize that we do not need to iterate over these steps and the combination stage (step 3) is the only step that requires communication between machines. Also note that sampling from each subposterior (step 2) can typically be done in the same way as one would sample from the full-data posterior. For example, when using the Metropolis-Hastings algorithm, one would compute the likelihood ratio as $\frac{p(\theta^*)^{\frac{1}{M}} p(x^{n_m}|\theta^*)}{p(\theta)^{\frac{1}{M}} p(x^{n_m}|\theta)}$ instead of $\frac{p(\theta^*)p(x^N|\theta^*)}{p(\theta)p(x^N|\theta)}$, where θ^* is the proposed move. In the next section, we show how the combination stage (step 3) is carried out to generate samples from the full-data posterior using the subposterior samples.

3 Combining Subposterior Samples

Our general idea is to combine the subposterior samples in such a way that we are implicitly sampling from an estimate of the subposterior density product function $\widehat{p_1 \cdots p_M}(\theta)$. If our density product estimator is consistent, then we can show that we are drawing asymptotically exact samples from the full posterior. Further, by studying the estimator error rate, we can explicitly analyze how quickly the distribution from which we are drawing samples is converging to the true posterior (and thus compare different combination algorithms).

In the following three sections we present procedures that yield samples from different estimates of the density product. Our first example is based on a simple parametric estimator motivated by the Bernstein-von Mises theorem [12]; this procedure generates approximate (asymptotically biased) samples from the full posterior. Our second example is based on a nonparametric estimator, and produces asymptotically exact samples from the full posterior. Our third example is based on a semiparametric estimator, which combines beneficial aspects from the previous two estimators while also generating asymptotically exact samples.

3.1 Approximate posterior sampling with a parametric estimator

The first method for forming samples from the full posterior given subposterior samples involves using an approximation based on the Bernstein-von Mises (Bayesian central limit) theorem, an important result in Bayesian asymptotic theory. Assuming that a unique, true data-generating model exists and is denoted θ_0 , this theorem states that the posterior tends to a normal distribution concentrated around θ_0 as the number of observations grows. In particular, under suitable regularity conditions, the posterior $P(\theta|x^N)$ is well approximated by $\mathcal{N}_d(\theta_0, F_N^{-1})$ (where F_N is the fisher information of the data) when N is large [12]. Since we aim to perform posterior sampling when the number of observations is large, a normal parametric form often serves as a good posterior approximation. A similar approximation was used in [2] in order to facilitate fast, approximately correct sampling. We therefore estimate each subposterior density with $\widehat{p}_m(\theta) = \mathcal{N}_d(\theta|\widehat{\mu}_m, \widehat{\Sigma}_m)$ where $\widehat{\mu}_m$ and $\widehat{\Sigma}_m$ are the sample mean and covariance, respectively, of the subposterior samples. The product of the M subposterior densities will be proportional to a Gaussian pdf, and our estimate of the density product function $p_1 \cdots p_M(\theta) \propto p(\theta|x^N)$ is

$$\widehat{p_1 \cdots p_M}(\theta) = \widehat{p}_1 \cdots \widehat{p}_M(\theta) \propto \mathcal{N}_d\left(\theta|\widehat{\mu}_M, \widehat{\Sigma}_M\right),$$

where the parameters of this distribution are

$$\widehat{\Sigma}_M = \left(\sum_{m=1}^M \widehat{\Sigma}_m^{-1} \right)^{-1} \quad (2)$$

$$\widehat{\mu}_M = \widehat{\Sigma}_M \left(\sum_{m=1}^M \widehat{\Sigma}_m^{-1} \widehat{\mu}_m \right). \quad (3)$$

These parameters can be computed quickly and, if desired, online (as new subposterior samples arrive).

3.2 Asymptotically exact posterior sampling with nonparametric density product estimation

In the previous method we made a parametric assumption based on the Bernstein-von Mises theorem, which allows us to generate approximate samples from the full posterior. Although this parametric estimate has quick convergence, it generates asymptotically biased samples, especially in cases where the posterior is particularly non-Gaussian. In this section, we develop a procedure that implicitly samples from the product of nonparametric density estimates, which allows us to produce asymptotically exact samples from the full posterior. By constructing a consistent density product estimator from which we can generate samples, we ensure that the distribution from which we are sampling converges to the full posterior.

Given T samples¹ $\{\theta_{t_m}^m\}_{t_m=1}^T$ from a subposterior p_m , we can write the kernel density estimator $\widehat{p}_m(\theta)$ as,

$$\begin{aligned} \widehat{p}_m(\theta) &= \frac{1}{T} \sum_{t_m=1}^T \frac{1}{h^d} K \left(\frac{\|\theta - \theta_{t_m}^m\|}{h} \right) \\ &= \frac{1}{T} \sum_{t_m=1}^T \mathcal{N}_d(\theta | \theta_{t_m}^m, h^2 I_d), \end{aligned}$$

where we have used a Gaussian kernel with bandwidth parameter h . After we have obtained the kernel density estimator $\widehat{p}_m(\theta)$ for M subposteriors, we define our nonparametric density product estimator for the full posterior as

$$\begin{aligned} \widehat{p_1 \cdots p_M}(\theta) &= \widehat{p}_1 \cdots \widehat{p}_M(\theta) \\ &= \frac{1}{T^M} \prod_{m=1}^M \sum_{t_m=1}^T \mathcal{N}_d(\theta | \theta_{t_m}^m, h^2 I_d) \\ &\propto \sum_{t_1=1}^T \cdots \sum_{t_M=1}^T w_t \cdot \mathcal{N}_d \left(\theta \mid \bar{\theta}_t, \frac{h^2}{M} I_d \right). \quad (4) \end{aligned}$$

¹For ease of description, we assume each machine generates the same number of samples, T . In practice, they do not have to be the same.

This estimate is the probability density function (pdf) of a mixture of T^M Gaussians with *unnormalized* mixture weights w_t . Here, we use $t \cdot = \{t_1, \dots, t_M\}$ to denote the set of indices for the M samples $\{\theta_{t_1}^1, \dots, \theta_{t_M}^M\}$ (each from a separate machine) associated with a given mixture component, and we define

$$\bar{\theta}_t = \frac{1}{M} \sum_{m=1}^M \theta_{t_m}^m \quad (5)$$

$$w_t = \prod_{m=1}^M \mathcal{N}_d(\theta_{t_m}^m | \bar{\theta}_t, h^2 I_d). \quad (6)$$

Although there are T^M possible mixture components, we can efficiently generate samples from this mixture by first sampling a mixture component (based on its unnormalized component weight w_t) and then sampling from this (Gaussian) component. In order to sample mixture components, we use an independent Metropolis within Gibbs (IMG) sampler. This is a form of MCMC, where at each step in the Markov chain, a single dimension of the current state is proposed (i.e. sampled) independently of its current value (while keeping the other dimensions fixed) and then is accepted or rejected. In our case, at each step, a new mixture component is proposed by redrawing one of the M current sample indices $t_m \in t \cdot$ associated with the component uniformly and then accepting or rejecting the resulting proposed component based on its mixture weight. We give the IMG algorithm for combining subposterior samples in Algorithm 1.²

In certain situations, Algorithm 1 may have a low acceptance rate and therefore may mix slowly. One way to remedy this is to perform the IMG combination algorithm multiple times, by first applying it to groups of $\tilde{M} < M$ subposteriors and then applying the algorithm again to the output samples from each initial application. For example, one could begin by applying the algorithm to all $\frac{M}{2}$ pairs (leaving one subposterior alone if M is odd), then repeating this process—forming pairs and applying the combination algorithm to pairs only—until there is only one set of samples remaining, which are samples from the density product estimate.

3.3 Asymptotically exact posterior sampling with semiparametric density product estimation

Our first example made use of a parametric estimator, which has quick convergence, but may be asymptotically biased, while our second example made use of

²Again for simplicity, we assume that we generate T samples to represent the full posterior, where T is the number of subposterior samples from each machine.

Algorithm 1 Asymptotically Exact Sampling via Non-parametric Density Product Estimation

Input: Subposterior samples: $\{\theta_{t_1}^1\}_{t_1=1}^T \sim p_1(\theta), \dots, \{\theta_{t_M}^M\}_{t_M=1}^T \sim p_M(\theta)$

Output: Posterior samples (asymptotically, as $T \rightarrow \infty$): $\{\theta_i\}_{i=1}^T \sim p_1 \cdots p_M(\theta) \propto p(\theta|x^N)$

- 1: Draw $t \cdot = \{t_1, \dots, t_M\} \stackrel{\text{iid}}{\sim} \text{Unif}(\{1, \dots, T\})$
 - 2: **for** $i = 1$ **to** T **do**
 - 3: Set $h \leftarrow i^{-1/(4+d)}$
 - 4: **for** $m = 1$ **to** M **do**
 - 5: Set $c \leftarrow t \cdot$
 - 6: Draw $c_m \sim \text{Unif}(\{1, \dots, T\})$
 - 7: Draw $u \sim \text{Unif}([0, 1])$
 - 8: **if** $u < w_{c \cdot} / w_{t \cdot}$ **then**
 - 9: Set $t \leftarrow c \cdot$
 - 10: **end if**
 - 11: **end for**
 - 12: Draw $\theta_i \sim \mathcal{N}_d(\bar{\theta}_{t \cdot}, \frac{h^2}{M} I_d)$
 - 13: **end for**
-

a nonparametric estimator, which is asymptotically exact, but may converge slowly when the number of dimensions is large. In this example, we implicitly sample from a semiparametric density product estimate, which allows us to leverage the fact that the full posterior has a near-Gaussian form when the number of observations is large, while still providing an asymptotically unbiased estimate of the posterior density, as the number of subposterior samples $T \rightarrow \infty$.

We make use of a semiparametric density estimator for p_m that consists of the product of a parametric estimator $\hat{f}_m(\theta)$ (in our case $\mathcal{N}_d(\theta|\hat{\mu}_m, \hat{\Sigma}_m)$ as above) and a nonparametric estimator $\hat{r}(\theta)$ of the correction function $r(\theta) = p_m(\theta)/\hat{f}_m(\theta)$ [6]. This estimator gives a near-Gaussian estimate when the number of samples is small, and converges to the true density as the number of samples grows. Given T samples $\{\theta_{t_m}^m\}_{t_m=1}^T$ from a subposterior p_m , we can write the estimator as

$$\begin{aligned} \hat{p}_m(\theta) &= \hat{f}_m(\theta) \hat{r}(\theta) \\ &= \frac{1}{T} \sum_{t_m=1}^T \frac{1}{h^d} K\left(\frac{\|\theta - \theta_{t_m}^m\|}{h}\right) \frac{\hat{f}_m(\theta)}{\hat{f}_m(\theta_{t_m}^m)} \\ &= \frac{1}{T} \sum_{t_m=1}^T \frac{\mathcal{N}_d(\theta|\theta_{t_m}^m, h^2 I_d) \mathcal{N}_d(\theta|\hat{\mu}_m, \hat{\Sigma}_m)}{\mathcal{N}_d(\theta_{t_m}^m|\hat{\mu}_m, \hat{\Sigma}_m)}, \end{aligned}$$

where we have used a Gaussian kernel with bandwidth parameter h for the nonparametric component of this estimator. Therefore, we define our semiparametric

density product estimator to be

$$\begin{aligned} \widehat{p_1 \cdots p_M}(\theta) &= \hat{p}_1 \cdots \hat{p}_M(\theta) \\ &= \frac{1}{T^M} \prod_{m=1}^M \sum_{t_m=1}^T \frac{\mathcal{N}_d(\theta|\theta_{t_m}^m, h I_d) \mathcal{N}_d(\theta|\hat{\mu}_m, \hat{\Sigma}_m)}{h^d \mathcal{N}_d(\theta_{t_m}^m|\hat{\mu}_m, \hat{\Sigma}_m)} \\ &\propto \sum_{t_1=1}^T \cdots \sum_{t_M=1}^T W_{t \cdot} \mathcal{N}_d(\theta|\mu_{t \cdot}, \Sigma_{t \cdot}). \end{aligned}$$

This estimate is proportional to the pdf of a mixture of T^M Gaussians with unnormalized mixture weights,

$$W_{t \cdot} = \frac{w_{t \cdot} \mathcal{N}_d(\bar{\theta}_{t \cdot}|\hat{\mu}_M, \hat{\Sigma}_M + \frac{h}{M} I_d)}{\prod_{m=1}^M \mathcal{N}_d(\theta_{t_m}^m|\hat{\mu}_m, \hat{\Sigma}_m)},$$

where $\bar{\theta}_{t \cdot}$ and $w_{t \cdot}$ are given in Eqs. 5 and 6. We can write the parameters of a given mixture component $\mathcal{N}_d(\theta|\mu_{t \cdot}, \Sigma_{t \cdot})$ as

$$\begin{aligned} \Sigma_{t \cdot} &= \left(\frac{M}{h} I_d + \hat{\Sigma}_M^{-1} \right)^{-1}, \\ \mu_{t \cdot} &= \Sigma_{t \cdot} \left(\frac{M}{h} I_d \bar{\theta}_{t \cdot} + \hat{\Sigma}_M^{-1} \hat{\mu}_M \right), \end{aligned}$$

where $\hat{\mu}_M$ and $\hat{\Sigma}_M$ are given by Eq. 2 and 3. We can sample from this semiparametric estimate using the IMG procedure outlined in Algorithm 1, replacing the component weights $w_{t \cdot}$ with $W_{t \cdot}$ and the component parameters $\bar{\theta}_{t \cdot}$ and $\frac{h}{M} I_d$ with $\mu_{t \cdot}$ and $\Sigma_{t \cdot}$.

We also have a second semiparametric procedure that may give higher acceptance rates in the IMG algorithm. We follow the above semiparametric procedure, where each component is a normal distribution with parameters $\mu_{t \cdot}$ and $\Sigma_{t \cdot}$, but we use the nonparametric component weights $w_{t \cdot}$ instead of $W_{t \cdot}$. This procedure is also asymptotically exact, since the semiparametric component parameters $\mu_{t \cdot}$ and $\Sigma_{t \cdot}$ approach the nonparametric component parameters $\bar{\theta}_{t \cdot}$ and $\frac{h}{M} I_d$ as $h \rightarrow 0$, and thus this procedure tends to the nonparametric procedure given in Algorithm 1.

4 Method Complexity

Given M data subsets, to produce T samples in d dimensions with the nonparametric or semiparametric asymptotically exact procedures (Algorithm 1) requires $O(dTM^2)$ operations. The variation on this algorithm that performs this procedure $M-1$ times on pairs of subposteriors (to increase the acceptance rate; detailed in Section 3.2) instead requires only $O(dTM)$ operations.

We have presented our method as a two step procedure, where first parallel MCMC is run to completion, and

then the combination algorithm is applied to the M sets of samples. We can instead perform an online version of our algorithm: as each machine generates a sample, it immediately sends it to a master machine, which combines the incoming samples³ and performs the accept or reject step (Algorithm 1, lines 3-12). This allows the parallel MCMC phase and the combination phase to be performed in parallel, and does not require transferring large volumes of data, as only a single sample is ever transferred at a time.

The total communication required by our method is transferring $O(dTM)$ scalars (T samples from each of M machines), and as stated above, this can be done online as MCMC is being carried out. Further, the communication is unidirectional, and each machine does not pause and wait for any information from other machines during the parallel sampling procedure.

5 Theoretical Results

Our second and third procedures aim to draw asymptotically exact samples by sampling from (fully or partially) nonparametric estimates of the density product. We prove the asymptotic correctness of our estimators, and bound their rate of convergence. This will ensure that we are generating asymptotically correct samples from the full posterior as the number of samples T from each subposterior grows.

5.1 Density product estimate convergence and risk analysis

To prove (mean-square) consistency of our estimator, we give a bound on the mean-squared error (MSE), and show that it tends to zero as we increase the number of samples drawn from each subposterior. To prove this, we first bound the bias and variance of the estimator. The following proofs make use of similar bounds on the bias and variance of the nonparametric and semiparametric density estimators, and therefore the theory applies to both the nonparametric and semiparametric density product estimators.

Throughout this analysis, we assume that we have T samples $\{\theta_{t_m}^m\}_{t_m=1}^T \subset \mathcal{X} \subset \mathbb{R}^d$ from each subposterior ($m = 1, \dots, M$), and that $h \in \mathbb{R}_+$ denotes the bandwidth of the nonparametric density product estimator (which is annealed to zero as $T \rightarrow \infty$ in Algorithm 1). Let Hölder class $\Sigma(\beta, L)$ on \mathcal{X} be defined as the set of all $\ell = \lfloor \beta \rfloor$ times differentiable functions $f : \mathcal{X} \rightarrow \mathbb{R}$ whose derivative $f^{(\ell)}$ satisfies

$$|f^{(\ell)}(\theta) - f^{(\ell)}(\theta')| \leq L |\theta - \theta'|^{\beta - \ell} \quad \text{for all } \theta, \theta' \in \mathcal{X}.$$

³For the semiparametric method, this will involve an online update of mean and variance Gaussian parameters.

We also define the class of densities $\mathcal{P}(\beta, L)$ to be

$$\mathcal{P}(\beta, L) = \left\{ p \in \Sigma(\beta, L) \mid p \geq 0, \int p(\theta) d\theta = 1 \right\}.$$

We also assume that all subposterior densities p_m are bounded, i.e. that there exists some $b > 0$ such that $p_m(\theta) \leq b$ for all $\theta \in \mathbb{R}^d$ and $m \in \{1, \dots, M\}$.

First, we bound the bias of our estimator. This shows that the bias tends to zero as the bandwidth shrinks.

Lemma 5.1. *The bias of the estimator $p_1 \widehat{\cdots} p_M(\theta)$ satisfies*

$$\sup_{p_1, \dots, p_M \in \mathcal{P}(\beta, L)} |\mathbb{E}[p_1 \widehat{\cdots} p_M(\theta)] - p_1 \cdots p_M(\theta)| \leq \sum_{m=1}^M c_m h^{m\beta}$$

for some $c_1, \dots, c_M > 0$.

Proof. For all $p_1, \dots, p_M \in \mathcal{P}(\beta, L)$,

$$\begin{aligned} |\mathbb{E}[p_1 \widehat{\cdots} p_M] - p_1 \cdots p_M| &= |\mathbb{E}[\widehat{p}_1 \cdots \widehat{p}_M] - p_1 \cdots p_M| \\ &= |\mathbb{E}[\widehat{p}_1] \cdots \mathbb{E}[\widehat{p}_M] - p_1 \cdots p_M| \\ &\leq |(p_1 + \tilde{c}_1 h^\beta) \cdots (p_M + \tilde{c}_M h^\beta) - p_1 \cdots p_M| \\ &\leq |c_1 h^\beta + \dots + c_M h^{M\beta}| \\ &\leq |c_1 h^\beta| + \dots + |c_M h^{M\beta}| \\ &= \sum_{m=1}^M c_m h^{m\beta} \end{aligned}$$

where we have used the fact that $|\mathbb{E}[\widehat{p}_m] - p| \leq \tilde{c}_m h^\beta$ for some $\tilde{c}_m > 0$. \square

Next, we bound the variance of our estimator. This shows that the variance tends to zero as the number of samples grows large and the bandwidth shrinks.

Lemma 5.2. *The variance of the estimator $p_1 \widehat{\cdots} p_M(\theta)$ satisfies*

$$\sup_{p_1, \dots, p_M \in \mathcal{P}(\beta, L)} \mathbb{V}[p_1 \widehat{\cdots} p_M(\theta)] \leq \sum_{m=1}^M \binom{M}{m} \frac{c_m}{T^m h^{dm}}$$

for some $c_1, \dots, c_M > 0$ and $0 < h \leq 1$.

Proof. For all $p_1, \dots, p_M \in \mathcal{P}(\beta, L)$,

$$\begin{aligned} \mathbb{V}[p_1 \widehat{\cdots} p_M] &= \mathbb{E}[\widehat{p}_1^2] \cdots \mathbb{E}[\widehat{p}_M^2] - \mathbb{E}[\widehat{p}_1]^2 \cdots \mathbb{E}[\widehat{p}_M]^2 \\ &= \left(\prod_{m=1}^M \mathbb{V}[\widehat{p}_m] + \mathbb{E}[\widehat{p}_m]^2 \right) - \left(\prod_{m=1}^M \mathbb{E}[\widehat{p}_m]^2 \right) \\ &\leq \sum_{m=0}^{M-1} \binom{M}{m} \frac{\tilde{c}^m c^{M-m}}{T^{M-m} h^{d(M-m)}} \\ &\leq \sum_{m=1}^M \binom{M}{m} \frac{c_m}{T^m h^{dm}} \end{aligned}$$

where we have used the facts that $\mathbb{V}[\widehat{p}_m] \leq \frac{c}{T h^d}$ for some $c > 0$ and $\mathbb{E}[\widehat{p}_m]^2 \leq \tilde{c}$ for some $\tilde{c} > 0$. \square

Finally, we use the bias and variance bounds to bound the MSE, which shows that our estimator is consistent.

Theorem 5.3. *If $h \asymp T^{-1/(2\beta+d)}$, the mean-squared error of the estimator $\widehat{p_1 \cdots p_M}(\theta)$ satisfies*

$$\begin{aligned} \sup_{p_1, \dots, p_M \in \mathcal{P}(\beta, L)} \mathbb{E} \left[\int (\widehat{p_1 \cdots p_M}(\theta) - p_1 \cdots p_M(\theta))^2 d\theta \right] \\ \leq \frac{c}{T^{2\beta/(2\beta+d)}} \end{aligned}$$

for some $c > 0$ and $0 < h \leq 1$.

Proof. For all $p_1, \dots, p_M \in \mathcal{P}(\beta, L)$, using the fact that the mean-squared error is equal to the variance plus the bias squared, we have that

$$\begin{aligned} \mathbb{E} \left[\int (\widehat{p_1 \cdots p_M}(\theta) - p_1 \cdots p_M(\theta))^2 d\theta \right] \\ \leq \left(\sum_{m=1}^M c_m h^{m\beta} \right)^2 + \sum_{m=1}^M \binom{M}{m} \frac{\tilde{c}_m}{T^m h^{dm}} \\ \leq k T^{-2\beta/(2\beta+d)} + \frac{\tilde{k}}{T^{1-d(2\beta+d)}} \quad (\text{for some } k, \tilde{k} > 0) \\ \leq \frac{c}{T^{2\beta/(2\beta+d)}} \end{aligned}$$

for some $c_1, \dots, c_M > 0$ and $\tilde{c}_1, \dots, \tilde{c}_M > 0$. \square

6 Method Scope

The theoretical results and algorithms in this paper hold for posterior distributions over finite-dimensional real spaces. These include generalized linear models (e.g. linear, logistic, or Poisson regression), mixture models with known weights, hierarchical models, and (more generally) finite-dimensional graphical models with unconstrained variables. This also includes both unimodal and multimodal posterior densities (such as in Section 8.2). However, the methods and theory presented here do not yet extend to cases such as infinite dimensional models (e.g. nonparametric Bayesian models [5]) nor to distributions over the simplex (e.g. topics in latent Dirichlet allocation [3]). In the future, we hope to extend this work to these domains.

7 Related Work

In [19, 2, 16], the authors develop a way to sample approximately from a posterior distribution when only a small randomized mini-batch of data is used at each step. In [9], the authors used a hypothesis test to decide whether to accept or reject proposals using a

small set of data (adaptively) as opposed to the exact Metropolis-Hastings rule. This reduces the amount of time required to compute the acceptance ratio. Since all of these algorithms are still sequential, they can be directly used in our algorithm to generate subposterior samples to further speed up the entire sampling process.

Several parallel MCMC algorithms have been designed for specific models, such as for topic models [18, 14] and nonparametric mixture models [21]. These approaches still require synchronization to be correct (or approximately correct), while ours aims for more general model settings and does not need synchronization until the final combination stage.

Consensus Monte Carlo [17] is perhaps the most relevant work to ours. In this algorithm, data is also partitioned into different machines and MCMC is performed independently on each machine. Thus, it roughly has the same time complexity as our algorithm. However, the prior is not explicitly reweighted during sampling as we do in Eq 1, and final samples for the full posterior are generated by averaging subposterior samples. Furthermore, this algorithm has few theoretical guarantees. We find that this algorithm can be viewed as a relaxation of our nonparametric, asymptotically exact sampling procedure, where samples are generated from an evenly weighted mixture (instead of each component having weight w_t) and where each sample is set to $\hat{\theta}_t$ instead of being drawn from $\mathcal{N}(\hat{\theta}_t, \frac{h}{M} I_d)$. This algorithm is one of our experimental baselines.

8 Empirical Study

In the following sections, we demonstrate empirically that our method allows for quicker, MCMC-based estimation of a posterior distribution, and that our consistent-estimator-based procedures yield asymptotically exact results. We show our method on a few Bayesian models using both synthetic and real data. In each experiment, we compare the following strategies for parallel, communication-free sampling:⁴

- **Single chain full-data posterior samples (regularChain)**—Typical, single-chain MCMC for sampling from the full-data posterior.
- **Parametric subposterior density product estimate (parametric)**—For M sets of subposterior samples, the combination yielding samples from the parametric density product estimate.
- **Nonparametric subposterior density product estimate (nonparametric)**—For M sets of

⁴We did not directly compare with the algorithms that require synchronization since the setup of these experiments can be rather different. We plan to explore these comparisons in the extended version of this paper.

subposterior samples, the combination yielding samples from the nonparametric density product estimate.

- **Semiparametric subposterior density product estimate** (`semiparametric`)—For M sets of subposterior samples, the combination yielding samples from the semiparametric density product estimate.
- **Subposterior sample average** (`subpostAvg`)—For M sets of subposterior samples, the average of M samples consisting of one sample taken from each subposterior.
- **Subposterior sample pooling** (`subpostPool`)—For M sets of subposterior samples, the union of all sets of samples.
- **Duplicate chains full-data posterior sample pooling** (`duplicateChainsPool`)—For M sets of samples from the full-data posterior, the union of all sets of samples.

To assess the performance of our sampling and combination strategies, we ran a single chain of MCMC on the full data for 500,000 iterations, removed the first half as burn-in, and considered the remaining samples the “groundtruth” samples for the true posterior density. We then needed a general method to compare the distance between two densities given samples from each, which holds for general densities (including multimodal densities, where it is ineffective to compare moments such as the mean and variance⁵). Following work in density-based regression [15], we use an estimate of the L_2 distance, $d_2(p, \hat{p})$, between the groundtruth posterior density p and a proposed posterior density \hat{p} , where $d_2(p, \hat{p}) = \|p - \hat{p}\|_2 = (\int (p(\theta) - \hat{p}(\theta))^2 d\theta)^{1/2}$.

In the following experiments involving timing, to compute the posterior L_2 error at each time point, we collected all samples generated before a given number of seconds, and added the time taken to transfer the samples and combine them using one of the proposed methods. In all experiments and methods, we followed a fixed rule of removing the first $\frac{1}{6}$ samples for burn-in (which, in the case of combination procedures, was applied to each set of subposterior samples before the combination was performed).

Experiments were conducted with a standard cluster system. We obtained subposterior samples by submitting batch jobs to each worker since these jobs are all independent. We then saved the results to the disk of each worker and transferred them to the same machine which performed the final combination.

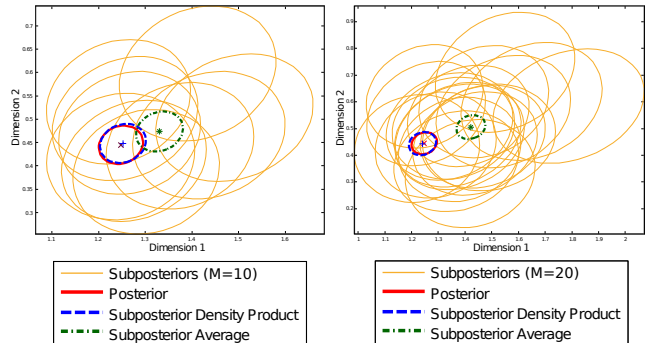


Figure 1: Bayesian logistic regression posterior ovals. We show the posterior 90% probability mass ovals for the first 2-dimensional marginal of the posterior, the M subposteriors, the subposterior density product (via the parametric procedure), and the subposterior average (via the `subpostAvg` procedure). We show $M=10$ subsets (left) and $M=20$ subsets (right). The subposterior density product generates samples that are consistent with the true posterior, while the `subpostAvg` produces biased results, which grow in error as M increases.

8.1 Generalized Linear Models

Generalized linear models are widely used for many regression and classification problems. Here we conduct experiments, using logistic regression as a test case, on both synthetic and real data to demonstrate the speed of our parallel MCMC algorithm compared with typical MCMC strategies.

8.1.1 Synthetic data

Our synthetic dataset contains 50,000 observations in 50 dimensions. To generate the data, we drew each element of the model parameter β and data matrix X from a standard normal distribution, and then drew each outcome as $y_i \sim \text{Bernoulli}(\text{logit}^{-1}(X_i \beta))$ (where X_i denotes the i^{th} row of X)⁶. We use Stan, an automated Hamiltonian Monte Carlo (HMC) software package,⁷ to perform sampling for both the true posterior (for groundtruth and comparison methods) and for the subposteriors on each machine. One advantage of Stan is that it is implemented with C++ and uses the No-U-Turn sampler for HMC, which does not require any user-provided parameters [8].

In Figure 1, we illustrate results for logistic regression, showing the subposterior densities, the subposterior density product, the subposterior sample average, and the true posterior density, for the number of subsets M set to 10 (left) and 20 (right). Samples generated by our approach (where we draw samples from the

⁵In these cases, dissimilar densities might have similar low-order moments. See Section 8.2 for an example.

⁶Note that we did not explicitly include the intercept term in our logistic regression model.

⁷<http://mc-stan.org>

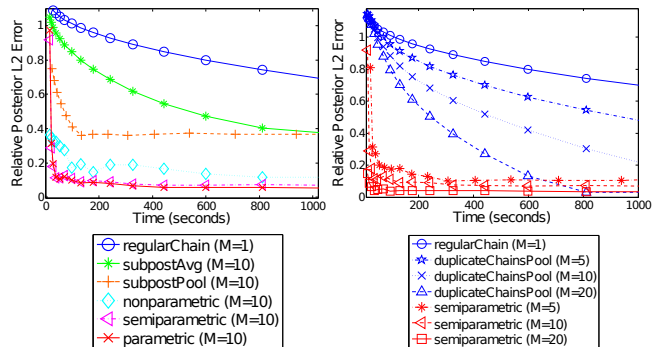


Figure 2: Posterior L_2 error vs time for logistic regression. Left: the three combination strategies proposed in this paper (**parametric**, **nonparametric**, and **semiparametric**) reduce the posterior error much more quickly than a single full-data Markov chain; the **subpostAvg** and **subpostPool** procedures yield biased results. Right: we compare with multiple full-data Markov chains (**duplicateChainsPool**); our method yields faster convergence to the posterior even though only a fraction of the data is being used by each chain.

subposterior density product via the **parametric** procedure) overlap with the true posterior much better than those generated via the **subpostAvg** (subposterior sample average) procedure—averaging of samples appears to create systematic biases. Further, the error in averaging appears to increase as M grows. In Figure 2 (left) we show the posterior error vs time. A regular full-data chain takes much longer to converge to low error compared with our combination methods, and simple averaging and pooling of subposterior samples gives biased solutions.

We next compare our combination methods with multiple independent “duplicate” chains each run on the full dataset. Even though our methods only require a fraction of the data storage on each machine, we are still able to achieve a significant speed-up over the full-data chains. This is primarily because the duplicate chains cannot parallelize burn-in (i.e. each chain must still take some n steps before generating reasonable samples, and the time taken to reach these n steps does not decrease as more machines are added). However, in our method, each subposterior sampler can take each step more quickly, effectively allowing us to decrease the time needed for burn-in as we increase M . We show this empirically in Figure 2 (right), where we plot the posterior error vs time, and compare with full duplicate chains as M is increased.

Using a Matlab implementation of our combination algorithms, all (batch) combination procedures take under twenty seconds to complete on a 2.5GHz Intel Core i5 with 16GB memory.

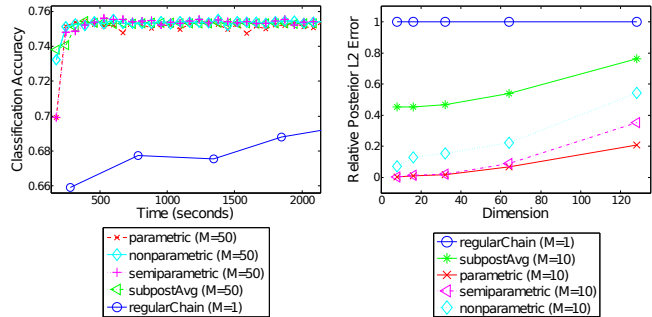


Figure 3: Left: Bayesian logistic regression classification accuracy vs time for the task of predicting forest cover type. Right: Posterior error vs dimension on synthetic data at 1000 seconds, normalized so that **regularChain** error is fixed at 1.

8.1.2 Real-world data

Here, we use the *covtype* (predicting forest cover types)⁸ dataset, containing 581,012 observations in 54 dimensions. A single chain of HMC running on this entire dataset takes an average of 15.76 minutes per sample; hence, it is infeasible to generate groundtruth samples for this dataset. Instead we show classification accuracy vs time. For a given set of samples, we perform classification using a sample estimate of the posterior predictive distribution for a new label y with associated datapoint x , i.e.

$$P(y|x, y^N, x^N) = \int P(y|x, \beta, y^N, x^N) P(\beta|x^N, y^N) \\ \approx \frac{1}{S} \sum_{s=1}^S P(y|x, \beta_s)$$

where x^N and y^N denote the N observations, and $P(y|x, \beta_s) = \text{Bernoulli}(\text{logit}^{-1}(x^\top \beta_s))$. Figure 3 (left) shows the results for this task, where we use $M=50$ splits. The parallel methods achieve a higher accuracy much faster than the single-chain MCMC algorithm.

8.1.3 Scalability with dimension

We investigate how the errors of our methods scale with dimensionality, to compare the different estimators implicit in the combination procedures. In Figure 3 (right) we show the relative posterior error (taken at 1000 seconds) vs dimension, for the synthetic data with $M=10$ splits. The errors at each dimension are normalized so that the **regularChain** error is equal to 1. Here, the **parametric** (asymptotically biased) procedure scales best with dimension, and the **semiparametric** (asymptotically exact) procedure is a close second. These results also demonstrate that, although the **nonparametric** method can be viewed as

⁸<http://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets>

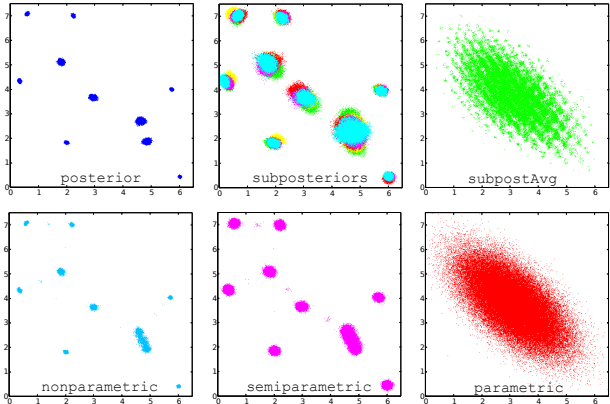


Figure 4: Gaussian mixture model posterior samples. We show 100,000 samples from a single 2-d marginal (corresponding to the posterior over a single mean parameter) of the full-data posterior (top left), all subposteriors (top middle—each one is given a unique color), the subposterior average via the `subpostAvg` procedure (top right), and the subposterior density product via the `nonparametric` procedure (bottom left), `semiparametric` procedure (bottom middle), and `parametric` procedure (bottom right).

implicitly sampling from a nonparametric density estimate (which is usually restricted to low-dimensional densities), the performance of our method does not suffer greatly when we perform parallel MCMC on posterior distributions in much higher-dimensional spaces.

8.2 Gaussian mixture models

In this experiment, we aim to show correct posterior sampling in cases where the full-data posterior, as well as the subposteriors, are multimodal. We will see that the combination procedures that are asymptotically biased suffer greatly in these scenarios. To demonstrate this, we perform sampling in a Gaussian mixture model. We generate 50,000 samples from a ten component mixture of 2-d Gaussians. The resulting posterior is multimodal; this can be seen by the fact that the component labels can be arbitrarily permuted and achieve the same posterior value. For example, we find after sampling that the posterior distribution over each component mean has ten modes. To sample from this multimodal posterior, we used the Metropolis-Hastings algorithm, where the component labels were permuted before each step (note that this permutation results in a move between two points in the posterior distribution with equal probability).

In Figure 4 we show results for $M=10$ splits, showing samples from the true posterior, overlaid samples from all five subposteriors, results from averaging the subposterior samples, and the results after applying our three subposterior combination procedures. This figure

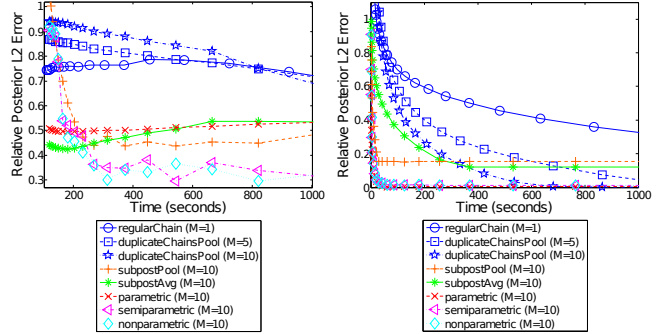


Figure 5: Left: Gaussian mixture model posterior error vs time results. Right: Poisson-gamma hierarchical model posterior error vs time results.

shows the 2-d marginal of the posterior corresponding to the posterior over a single mean component. The `subpostAvg` and `parametric` procedures both give biased results, and cannot capture the multimodality of the posterior. We show the posterior error vs time in Figure 5 (left), and see that our asymptotically exact methods yield quick convergence to low posterior error.

8.3 Hierarchical models

We show results on a hierarchical Poisson-gamma model of the following form

$$a \sim \text{Exponential}(\lambda) \quad b \sim \text{Gamma}(\alpha, \beta) \\ q_i \sim \text{Gamma}(a, b) \quad x_i \sim \text{Poisson}(q_i t_i) \quad i = 1, \dots, N$$

for $N=50,000$ observations. We draw $\{x_i\}_{i=1}^N$ from the above generative process (after fixing values for a, b, λ , and $\{t_i\}_{i=1}^N$), and use $M=10$ splits. We again perform MCMC using the Stan software package.

In Figure 5 (right) we show the posterior error vs time, and see that our combination methods complete burn-in and converge to a low posterior error very quickly relative to the `subpostAvg` and `subpostPool` procedures and full-data chains.

9 Discussion and Future Work

In this paper, we present an embarrassingly parallel MCMC algorithm and provide theoretical guarantees about the samples it yields. Experimental results demonstrate our method’s potential to speed up burn-in and perform faster asymptotically correct sampling. Further, it can be used in settings where data are partitioned onto multiple machines that have little intercommunication—this is ideal for use in a MapReduce setting. Currently, our algorithm works primarily when the posterior samples are real, unconstrained values and we plan to extend our algorithm to more general settings in future work.

References

- [1] Alekh Agarwal and John C Duchi, *Distributed delayed stochastic optimization*, Decision and Control (CDC), 2012 IEEE 51st Annual Conference on, IEEE, 2012, pp. 5451–5452.
- [2] Sungjin Ahn, Anoop Korattikara, and Max Welling, *Bayesian posterior sampling via stochastic gradient fisher scoring*, Proceedings of the 29th International Conference on Machine Learning, 2012, pp. 1591–1598.
- [3] David M Blei, Andrew Y Ng, and Michael I Jordan, *Latent dirichlet allocation*, The Journal of Machine Learning Research **3** (2003), 993–1022.
- [4] Jeffrey Dean and Sanjay Ghemawat, *Mapreduce: simplified data processing on large clusters*, Communications of the ACM **51** (2008), no. 1, 107–113.
- [5] Samuel J Gershman and David M Blei, *A tutorial on bayesian nonparametric models*, Journal of Mathematical Psychology **56** (2012), no. 1, 1–12.
- [6] Nils Lid Hjort and Ingrid K Glad, *Nonparametric density estimation with a parametric start*, The Annals of Statistics (1995), 882–904.
- [7] Qirong Ho, James Cipar, Henggang Cui, Seunghak Lee, Jin Kyu Kim, Phillip B. Gibbons, Gregory R. Ganger, Garth Gibson, and Eric P. Xing, *More effective distributed ml via a stale synchronous parallel parameter server*, Advances in Neural Information Processing Systems, 2013.
- [8] Matthew D Hoffman and Andrew Gelman, *The no-urn sampler: Adaptively setting path lengths in hamiltonian monte carlo*, arXiv preprint arXiv:1111.4246 (2011).
- [9] Anoop Korattikara, Yutian Chen, and Max Welling, *Austerity in MCMC land: Cutting the Metropolis-Hastings budget*, arXiv preprint arXiv:1304.5299 (2013).
- [10] John Langford, Alex J Smola, and Martin Zinkevich, *Slow learners are fast*, Advances in Neural Information Processing Systems, 2009.
- [11] Kathryn Blackmond Laskey and James W Myers, *Population Markov chain Monte Carlo*, Machine Learning **50** (2003), no. 1-2, 175–196.
- [12] Lucien Le Cam, *Asymptotic methods in statistical decision theory*, New York (1986).
- [13] Lawrence Murray, *Distributed Markov chain Monte Carlo*, Proceedings of Neural Information Processing Systems Workshop on Learning on Cores, Clusters and Clouds, vol. 11, 2010.
- [14] David Newman, Arthur Asuncion, Padhraic Smyth, and Max Welling, *Distributed algorithms for topic models*, The Journal of Machine Learning Research **10** (2009), 1801–1828.
- [15] Junier Oliva, Barnabás Póczos, and Jeff Schneider, *Distribution to distribution regression*, Proceedings of The 30th International Conference on Machine Learning, 2013, pp. 1049–1057.
- [16] Sam Patterson and Yee Whye Teh, *Stochastic gradient riemannian langevin dynamics on the probability simplex*, Advances in Neural Information Processing Systems, 2013.
- [17] Steven L. Scott, Alexander W. Blocker, and Fernando V. Bonassi, *Bayes and big data: The consensus monte carlo algorithm*, Bayes 250, 2013.
- [18] Alexander Smola and Shравan Narayanamurthy, *An architecture for parallel topic models*, Proceedings of the VLDB Endowment **3** (2010), no. 1-2, 703–710.
- [19] Max Welling and Yee W Teh, *Bayesian learning via stochastic gradient Langevin dynamics*, Proceedings of the 28th International Conference on Machine Learning, 2011, pp. 681–688.
- [20] Darren J Wilkinson, *Parallel Bayesian computation*, Statistics Textbooks and Monographs **184** (2006), 477.
- [21] Sinead Williamson, Avinava Dubey, and Eric P Xing, *Parallel Markov chain Monte Carlo for nonparametric mixture models*, Proceedings of the 30th International Conference on Machine Learning, 2013, pp. 98–106.