

# Firefly Monte Carlo: Exact MCMC with Subsets of Data

**Dougal Maclaurin**

Department of Physics  
Harvard University  
Cambridge, MA 02138

**Ryan P. Adams**

School of Engineering and Applied Sciences  
Harvard University  
Cambridge, MA 02138

## Abstract

Markov chain Monte Carlo (MCMC) is a popular tool for Bayesian inference. However, MCMC cannot be practically applied to large data sets because of the prohibitive cost of evaluating every likelihood term at every iteration. Here we present *Firefly Monte Carlo* (FlyMC) MCMC algorithm with auxiliary variables that only queries the likelihoods of a subset of the data at each iteration yet simulates from the exact posterior distribution. FlyMC is compatible with modern MCMC algorithms, and only requires a lower bound on the per-datum likelihood factors. In experiments, we find that FlyMC generates samples from the posterior more than an order of magnitude faster than regular MCMC, allowing MCMC methods to tackle larger datasets than were previously considered feasible.

## 1 Introduction

The Bayesian approach to probabilistic modeling is appealing for several reasons: the generative framework allows one to separate out modeling assumptions from inference procedures, outputs include estimates of uncertainty that can be used for decision making and prediction, and it provides clear ways to perform model selection and complexity control. Unfortunately, the fully-Bayesian approach to modeling is often very computationally challenging. It is unusual for non-trivial models of real data to have closed-form posterior distributions. Instead, one uses approximate inference via Monte Carlo, variational approximations, Laplace approximations, or other tools.

Bayesian inference can be costly. In particular, evaluating a new hypothesis usually requires examination of an entire data. For example, when running Metropolis–Hastings (MH) it is necessary to evaluate the target posterior density for each proposed parameter update, and this posterior will usually contain a factor for each datum. Similarly, typical variational Bayesian procedures need to build local approximations for each of the data in order to update the approximation to any global parameters. In both cases, these data-intensive computations are performed many times as part of an iterative procedure.

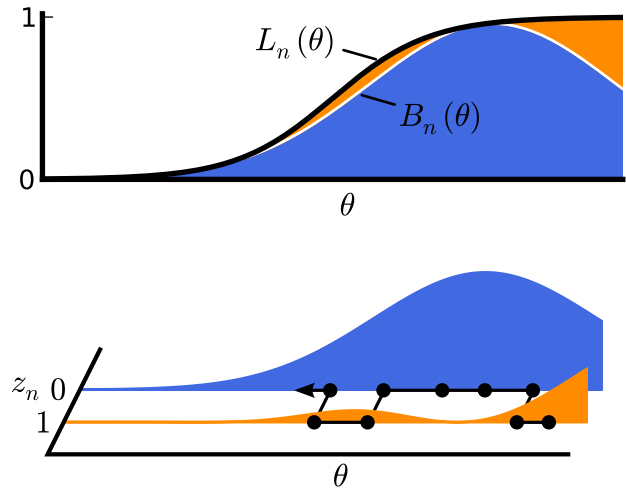


Figure 1: Illustration of the auxiliary variable representation of a single likelihood for a one-dimensional logistic regression model. The top panel shows how the likelihood function,  $L_n(\theta)$ , corresponding to a single datum  $n$ , can be partitioned into two parts: a lower bound,  $B_n(\theta)$ , shaded blue, and the remainder, shaded orange. The bottom panel shows that we can introduce a Bernoulli random variable  $z_n$  and construct a Markov chain in this new, higher dimensional space, such that marginalizing out (i.e. ignoring) the  $z_n$  recovers the original likelihood. If  $B_n(\theta) \gg L_n(\theta) - B_n(\theta)$ , the Markov chain will tend to occupy  $z_n = 0$  and we can avoid evaluating  $L_n(\theta)$  at each iteration.

Recent methods have been proposed to partially overcome these difficulties. Stochastic and online variational approximation procedures [Hoffman *et al.*, 2010, 2013] can use subsets of the data to make approximations to global parameters. As these procedures are optimizations, it is possible to build on convergence results from the stochastic optimization literature and achieve guarantees on the resulting approximation. For Markov chain Monte Carlo, the situation is somewhat murkier. Recent work has shown that approximate transition operators based on subsets of data can be used for predic-

tive prefetching to help parallelize MCMC [Angelino *et al.*, 2014]. Other work uses approximate transition operators directly, for Metropolis-Hastings (MH) and related algorithms [Welling and Teh, 2011]. Korattikara *et al.* [2014] and Bardenet *et al.* [2014] have shown that such approximate MH moves can lead to stationary distributions which are approximate but that have bounded error, albeit under strong conditions of rapid mixing.

In this paper we present a Markov chain Monte Carlo algorithm, *Firefly Monte Carlo* (FlyMC), that is in line with these latter efforts to exploit subsets of data. Our transition operator, however, is *exact*: it leaves the true posterior distribution invariant. FlyMC is a latent variable model which introduces a collection of Bernoulli variables – one for each datum – with conditional distributions chosen so that they effectively turn on and off data points in the posterior, hence “firefly”. The introduction of these latent variables does not alter the marginal distribution of the parameters of interest. Our only requirement is that it be possible to provide a “collapsible” lower bound for each likelihood term. FlyMC can lead to dramatic performance improvements in MCMC, as measured in wallclock time.

## 2 Firefly Monte Carlo

The Firefly Monte Carlo algorithm tackles the problem of sampling from the posterior distribution of a probabilistic model. We denote the parameters of interest as  $\theta$  and assume that they have prior  $p(\theta)$ . We assume that  $N$  data have been observed  $\{x_n\}_{n=1}^N$  and that these data are conditionally independent given  $\theta$  under a likelihood  $p(x_n | \theta)$ . Our target distribution is therefore

$$p(\theta | \{x_n\}_{n=1}^N) \propto p(\theta, \{x_n\}_{n=1}^N) = p(\theta) \prod_{n=1}^N p(x_n | \theta). \quad (1)$$

We will write the  $n$ th likelihood term as a function of  $\theta$  as

$$L_n(\theta) = p(x_n | \theta).$$

An MCMC sampler makes transitions from a given  $\theta$  to a new  $\theta'$  such that posterior distribution remains invariant. Conventional algorithms, such as Metropolis–Hastings, require evaluation of the unnormalized posterior in full at every iteration. When the data set is large, evaluating all  $N$  likelihoods is a computational bottleneck. This is the problem that we seek to solve with FlyMC.

For each data point,  $n$ , we introduce a binary auxiliary variable,  $z_n \in \{0, 1\}$ , and a function  $B_n(\theta)$  which is a strictly positive lower bound on the  $n$ th likelihood:  $0 < B_n(\theta) \leq L_n(\theta)$ . Each  $z_n$  has the following Bernoulli distribution conditioned on the parameters:

$$p(z_n | x_n, \theta) = \left[ \frac{L_n(\theta) - B_n(\theta)}{L_n(\theta)} \right]^{z_n} \left[ \frac{B_n(\theta)}{L_n(\theta)} \right]^{1-z_n}.$$

We now augment the posterior distribution with these  $N$  variables:

$$\begin{aligned} p(\theta, \{z_n\}_{n=1}^N | \{x_n\}_{n=1}^N) &\propto p(\theta, \{x_n, z_n\}_{n=1}^N) \\ &= p(\theta) \prod_{n=1}^N p(x_n | \theta) p(z_n | x_n, \theta). \end{aligned}$$

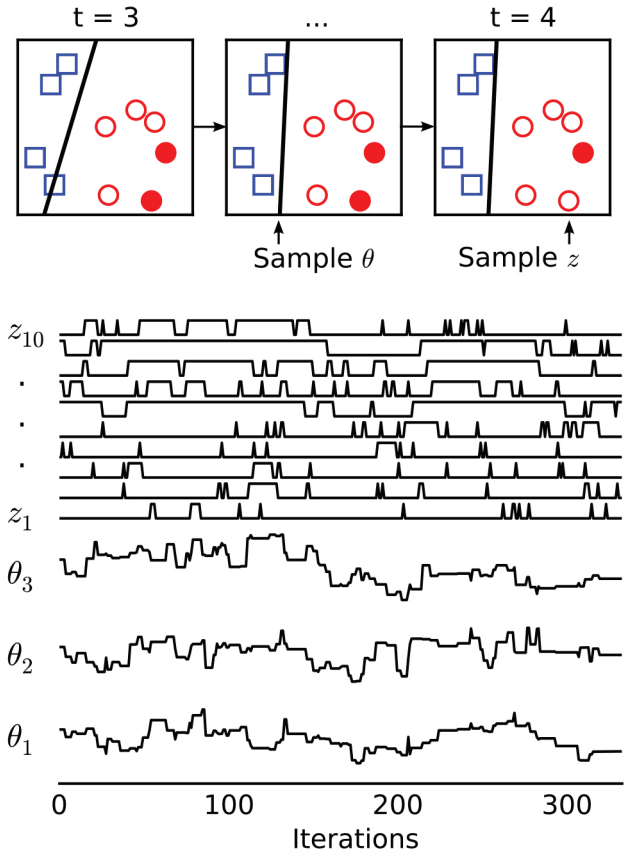


Figure 2: Illustration of the FlyMC algorithm operating on a logistic regression model of a toy synthetic data set, a two-class classification problem in two dimensions (and one bias dimension). The top panel shows a single iteration of FlyMC, from  $t = 3$  to  $t = 4$ , which consists of two steps: first we sample  $\theta$ , represented by the line of equal class probability. Next we sample the  $z_n$ . In this case, we see one ‘bright’ (solid) data point become dark. The bottom panel shows the trajectories of all components of  $\theta$  and  $z$ .

This joint distribution has a remarkable property: to evaluate the probability density over  $\theta$ , conditioned on  $\{z_n\}_{n=1}^N$ , it is only necessary to evaluate those likelihood terms for which  $z_n = 1$ . Consider factor  $n$  from the product above.

$$p(x_n | \theta) p(z_n | x_n, \theta) = \begin{cases} L_n(\theta) - B_n(\theta) & \text{if } z_n = 1 \\ B_n(\theta) & \text{if } z_n = 0 \end{cases}$$

The likelihood term  $L_n(\theta)$  only appears in those factors for which  $z_n = 1$  and we can think of these data as forming a minibatch subsample of the full set. If  $z_n = 0$  for most  $n$ , transition updates for the parameters  $\theta$  will be much cheaper, as these are applied to  $p(\theta | \{x_n, z_n\}_{n=1}^N)$ .

Of course, we do have to evaluate all  $N$  bounds  $B_n(\theta)$  at each iteration. We seem to have just shifted the computational burden from evaluating the  $L_n(\theta)$  to evaluating the  $B_n(\theta)$ . However, if we choose  $B_n(\theta)$  to have a convenient form such as a scaled Gaussian, then the full product  $\prod_{n=1}^N B_n(\theta)$  can

be computed for each new  $\theta$  in  $O(1)$  time using the sufficient statistics of the distribution, which only need to be computed once. To make this clearer, we can rearrange the joint distribution in terms of a “pseudo-prior,”  $\tilde{p}(\theta)$  and “pseudo-likelihood,”  $\tilde{L}_n(\theta)$  as follows:

$$p(\theta, \{z_n\}_{n=1}^N | \{x_n\}_{n=1}^N) \propto \tilde{p}(\theta) \prod_{n:z_n=1} \tilde{L}_n(\theta) \quad (2)$$

where the product only runs over those  $n$  for which  $z_n = 1$ , and we have defined

$$\tilde{p}(\theta) = p(\theta) \prod_{n=1}^N B_n(\theta) \quad \tilde{L}_n(\theta) = \frac{L_n(\theta) - B_n(\theta)}{B_n(\theta)}.$$

We can generate a Markov chain for the joint distribution in Equation (2) by alternating between updates of  $\theta$  conditional on  $\{z_n\}_{n=1}^N$ , which can be done with any conventional MCMC algorithm, and updates of  $\{z_n\}_{n=1}^N$  conditional on  $\theta$  for which we discuss efficient methods in Section 3.2. We emphasize that the marginal distribution over  $\theta$  is still the correct posterior distribution given in Equation (1).

At a given iteration, the  $z_n = 0$  data points are “dark”: we simulate the Markov chain without computing their likelihoods. Upon a Markov transition in the space of  $\{z_n\}_{n=1}^N$ , a smattering of these dark data points become “bright” with their  $z_n = 1$ , and we include their likelihoods in subsequent iterations. The evolution of the chain evokes an image of fireflies, as the individual data blink on and off.

The details of choosing a lower bound and efficiently sampling the  $\{z_n\}$  are treated in the proceeding sections, but the high-level picture is now complete. Figure 1 illustrates the augmented space, and a simple version of the algorithm is shown in Algorithm 1. Figure 2 shows several steps of Firefly Monte Carlo on a toy logistic regression model.

### 3 Implementation Considerations

In this section we discuss two practical matters for implementing an effective FlyMC algorithm: how to choose and compute lower bounds and how to sample the brightness variables  $z_n$ . We assume a data set with  $N$  data points and a parameter vector,  $\theta$ , of dimension  $D \ll N$ . We also assume that it takes at least  $O(ND)$  time to evaluate the likelihoods at some  $\theta$  for the whole data set and that evaluating this set of likelihoods at each iteration is the computational bottleneck for MCMC.

The goal of an effective implementation of FlyMC is to construct a Markov chain with similar convergence and mixing properties to that of regular MCMC, while only evaluating a subset of the data points on average at each iteration. If the average number of “bright” data points is  $M$ , we would like this to achieve a computational speedup of nearly  $N/M$  over regular MCMC.

#### 3.1 Choosing a lower bound

The lower bounds,  $B_n(\theta)$  of each data point’s likelihood  $L_n(\theta)$  should satisfy two properties. They should be relatively tight, and it should be possible to efficiently summarize a product of lower bounds  $\prod_n B_n(\theta)$  in time independent of  $N$  (after setup).

The tightness of the bounds is important because it determines the number of bright data points at each iteration, which determines the time it takes to evaluate the joint posterior. For a burned-in chain, the average number of bright data points,  $M$ , will be:

$$M = \sum_{n=1}^N \langle z_n \rangle = \sum_{n=1}^N \int p(\theta | \{x_n\}_{n=1}^N) \frac{L_n(\theta) - B_n(\theta)}{L_n(\theta)} d\theta.$$

Therefore it is important that the bounds are tight at values of  $\theta$  where the posterior puts the bulk of its mass.

The second required property is that the product of the lower bounds must be easy to compute and represent. This emerges naturally if we use scaled exponential-family lower bounds so that their product can be summarized with a set of sufficient statistics. The individual bounds  $B_n(\theta)$  should be easy to compute themselves, since these are computed alongside  $L_n(\theta)$  for all the bright points at each iteration.

Let’s consider a concrete example. The logistic regression likelihood is

$$L_n(\theta) = \text{logit}^{-1}(t_n \theta^\top x_n) = \frac{1}{1 + \exp\{-t_n \theta^\top x_n\}},$$

where  $x_n \in \mathbb{R}^D$  is the set of features for the  $n$ th data point and  $t_n \in \{-1, 1\}$  is its class. The logistic function has a family of scaled Gaussian lower bounds, described in Jaakkola and Jordan [1997], parameterized by  $\xi$ , the location at which the bound is tight:

$$\log(B_n(\theta)) = a(t_n \theta^\top x_n)^2 + b(t_n \theta^\top x_n) + c$$

where:

$$a = \frac{-1}{4\xi} \left( \frac{e^\xi - 1}{e^\xi + 1} \right), \quad b = \frac{1}{2}, \quad c = -a\xi^2 + \frac{\xi}{2} - \log(e^\xi + 1)$$

This is the bound shown in Fig. 1. The product of these bounds can be computed for a given  $\theta$  in  $O(D^2)$  time, provided we have precomputed the moments of the data, at a one-time setup cost of  $O(ND^2)$ .

This bound can be quite tight. For example, if we choose  $\xi = 1.5$  the probability of a data point being bright is less than 0.02 in the region where  $0.1 < L_n(\theta) < 0.9$ . With a bit of up-front work, we can do even better than this by choosing bounds that are tight in the right places. For example, we can perform a quick optimization to find an approximate maximum *a posteriori* (MAP) value of  $\theta$  and construct the bounds to be tight there. We explore this idea further in Section 4.

#### 3.2 Sampling and handling the auxiliary brightness variables

The resampling of the  $z_n$  variables, as shown in lines 3 to 5 of Algorithm 1, takes a step by explicitly sampling  $z_n$  from its conditional distribution for a random fixed-size subset of the data. We call this approach *explicit resampling* and it has a clear drawback: if the fixed fraction is  $\alpha$  (shown as RESAMPLEFRACTION in Algorithm 1), then the chain cannot have a mixing time faster than  $1/\alpha$ , as each data point is only visited a fraction of the time.

Nevertheless, explicit resampling works well in practice since the bottleneck for mixing is usually the exploration of  $\theta$ ,

---

```

1:  $\theta_0 \sim \text{INITIALDIST}$  ▷ Initialize the Markov chain state.
2: for  $i \leftarrow 1 \dots \text{ITERS}$  do ▷ Iterate the Markov chain.
3:   for  $j \leftarrow 1 \dots \lceil N \times \text{RESAMPLEFRACTION} \rceil$  do
4:      $n \sim \text{RandInteger}(1, N)$  ▷ Select a random data point.
5:      $z_n \sim \text{Bernoulli}(1 - B_n(\theta_{i-1})/L_n(\theta_{i-1}))$  ▷ Biased coin-flip to determine whether  $n$  is bright or dark.
6:      $\theta' \leftarrow \theta_{i-1} + \eta$  where  $\eta \sim \text{Normal}(0, \epsilon^2 \mathbb{I}_D)$  ▷ Make a random walk proposal with step size  $\epsilon$ .
7:      $u \sim \text{Uniform}(0, 1)$  ▷ Draw the MH threshold.
8:     if  $\frac{\text{JOINTPOSTERIOR}(\theta'; \{z_n\}_{n=1}^N)}{\text{JOINTPOSTERIOR}(\theta; \{z_n\}_{n=1}^N)} > u$  then ▷ Evaluate MH ratio conditioned on auxiliary variables.
9:        $\theta_i \leftarrow \theta'$  ▷ Accept proposal.
10:    else
11:       $\theta_i \leftarrow \theta_{i-1}$  ▷ Reject proposal and keep current state.
12:    function  $\text{JOINTPOSTERIOR}(\theta; \{z_n\}_{n=1}^N)$  ▷ Modified posterior that conditions on auxiliary variables.
13:       $P \leftarrow p(\theta) \times \prod_{n=1}^N B_n(\theta)$  ▷ Evaluate prior and bounds. Collapse of bound product not shown.
14:      for each  $n$  for which  $z_n = 1$  do ▷ Loop over bright data only.
15:         $P \leftarrow P \times (L_n(\theta)/B_n(\theta) - 1)$  ▷ Include bound-corrected factor.
16:      return  $P$ 

```

---

not  $z_n$ . Explicit resampling is a simple, low-overhead algorithm that is easy to vectorize for speed. The variant shown in Algorithm 1 is the simplest: data points are chosen at random, with replacement. Another variant is to deterministically choose a subset from which to Gibbs sample at each iteration. Such an approach may be appropriate for data sets which are too large to fit into memory, since we would no longer need random access to all data points.

Explicitly sampling a subset of the  $z_n$  is wasteful if  $M \ll N$ , since most updates to  $z_n$  will leave it unchanged. We can do better by drawing each update for  $z_n$  from a pair of tunable Bernoulli proposal distributions  $q(z'_n = 1 | z_n = 0) = q_{d \rightarrow b}$  and  $q(z'_n = 0 | z_n = 1) = q_{b \rightarrow d}$ , and performing a Metropolis–Hastings accept/reject step with the true auxiliary probability  $p(z_n | x_n, \theta)$ . It is only necessary to evaluate  $p(z_n | x_n, \theta)$  – and therefore the likelihood function – for the subset of data points which are proposed to change state. That is, if a sample from the proposal distribution sends  $z_n = 0$  to  $z_n = 0$  then it doesn’t matter whether we accept or reject. If we use samples from a geometric distribution to choose the data points, it is not even necessary to explicitly sample all of the  $N$  proposals.

## 4 Experiments

For FlyMC to be a useful algorithm it must be able to produce effectively independent samples from posterior distributions more quickly than regular MCMC. We certainly expect it to iterate more quickly than regular MCMC since it evaluates fewer likelihoods per iteration. But we might also expect it to mix more slowly, since it has extra auxiliary variables. To see whether this trade-off works out in FlyMC’s favor we need to know how much faster it iterates and how much slower it mixes.

We conducted three experiments, each with a different data set, model, and parameter-update algorithm, to give an impression of how well FlyMC can be expected to perform. In

each experiment we compared FlyMC, with two choices of bound selection, to regular full-posterior MCMC. We looked at the average number of likelihoods queried at each iteration and the number of effective samples generated per iteration, accounting for autocorrelation. The results are summarized in Figure 3 and Table 1. The broad conclusion is that FlyMC offers a speedup of at least one order of magnitude compared with regular MCMC if the bounds are tuned according to a MAP-estimate of  $\theta$ . In the following subsections we describe the experiments in detail.

### 4.1 Logistic regression

We applied FlyMC to the logistic regression task described in Welling and Teh [2011] using the Jaakkola–Jordan bounds described earlier. The task is to classify MNIST 7s and 9s, using the first 50 principal components (and one bias) as features. We used a Gaussian prior over the weights and chose the scale of that prior by evaluating performance on a held-out test set. To sample over  $\theta$ , we used symmetric Metropolis–Hasting proposals, with step size chosen to yield an acceptance rate of 0.234 [Roberts *et al.*, 1997], optimized for each algorithm separately. We sampled the  $z_n$  using the implicit Metropolis–Hastings sampling algorithm.

We compared three different algorithms: regular MCMC, untuned FlyMC, and MAP-tuned FlyMC. For untuned FlyMC, we chose  $\xi = 1.5$  for all data points. To compute the bounds for the MAP-tuned algorithm, we performed stochastic gradient descent optimization to find a set of weights close to the MAP value and gave each data point its own  $\xi$  to make the bounds tight at the MAP parameters:  $L_n(\theta_{\text{MAP}}) = B_n(\theta_{\text{MAP}})$  for all  $n$ . For untuned FlyMC, and MAP-tuned FlyMC we used  $q_{d \rightarrow b} = 0.1$  and  $q_{b \rightarrow d} = 0.01$  respectively, chosen to be similar to the typical fraction of bright data points in each case.

The results are shown in Figure 3a and summarized in Table 1. On a per-iteration basis, the FlyMC algorithms mix and burn-in more slowly than regular MCMC by around a

		Algorithm	Average Likelihood queries per iteration	Effective Samples per 1000 iterations	Speedup relative to regular MCMC
Data set:	MNIST	Regular MCMC	12,214	3.7	(1)
Model:	Logistic regression	Untuned FlyMC	6,252	1.3	0.7
Updates:	Metropolis-Hastings	MAP-tuned FlyMC	207	1.4	22
Data set:	3-Class CIFAR-10	Regular MCMC	18,000	8.0	(1)
Model:	Softmax classification	Untuned FlyMC	8,058	4.2	1.2
Updates:	Langevin	MAP-tuned FlyMC	654	3.3	11
Data set:	OPV	Regular MCMC	18,182,764	1.3	(1)
Model:	Robust regression	Untuned FlyMC	2,753,428	1.1	5.7
Updates:	Slice sampling	MAP-tuned FlyMC	575,528	1.2	29

Table 1: Results from empirical evaluations. Three experiments are shown: logistic regression applied to MNIST digit classification, softmax classification for three categories of CIFAR-10, and robust regression for properties of organic photovoltaic molecules, sampled with random-walk Metropolis–Hastings, Langevin-adjusted Metropolis, and slice sampling, respectively. For each of these, the vanilla MCMC operator was compared with both untuned FlyMC and FlyMC where the bound was determined from a MAP estimate of the posterior parameters. We use likelihood evaluations as an implementation-independent measure of computational cost and report the number of such evaluations per iteration, as well as the resulting sample efficiency (computed via R-CODA [Plummer *et al.*, 2006]), and relative speedup.

factor of two, as illustrated by the autocorrelation plots. Even on a per-likelihood basis, the naïve FlyMC algorithm, with a fixed  $\xi$ , performs worse than regular MCMC, by a factor of 0.7, despite needing fewer likelihood evaluations per iteration. The MAP-tuned algorithm was much more impressive: after burn-in, it queried only 207 of the 12,2214 likelihoods per iteration on average, giving a speedup of more than 20, even taking into account the slower per-iteration mixing time. We initialized all chains with draws from the prior. Notice that the MAP-tuned algorithm performs poorly during burn-in, since the bounds are less tight during this time, whereas the reverse is true for the untuned algorithm.

## 4.2 Softmax classification

Logistic regression can be generalized to multi-class classification problems by softmax classification. The softmax likelihood of a data point belonging to class  $k$  of  $K$  classes is

$$L_n(\theta) = \frac{\exp(\theta_k^\top x_n)}{\sum_{k'=1}^K \exp(\theta_{k'}^\top x_n)},$$

where  $\theta$  is now a  $K \times D$  matrix. The Jaakkola-Jordan bound does not apply to this softmax likelihood, but we can use a related bound, due to Böhning [1992], whose log matches the value and gradient of the log of the softmax likelihood at some particular  $\theta$ , but has a tighter curvature. Murphy [2012] has the result in full in the chapter on variational inference.

We applied softmax classification to a three-class version of CIFAR-10 (airplane, automobile and bird) using 256 binary features discovered by Krizhevsky [2009] using a deep autoencoder. Once again, we used a Gaussian prior on the weights, chosen to maximize out-of-sample performance.

This time we used the Metropolis-adjusted Langevin algorithm (MALA, Roberts and Tweedie [1996]) for our parameter updates. We chose the step sizes to yield acceptance rates close to the optimal 0.57 [Roberts and Rosenthal, 1998]. Other parameters were tuned as in the logistic regression experiment.

The softmax experiment gave qualitatively similar results to the logistic regression experiment, as seen in Figure 3b and Table 1. Again, the MAP-tuned FlyMC algorithm dramatically outperformed both the lackluster untuned FlyMC and regular MCMC, offering an 11-fold speedup over the latter.

## 4.3 Robust sparse linear regression

Linear regression with Gaussian likelihoods yields a closed-form expression for the posterior. Non-Gaussian likelihoods, however, like heavy-tailed distributions used in so-called “robust regression” do not. Our final experiment was to perform inference over robust regression weights for a very large dataset of molecular features and computed electronic properties. The data set, described by Hachmann *et al.* [2011, 2014] consists of 1.8 million molecules, with 57 cheminformatic features each [Olivares-Amaya *et al.*, 2011; Amador-Bedolla *et al.*, 2013]. The task was to predict the HOMO-LUMO energy gap, which is useful for predicting photovoltaic efficiency.

We used a student-t distribution with  $\nu = 4$  for the likelihood function and we computed a Gaussian lower bound to this by matching the value and gradient of the t distribution probability density function value at some  $\xi$  ( $\xi = 0$  for the untuned case,  $\xi = \theta_{MAP}^\top x$  for the MAP-tuned case). We used a sparsity-inducing Laplace prior on the weights. As before,

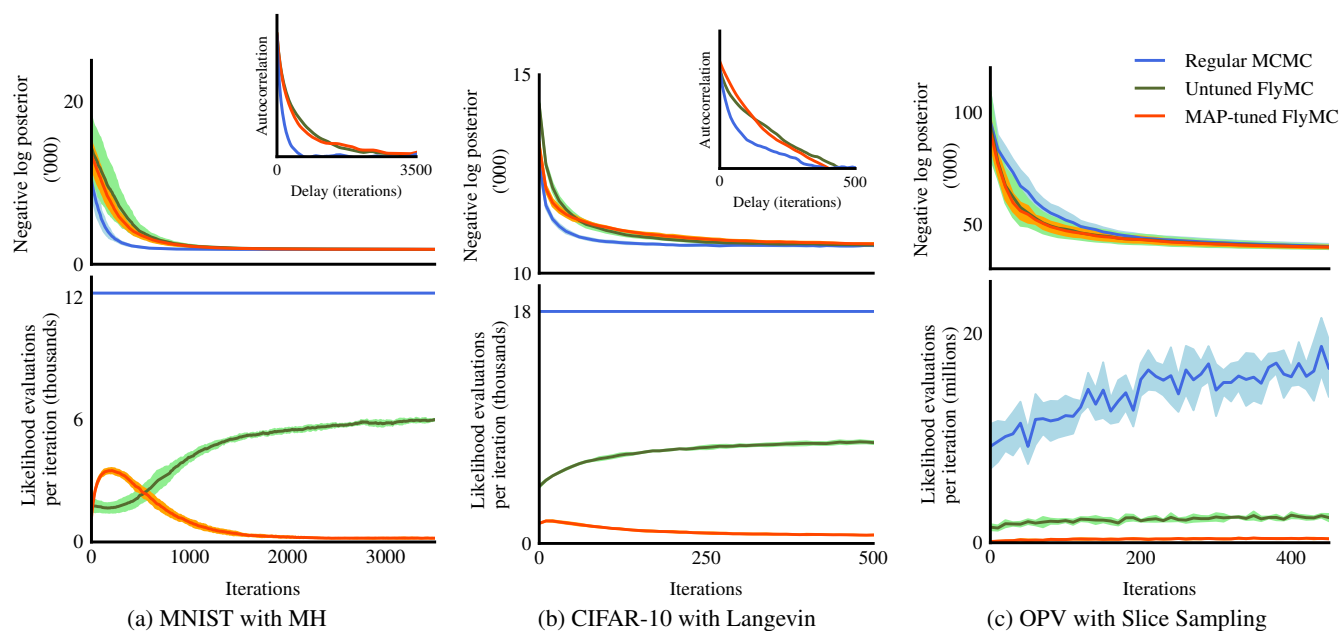


Figure 3: Tuned and untuned Firefly Monte Carlo compared to regular MCMC with three different operators, data sets, and models: (a) the digits 7 and 9 from the MNIST data are classified using logistic regression, with a random-walk Metropolis-Hastings operator; (b) softmax classification on three classes (airplane, automobile, and bird) from the CIFAR-10 image dataset, using Langevin-adjusted Metropolis; (c) robust regression on the HOMO-LUMO gap (as computed by density functional theory calculations) for a large set of organic photovoltaic molecules, using slice sampling. In each subfigure, the top shows the trace of the log posterior density to illustrate convergence, and the bottom shows the average number of likelihoods computed per iteration. One standard deviation is shown around the mean value, as computed from five runs of each. The blue lines are computed using the full-data posterior, and the green and orange lines show the untuned and tuned Firefly MC traces, respectively.

we chose the scales of the prior and the likelihood to optimize out-of-sample performance.

We performed parameter updates using slice sampling [Neal, 2003]. Note that slice sampling results in a variable number of likelihood evaluations per iteration, even for the regular MCMC algorithm. Again, we found that MAP-tuned FlyMC substantially outperformed regular MCMC, as shown in Figure 3c and Table 1.

## 5 Discussion

In this paper, we have presented Firefly Monte Carlo, a Markov chain Monte Carlo algorithm using subsets (mini-batches) of data. FlyMC is exact in the sense that it has the true full-data posterior as its target distribution. This is achieved by introducing binary latent variables whose states represent whether a given datum is bright (used to compute the posterior) or dark (not used in posterior updates). By carefully choosing the conditional distributions of these latent variables, the true posterior is left intact under marginalization. The primary requirement for this to be efficient is that the likelihoods term must have lower bounds that collapse in an efficient way.

There are several points that warrant additional discussion and future work. First, we recognize that useful lower bounds can be difficult to obtain for many problems. It would be helpful to produce such bounds automatically for a wider class of

problems. As variational inference procedures are most often framed in terms of lower bounds on the marginal likelihood, we expect that Firefly Monte Carlo will benefit from developments in so-called “black box” variational methods [Wang and Blei, 2013; Ranganath *et al.*, 2014]. Second, we believe we have only scratched the surface of what is possible with efficient data structures and latent-variable update schemes. For example, the MH proposals we consider here for  $z_n$  have a fixed global  $q_{d \rightarrow b}$ , but clearly such a proposal should vary for each datum. Third, it is often the case that larger state spaces lead to slower MCMC mixing. We have shown empirically that the slower mixing can be more than offset by the faster per-transition computational time. In future work we hope to show that fast-mixing Markov chains on the parameter space will continue to mix fast in the Firefly auxiliary variable representation.

Firefly Monte Carlo is closely related to recent ideas in using pseudo-marginal MCMC [Andrieu and Roberts, 2009] for sampling from challenging target distributions. If we sampled each of the variables  $\{z_n\}$  as a Bernoulli random variable with success probability 0.5, then the joint posterior we have been using becomes an unbiased estimator of the original posterior over  $\theta$ , up to normalization. Running pseudo-marginal MCMC using this unbiased estimator would be a special case of FlyMC: namely FlyMC with  $z$  and  $\theta$  updated simultaneously with Metropolis-Hastings updates.

## Acknowledgements

We thank Andrew Miller, Jasper Snoek, Michael Gelbart, Brenton Partridge, Oren Rippel and Elaine Angelino for helpful discussions. We also thank Alan Aspuru-Guzik, Johannes Hachmann, and Roberto Olivares-Amaya for the use of the OPV data set and introduction to the cheminformatic feature set. Partial funding was provided by Analog Devices (Lyric Labs), NSF IIS-1421780, and Samsung Advanced Institute of Technology.

## References

- Carlos Amador-Bedolla, Roberto R. Olivares-Amaya, Johannes Hachmann, and Alan Aspuru-Guzik. Organic photovoltaics. In K. Rajan, editor, *Informatics for Materials Science and Engineering*, pages 423–440. Elsevier, Amsterdam, 2013.
- Christophe Andrieu and Gareth O Roberts. The pseudo-marginal approach for efficient Monte Carlo computations. *The Annals of Statistics*, pages 697–725, 2009.
- Elaine Angelino, Eddie Kohler, Amos Waterland, Margo Seltzer, and Ryan P. Adams. Accelerating mcmc via parallel predictive prefetching. In *Proceedings of the 30th Conference on Uncertainty in Artificial Intelligence*, 2014.
- Rémi Bardenet, Arnaud Doucet, and Chris Holmes. Towards scaling up Markov chain Monte Carlo : an adaptive subsampling approach. In *Proceedings of the 31st International Conference on Machine Learning*, 2014.
- Dankmar Böhning. Multinomial logistic regression algorithm. *Annals of the Institute of Statistical Mathematics*, 44(1):197–200, 1992.
- Johannes Hachmann, Roberto Olivares-Amaya, Sule Atahan-Evrenk, Carlos Amador-Bedolla, Roel S. Sanchez-Carrera, Aryeh Gold-Parker, Leslie Vogt, Anna M. Brockway, and Alan Aspuru-Guzik. The harvard clean energy project: Large-scale computational screening and design of organic photovoltaics on the world community grid. *The Journal of Physical Chemistry Letters*, 2(17):2241–2251, 2011.
- Johannes Hachmann, Roberto Olivares-Amaya, Adrian Jinich, Anthony L. Appleton, Martin A. Blood-Forsythe, Laszlo R. Seress, Carolina Roman-Salgado, Kai Trepte, Sule Atahan-Evrenk, Suleyman Er, Supriya Shrestha, Rajib Mondal, Anatoliy Sokolov, Zhenan Bao, and Alan Aspuru-Guzik. Lead candidates for high-performance organic photovoltaics from high-throughput quantum chemistry - the harvard clean energy project. *Energy Environ. Sci.*, 7:698–704, 2014.
- Matthew D. Hoffman, David M. Blei, and Francis Bach. Online learning for latent Dirichlet allocation. In *Advances in Neural Information Processing Systems 23*, 2010.
- Matthew D. Hoffman, David M. Blei, Chong Wang, and John Paisley. Stochastic variational inference. *Journal of Machine Learning Research*, 14(1):1303–1347, 2013.
- Tommi S. Jaakkola and Michael I. Jordan. A variational approach to Bayesian logistic regression models and their extensions. In *Workshop on Artificial Intelligence and Statistics*, 1997.
- Anoop Korattikara, Yutian Chen, and Max Welling. Austerity in MCMC Land: Cutting the Metropolis-Hastings Budget. In *Proceedings of the 31st International Conference on Machine Learning*, 2014.
- Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, Department of Computer Science, University of Toronto, 2009.
- Kevin P. Murphy. *Machine learning: a probabilistic perspective*. MIT Press, Cambridge, MA, 2012.
- Radford M. Neal. Slice sampling. *The Annals of Statistics*, 31(3):705–767, 06 2003.
- Roberto Olivares-Amaya, Carlos Amador-Bedolla, Johannes Hachmann, Sule Atahan-Evrenk, Roel S. Sanchez-Carrera, Leslie Vogt, and Alan Aspuru-Guzik. Accelerated computational discovery of high-performance materials for organic photovoltaics by means of cheminformatics. *Energy Environ. Sci.*, 4:4849–4861, 2011.
- Martyn Plummer, Nicky Best, Kate Cowles, and Karen Vines. CODA: Convergence diagnosis and output analysis for MCMC. *R News*, 6(1):7–11, 2006.
- Rajesh Ranganath, Sean Gerrish, and David M. Blei. Black box variational inference. In *Proceedings of the 17th International Conference on Artificial Intelligence and Statistics*, 2014.
- Gareth O Roberts and Jeffrey S Rosenthal. Optimal scaling of discrete approximations to Langevin diffusions. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 60(1):255–268, 1998.
- Gareth O Roberts and Richard L Tweedie. Exponential convergence of Langevin distributions and their discrete approximations. *Bernoulli*, pages 341–363, 1996.
- Gareth O. Roberts, Andrew Gelman, and Walter R. Gilks. Weak convergence and optimal scaling of random walk Metropolis algorithms. *Annals of Applied Probability*, 7:110–120, 1997.
- Chong Wang and David M Blei. Variational inference in non-conjugate models. *The Journal of Machine Learning Research*, 14(1):1005–1031, 2013.
- Max Welling and Yee Whye Teh. Bayesian learning via stochastic gradient Langevin dynamics. In *Proceedings of the 28th International Conference on Machine Learning*, 2011.