

FIREFLY MONTE CARLO: EXACT MCMC WITH SUBSETS OF DATA

MCMC

We have a unnormalized distribution from which we want to sample

- Metropolis-Hastings algorithm for MCMC

With probability $p = \min \left[1, \exp \left(\frac{U(\mathbf{x}_n) - U(\mathbf{x}'_n)}{T} \right) \right]$ system **accepts new state** (jumps to the new state): $\mathbf{x}_{n+1} = \mathbf{x}'_n$ and with probability $1 - p$ **new state is rejected**: $\mathbf{x}_{n+1} = \mathbf{x}_n$.

- Hamiltonian Monte Carlo

$$E(\mathbf{x}, \mathbf{v}) = U(\mathbf{x}) + K(\mathbf{v}), \quad K(\mathbf{v}) = \sum_i \frac{m v_i^2}{2}$$

Thus, velocities \mathbf{v} and positions \mathbf{x} have **independent** canonical distributions:

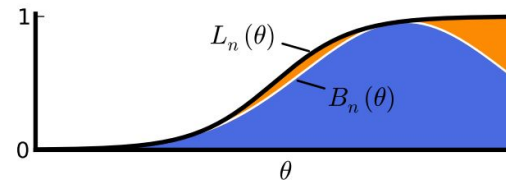
$$p(\mathbf{x}, \mathbf{v}) \propto \exp \left(\frac{-E(\mathbf{x}, \mathbf{v})}{T} \right) = \exp \left(\frac{-U(\mathbf{x})}{T} \right) \exp \left(\frac{-K(\mathbf{v})}{T} \right) \propto p(\mathbf{x}) p(\mathbf{v}).$$

MCMC problem

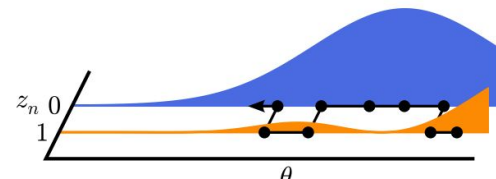
MCMC cannot be practically applied to large data sets because of the prohibitive cost of evaluating every likelihood term at every iteration.

Solution

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



$$\begin{aligned} & p(x_n | \theta) p(z_n | x_n, \theta) \\ &= L_n(\theta) \left[\frac{L_n(\theta) - B_n(\theta)}{L_n(\theta)} \right]^{z_n} \left[\frac{B_n(\theta)}{L_n(\theta)} \right]^{1-z_n} \\ &= \begin{cases} L_n(\theta) - B_n(\theta) & \text{if } z_n = 1 \\ B_n(\theta) & \text{if } z_n = 0 \end{cases} . \end{aligned}$$



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:

Algorithm 1 Firefly Monte Carlo

Note: Using simple random-walk MH for clarity.

```
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:   end for
7:    $\theta' \leftarrow \theta_{i-1} + \eta$  where  $\eta \sim \text{Normal}(0, \epsilon^2 \mathbb{I}_D)$  ▷ Make a random walk proposal with step size  $\epsilon$ .
8:    $u \sim \text{Uniform}(0, 1)$  ▷ Draw the MH threshold.
9:   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.
10:     $\theta_i \leftarrow \theta'$  ▷ Accept proposal.
11:   else
12:     $\theta_i \leftarrow \theta_{i-1}$  ▷ Reject proposal and keep current state.
13:   end if
14: end for
15:
16: function  $\text{JOINTPOSTERIOR}(\theta; \{z_n\}_{n=1}^N)$  ▷ Modified posterior that conditions on auxiliary variables.
17:    $P \leftarrow p(\theta) \times \prod_{n=1}^N B_n(\theta)$  ▷ Evaluate prior and bounds. Collapse of bound product not shown.
18:   for each  $n$  for which  $z_n = 1$  do ▷ Loop over bright data only.
19:      $P \leftarrow P \times (L_n(\theta)/B_n(\theta) - 1)$  ▷ Include bound-corrected factor.
20:   end for
21:   return  $P$ 
22: end function
```

In the case of tight bound, explicit sampling doesn't effective. So the author's propose the following solution

Algorithm 2 Implicit z_n sampling

1: for $n \leftarrow 1 \dots N$ do	▷ Loop over all the auxiliary variables.
2: if $z_n = 1$ then	▷ If currently bright, propose going dark.
3: $u \sim \text{Uniform}(0, 1)$	▷ Sample the MH threshold.
4: if $\frac{q_{d \rightarrow b}}{\tilde{L}_n(\theta)} > u$ then	▷ Compute MH ratio with $\tilde{L}_n(\theta)$ cached from θ update.
5: $z_n \leftarrow 0$	▷ Flip from bright to dark.
6: end if	
7: else	▷ Already dark, consider proposing to go bright.
8: if $v < q_{d \rightarrow b}$ where $v \sim \text{Uniform}(0, 1)$ then	▷ Flip a biased coin with probability $q_{d \rightarrow b}$.
9: $u \sim \text{Uniform}(0, 1)$	▷ Sample the MH threshold.
10: if $\frac{\tilde{L}_n(\theta)}{q_{d \rightarrow b}} < u$ then	▷ Compute MH ratio.
11: $z_n \leftarrow 1$	▷ Flip from dark to bright.
12: end if	
13: end if	
14: end if	
15: end for	

Results

		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.