

Piecewise deterministic Markov processes for scalable Monte Carlo on restricted domains

Bierkens, Joris; Bouchard-Côté, Alexandre; Doucet, Arnaud ; Duncan, Andrew B.; Fearnhead, Paul ; Lienart, Thibaut ; Roberts, Gareth; Vollmer, Sebastian J.

DOI

[10.1016/j.spl.2018.02.021](https://doi.org/10.1016/j.spl.2018.02.021)

Publication date

2018

Document Version

Accepted author manuscript

Published in

Statistics and Probability Letters

Citation (APA)

Bierkens, J., Bouchard-Côté, A., Doucet, A., Duncan, A. B., Fearnhead, P., Lienart, T., Roberts, G., & Vollmer, S. J. (2018). Piecewise deterministic Markov processes for scalable Monte Carlo on restricted domains. *Statistics and Probability Letters*, 136, 148-154. <https://doi.org/10.1016/j.spl.2018.02.021>

Important note

To cite this publication, please use the final published version (if applicable).
Please check the document version above.

Copyright

Other than for strictly personal use, it is not permitted to download, forward or distribute the text or part of it, without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license such as Creative Commons.

Takedown policy

Please contact us and provide details if you believe this document breaches copyrights.
We will remove access to the work immediately and investigate your claim.

Piecewise Deterministic Markov Processes for Scalable Monte Carlo on Restricted Domains

Joris Bierkens^{a,*}, Alexandre Bouchard-Côté^b, Arnaud Doucet^c, Andrew B. Duncan^d,
Paul Fearnhead^e, Thibaut Lienart^c, Gareth Roberts^f, Sebastian J. Vollmer^f

^a*Delft Institute of Applied Mathematics, TU Delft, Netherlands*

^b*Department of Statistics, University of British Columbia, Canada*

^c*Department of Statistics, University of Oxford, UK*

^d*School of Mathematical and Physical Sciences, University of Sussex, UK*

^e*Department of Mathematics and Statistics, Lancaster University, UK*

^f*Department of Statistics, University of Warwick, UK*

Abstract

Piecewise Deterministic Monte Carlo algorithms enable simulation from a posterior distribution, whilst only needing to access a sub-sample of data at each iteration. We show how they can be implemented in settings where the parameters live on a restricted domain.

Keywords: MCMC, Bayesian statistics, piecewise deterministic Markov processes, logistic regression

1. Introduction

Markov chain Monte Carlo (MCMC) methods have been central to the wide-spread use of Bayesian methods. However their applicability to some modern applications has been limited due to their high computational cost, particularly in big-data, high-dimensional settings. This has led to interest in new MCMC methods, particularly non-reversible methods which can mix better than standard reversible MCMC [9, 23], and variants of MCMC that require accessing only small subsets of the data at each iteration [25].

One of the main technical challenges associated with likelihood-based inference for big data is the fact that likelihood calculation is computationally expensive (typically $O(N)$ for data sets of size N). MCMC methods built from piecewise deterministic Markov processes (PDMPs) offer considerable promise for reducing this $O(N)$ burden, due to their ability to use sub-sampling

*Corresponding author

Email address: joris.bierkens@tudelft.nl (Joris Bierkens)

techniques, whilst still being guaranteed to target the true posterior distribution [5, 7, 11, 12, 18]. Furthermore, factor graph decompositions of the target distribution can be leveraged to perform sparse updates of the variables [7, 17, 21].

PDMPs explore the state space according to constant velocity dynamics, but where the velocity changes at random event times. The rate of these event times, and the change in velocity at each event, are chosen so that the position of the resulting process has the posterior distribution as its invariant distribution. We will refer to this family of sampling methods as Piecewise Deterministic Monte Carlo methods (PDMC).

Existing PDMC algorithms can only be used to sample from posteriors where the parameters can take any value in \mathbb{R}^d . In this paper (Section 2) we show how to extend PDMC methodology to deal with constraints on the parameters. Such models are ubiquitous in machine learning and statistics. For example, many popular models used for binary, ordinal and polychotomous response data are multivariate real-valued latent variable models where the response is given by a deterministic function of the latent variables [1, 10, 22]. Under the posterior distribution, the domain of the latent variables is then constrained based on the values of the responses. Additional examples arise in regression where prior knowledge restricts the signs of marginal effects of explanatory variables such as in econometrics [13], image processing and spectral analysis [3], [14] and non-negative matrix factorization [15]. A few methods for dealing with restricted domains are available but these either target an approximation of the correct distribution [20] or are limited in scope [19].

2. Piecewise Deterministic Monte Carlo on Restricted Domains

Here we present the general PDMC algorithm in a restricted domain. Specific implementations of PDMC algorithms can be derived as continuous-time limits of familiar discrete-time MCMC algorithms [6, 21], and these derivations convey much of the intuition behind why the algorithms have the correct stationary distribution. Our presentation of these methods is different, and more general. We first define a simple class of PDMPs and show how these can be simulated. We then give simple recipes for how to choose the dynamics of the PDMP so that it will have the correct stationary distribution.

Our objective is to compute expectations with respect to a probability distribution π on

$\mathcal{O} \subseteq \mathbb{R}^d$ which is assumed to have a smooth density, also denoted $\pi(x)$, with respect to the Lebesgue measure on \mathcal{O} . With this objective in mind, we will construct a continuous-time Markov process $Z_t = (X_t, V_t)_{t \geq 0}$ taking values in the domain $E = \mathcal{O} \times \mathcal{V}$, where \mathcal{O} and \mathcal{V} are subsets of \mathbb{R}^d , such that \mathcal{O} is open, pathwise connected and with Lipschitz boundary $\partial\mathcal{O}$. In particular, if $\mathcal{O} = \mathbb{R}^d$ then $\partial\mathcal{O} = \emptyset$. The dynamics of Z_t are easy to describe if one views X_t as position and V_t as velocity. The position process X_t moves deterministically, with constant velocity V_t between a discrete set of *switching times* which are simulated according to N inhomogeneous Poisson processes, with respective intensity functions $\lambda_i(X_t, V_t)$, $i = 1, \dots, N$, depending on the current state of the system. At each switching time the position stays the same, but the velocity is updated according to a specified transition kernel. More specifically, suppose the next switching event occurs from the i^{th} Poisson process, then the velocity immediately after the switch is sampled randomly from the probability distribution $Q_i(x, v, \cdot)$ given the current position x and velocity v . The switching times are random, and designed in conjunction with the kernels $(Q_i)_{i=1}^N$ so that the invariant distribution of the process coincides with the target distribution π .

To ensure that X_t remains confined within \mathcal{O} the velocity of the process is updated whenever X_t hits $\partial\mathcal{O}$ so that the process moves back into \mathcal{O} . We shall refer to such updates as *reflections* even though they need not be specular reflections.

The resulting stochastic process is a Piecewise Deterministic Markov Process (PDMP, [8]). For it to be useful as the basis of a Piecewise Deterministic Monte Carlo (PDMC) algorithm we need to (i) be able to easily simulate this process; and (ii) have simple recipes for choosing the intensities, $(\lambda_i)_{i=1}^N$, and transition kernels, $(Q_i)_{i=1}^N$, such that the resulting process has $\pi(x)$ as its marginal stationary distribution. We will tackle each of these problems in turn.

2.1. Simulation

The key challenge in simulating our PDMP is simulating the event times. The intensity of events is a function of the state of the process. But as the dynamics between event times are deterministic, we can easily represent the intensity for the next event as a deterministic function of time. Suppose that the PDMP is driven by a single inhomogeneous Poisson process with intensity function

$$\tilde{\lambda}(u; X_t, V_t) = \lambda(X_t + uV_t, V_t), \quad u \geq 0.$$

We can simulate the first event time directly if we have an explicit expression for the inverse function of the monotonically increasing function

$$u \mapsto \int_0^u \tilde{\lambda}(s; X_t, V_t) ds. \quad (1)$$

In this case the time until the next event is obtained by (i) simulating a realization, y say, of an exponential random variable with rate 1; and (ii) setting the time until the next event as the value τ that solves $\int_0^\tau \tilde{\lambda}(s; X_t, V_t) ds = y$.

Inverting (1) is often not practical. In such cases simulation can be carried out via *thinning* [16]. This requires finding a tractable upper bound on the rate, $\bar{\lambda}(u) \geq \tilde{\lambda}(u; X_t, V_t)$ for all $u > 0$. Such an upper bound will typically take the form of a piecewise linear function or a step function. Note that the upper bound $\bar{\lambda}$ is only required to be valid along the trajectory $u \mapsto (X_t + uV_t, V_t)$ in $\mathcal{O} \times \mathcal{V}$. Therefore the upper bound may depend on the starting point (X_t, V_t) of the line segment we are currently simulating. We then propose potential events by simulating events from an inhomogenous Poisson process with rate $\bar{\lambda}(u)$, and accept an event at time u with probability $\tilde{\lambda}(u; X_t, V_t)/\bar{\lambda}(u)$. The time of the first accepted event will be the time until the next event in our PDMP.

To handle boundary reflections, at every given time t , we also keep track of the next reflection event in the absence of a switching event, i.e. we compute

$$\tau_b = \inf \{u > 0 : X_t + uV_t \notin \mathcal{O}\}.$$

If the boundary $\partial\mathcal{O}$ can be represented as a finite set of M hyper-planes in \mathbb{R}^d , then the cost of computing τ_b is $O(Md)$. When generating the switching event times and positions for Z_t we determine whether a boundary reflection will occur before the next potential switching event. If so, then we induce a switching event at time $t + \tau_b$ where $X_{t+\tau_b} \in \partial\mathcal{O}$ and sample a new velocity from the transition kernel Q_b , i.e. $V_{t+\tau_b} \sim Q_b(X_{t+\tau_b}, V_t, \cdot)$.

Although theoretically we may choose a new velocity pointing outwards and have an immediate second jump, we will for algorithmic purposes assume that the probability measure $Q_b(x, u, \cdot)$ for $(x, u) \in \partial\mathcal{O} \times \mathcal{V}$ is concentrated on those directions v for which $(v \cdot n(x)) \leq 0$, where $n(x)$ is the outward normal at $x \in \partial\mathcal{O}$.

For a PDMP driven by N inhomogeneous Poisson processes with intensities $(\lambda_i)_{i=1}^N$ the

previous steps lead to the following algorithm for simulating the next event of our PDMP. This algorithm can be iterated to simulate the PDMP for a chosen number of events or a pre-specified time-interval.

- (0) **Initialize:** Set t to the current time and (X_t, V_t) to the current position and velocity.
- (1) **Determine bound:** For each $i \in 1, \dots, N$, find a convenient function $\bar{\lambda}_i$ satisfying $\bar{\lambda}_i(u) \geq \tilde{\lambda}_i(u; X_t, V_t)$ for all $u \geq 0$, depending on the initial point (X_t, V_t) from which we are departing.
- (2) **Propose event:** For $i = 1, \dots, N$ simulate the first event times τ'_i of a Poisson process with rate function $\bar{\lambda}_i$. Compute the next boundary reflection time τ_b .
- (3) Let $i_{\min} = \arg \min_{j=1, \dots, N} \tau'_j$ and $\tau' = \tau'_{i_{\min}}$.
- (4) **Accept/Reject event:**
 - (4.1) If $\tau_b < \tau'$ then set $\tau = \tau_b$; set $X_{t+\tau} = X_t + \tau V_t$; sample a new velocity $V_{t+\tau} \sim Q_b(X_{t+\tau}, V_t, \cdot)$.
 - (4.2) Otherwise with probability

$$\frac{\tilde{\lambda}_{i_{\min}}(\tau'; X_t, V_t)}{\bar{\lambda}_{i_{\min}}(\tau')}$$
 accept the event at time $\tau = \tau'$.
 - (4.2.1) **Upon acceptance:** set $X_{t+\tau} = X_t + \tau V_t$; sample a new velocity $V_{t+\tau} \sim Q_{i_{\min}}(X_{t+\tau}, V_t, \cdot)$.
 - (4.2.2) **Upon rejection:** set $X_{t+\tau} = X_t + \tau V_t$ and set $V_{t+\tau} = V_t$.
- (5) **Update:** Record the time $t + \tau'$ and state $(X_{t+\tau'}, V_{t+\tau'})$.

2.2. Output of PDMC algorithms

The output of these algorithms will be a sequence of event times $t_1, t_2, t_3, \dots, t_K$ and associated states $(X_1, V_1), (X_2, V_2), \dots, (X_K, V_K)$. To obtain the value of the process at times $t \in [t_k, t_{k+1})$, we can linearly interpolate the continuous path of the process between event times, i.e. $X_t = X_{t_k} + V_k(t - t_k)$. Time integrals $\int_0^t f(X_s) ds$ of a function f of the process X_t can often be computed analytically from the output of the above algorithm. If not they can be

approximated by numerically integrating the one dimensional integral along the piecewise linear trajectory of the PDMP. Alternatively we can sample the PDMP at a set of evenly spaced time points along the trajectory and use this collection as an approximate sample from our target distribution.

Under the assumption that the resulting PDMP is ergodic (for sufficient conditions see e.g. [5, 7]) and that the marginal density on \mathcal{O} of the stationary distribution of (X_t, V_t) is equal to π , we have the following version of the law of large numbers for the PDMP $(X_t, V_t)_{t \geq 0}$: For all $f \in L^2(\pi)$ we have that, with probability one,

$$\int_{\mathbb{R}^d} f(x) \pi(x) dx = \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T f(X_s) ds.$$

It is this formula which allows us to use PDMPs for Monte Carlo purposes.

2.3. Choosing the intensity and transition kernels

Assume, as most existing PDMC methods do [5, 7, 21], that the target density, $\pi(x) : \mathcal{O} \rightarrow (0, \infty)$ is differentiable. Under this condition we can provide criteria on the switching intensities (λ_i) and transition kernels Q_i and Q_b which must hold for a given probability distribution to be a stationary distribution of Z_t . We shall consider stationary distributions for which x and v are independent, i.e. distributions of the form $\pi(x)dx \otimes \rho(dv)$ on E . Furthermore we assume that $\pi(x) \propto \exp(-U(x))$ where U is continuously differentiable.

We impose the condition that

$$\int_{v \in \mathcal{V}} \sum_{i=1}^N \lambda_i(x, v) Q_i(x, v, du) \rho(dv) = \int_{v \in \mathcal{V}} \sum_{i=1}^N \lambda_i(x, v) Q_i(x, u, dv) \rho(du), \quad x \in \mathcal{O}. \quad (2)$$

A sufficient condition for (2) is that each Q_i is reversible with respect to ρ , i.e. for every $i = 1, \dots, N$ and $x \in \mathcal{O}$, we have that $Q_i(x, v, du) \rho(dv) = Q_i(x, u, dv) \rho(du)$.

Moreover, we shall require the following condition which relates the probability flow with the switching intensities λ_i :

$$\sum_{i=1}^N \int_{\mathcal{V}} \lambda_i(x, v) Q_i(x, u, dv) - \sum_{i=1}^N \lambda_i(x, u) = -u \cdot \nabla U(x), \quad (x, u) \in \mathcal{O} \times \mathcal{V}. \quad (3)$$

Finally, the boundary transition kernel should satisfy

$$Q_b(x, u, dv) \rho(du) = Q_b(x, v, du) \rho(dv), \quad x \in \partial \mathcal{O}, \quad (4)$$

and

$$\int_{\mathcal{V}} (n(x) \cdot u) Q_b(x, v, du) = -v \cdot n(x), \quad (x, v) \in \partial\mathcal{O} \times \mathcal{V}, \quad (5)$$

where for $x \in \partial\mathcal{O}$, we denote by $n(x)$ the outward unit normal of $\partial\mathcal{O}$.

Proposition 1. *Consider the process Z_t on $\mathcal{O} \times \mathcal{V}$ where \mathcal{O} is an open, pathwise connected subset of \mathbb{R}^d with Lipschitz boundary $\partial\mathcal{O}$. Suppose that conditions (2),(3), (4) and (5) are satisfied. Then $\pi(x) dx \otimes \rho(dv)$ is an invariant distribution for the process Z_t .*

The proof of this result relies on verifying that $\mathbb{E}_{\pi \otimes \rho}[\mathcal{L}f(X, V)] = 0$ where \mathcal{L} denotes the generator of our PDMP and is deferred to the supplementary material, Section 1.

In practice we only have to satisfy (4) and (5) on the exit region $\Gamma \subset \mathcal{O} \times \mathcal{V}$. For example if $\mathcal{O} = (a, b) \subset \mathbb{R}^1$ and $\mathcal{V} = \{-1, +1\}$, then $\Gamma = \{b, +1\} \cup \{a, -1\}$. The specification of Q_b on $(\partial\mathcal{O} \times \mathcal{V}) \setminus \Gamma$ is irrelevant as these points are never reached by Z_t . On this irrelevant set, we may choose Q_b as desired to satisfy (5).

2.4. Example: The Bouncy Particle Sampler

Current PDMC algorithms differ in terms of how the Q_i and λ_i are chosen such that the above equation holds for some simple distribution for the velocity. Here we discuss how the Bouncy Particle Sampler (BPS), introduced in [21] and explored in [7], is an example of the framework introduced here. In the supplementary material, Section 1.1, the Zig-Zag sampler is described as a second example. In the following example δ_x denotes the Dirac-measure centered in x .

The Bouncy Particle Sampler is obtained setting $N = 1$ and $\rho = \mathcal{N}(0, I)$ on \mathbb{R}^d or $\rho = \mathcal{U}(S^{d-1})$, i.e. the uniform distribution on the unit sphere. The single switching rate is chosen to be $\lambda_{\text{BPS}}(x, v) = \max(v \cdot \nabla U(x), 0)$, with corresponding switching kernel Q which reflects v with respect to the orthogonal complement of ∇U with probability 1:

$$Q(x, v, dv') = \delta_{(I - 2P_{\nabla U})v}(dv'),$$

where $P_y : z \mapsto \frac{z \cdot y}{\|y\|^2} y$ denotes orthogonal projection along the one dimensional subspace spanned by y .

As noted in [7] this algorithm suffers from reducibility issues. These can be overcome by refreshing the velocity by drawing a new velocity independently from $\rho(dv)$. In the simplest case

the refreshment times come from an independent Poisson process with constant rate λ_{ref} . This also fits in the framework above by choosing $\tilde{\lambda} = \lambda_{\text{BPS}} + \lambda_{\text{ref}}$ and

$$Q(x, u, dv) = \frac{\lambda_{\text{BPS}}}{\lambda_{\text{BPS}} + \lambda_{\text{ref}}} \delta_{(I - 2P_{\nabla U})u}(dv) + \frac{\lambda_{\text{ref}}}{\lambda_{\text{BPS}} + \lambda_{\text{ref}}} \rho(dv).$$

As boundary transition kernel it is natural to choose

$$Q_b(s, v, du) = \delta_{(I - 2P_{n(s)})v}(du),$$

for $s \in \partial\mathcal{O}$, so that the process X_t reflects specularly at the boundary (i.e. angle of incidence equals angle of reflection of process with respect to the boundary normal). It is straightforward to check that condition (2) holds at the boundary and that (5) is satisfied.

As a generalization of the BPS, one can consider a *preconditioned* version, which is obtained by introducing a constant positive definite symmetric matrix M to rescale the velocity process. The choice of M plays a very similar role to the mass matrix in HMC, and careful tuning can give rise to dramatic increases in performance [12, 18].

3. Subsampling

When using PDMC to sample from a posterior, we can use sub-samples of data at each iteration of the algorithm, as described in [5, 7], which reduces the computational complexity of the algorithm from $O(N)$ to $O(1)$, where N is the size of the data, without affecting the theoretical validity of the algorithm. In the following we will assume that we can write the posterior as $\pi(x) \propto \prod_{i=1}^N f(y_i; x)$, for some function f . For example this would be the likelihood for a single IID data point times the $1/N$ th power of the prior.

The idea of using sub-sampling, within say the Bouncy Particle Sampler (BPS), is that at each iteration of our PDMC algorithm we can replace $\nabla U(x)$ by an unbiased estimator in step (3). We need to use the same estimate both when calculating the actual event rate in the accept/reject step and, if we accept, when simulating the new velocity. The only further alteration we need to the algorithm is to choose an upper bound $\bar{\lambda}$ that holds for all realizations of $\widehat{\nabla U}$. A more comprehensive explanation of this argument can be found in [5, 11] in the context of the Zig-Zag sampler, and in [7, 12] for the bouncy particle sampler.

We first present a way for estimating ∇U unbiasedly using control variates [2, 5]. For any

$x, \hat{x} \in \mathcal{O}$ we note that $\nabla U(x) = \nabla U(\hat{x}) + [\nabla U(x) - \nabla U(\hat{x})]$. We can then introduce the estimator $\widehat{\nabla U}(x)$ of $\nabla U(x)$ by

$$\widehat{\nabla U}(x) = \nabla U(\hat{x}) + N [\nabla \log f(y_I; x) - \nabla \log f(y_I; \hat{x})], \quad (6)$$

where I is drawn uniformly from $\{1, \dots, N\}$.

It is straightforward to show that the resulting BPS algorithm uses an event rate that is $\mathbb{E} \left[\max \left(0, (\widehat{\nabla U}(x) \cdot v) \right) \right]$, and that this rate and the resulting transition probability Q at events satisfies Proposition 1. Hence this algorithm still targets $\pi(x)$, but only requires access to one data point at each accept-reject decision.

Note that this gain in computational efficiency does not come for free, as it follows from Jensen's inequality that the overall rate of events will be higher. This makes mixing of the PDMC process slower. It is also immediate that the bound, $\bar{\lambda}$, we will have to use will be higher. However [5] show that if our estimator of $\widehat{\nabla U}(x)$ has sufficiently small variance, then we can still gain substantially in terms of efficiency. In particular they give an example where the CPU cost effective sample size does not grow with N – by comparison all standard MCMC algorithms would have a cost that is at least linear in N .

To obtain such a low-variance estimator requires a good choice of \hat{x} , so that with high probability x will be closer to \hat{x} . This involves a preprocessing step to find a value \hat{x} close to the posterior mode, a preprocessing step to then calculate $\nabla U(\hat{x})$ is also needed.

We now illustrate how to find an upper bound on the event rate. Following [5], if we assume L is a uniform (in space and i) upper bound on the largest eigenvalue of the Hessian of U^i , and if $\|v\| = 1$:

$$\begin{aligned} & \max \left(0, (\nabla U(\hat{x}) + N(\nabla U^i(X_t) - \nabla U^i(\hat{x}))) \cdot v \right) \\ & \leq \max \left(0, \nabla U(\hat{x}) \cdot v \right) + N \left\| \nabla U^i(x) - \nabla U^i(\hat{x}) \nabla U^i(x) - \nabla^i U(X_t) \right\| \\ & \leq \max \left(0, \nabla U(\hat{x}) \cdot v \right) + NL \|x - \hat{x}\| + NLt \end{aligned} \quad (7)$$

Thus the upper bound on the intensity is of the form $\bar{\lambda}(\tau) = a + b \cdot \tau$ with $a, b \geq 0$. In this case the first arrival time can be simulated as follows

$$\tau' = -a/b + \sqrt{\left(\frac{a}{b}\right)^2 + 2 \cdot \frac{R}{b}} \text{ with } R \sim \text{Exp}(1). \quad (8)$$

An alternative and complementary approach to improve the efficiency of this subsampling procedure is to use an estimator of the gradient (3) where I is drawn according to a distribution dependent on the observations [7, 12].

4. Software and Numerical Experiments

A open-source Julia package `PDMP.jl` has been developed to provide efficient implementations of various recently developed piecewise deterministic Monte Carlo methods for sampling in (possibly restricted) continuous spaces. A variety of algorithms are implemented including the Zig-Zag sampler and the Bouncy Particle Sampler with full and local refreshment along with control variate based sub-sampling for these methods. The package has been specifically designed with extensibility in mind, permitting rapid implementation of new PDMP based methods. The library along with code and documentation is available at github.com/alan-turing-institute/PDSampler.jl.

We use Bayesian binary logistic regression as a testbed for our newly proposed methodology and perform a simulation study. The data $y_i \in \{-1, 1\}$ is modelled by

$$p(y_i | \xi_i, x) = f(y_i x^T \xi_i) \tag{9}$$

where $\xi \in \mathbb{R}^{p \times n}$ are fixed covariates and $f(z) = \frac{1}{1 + \exp(-z)} \in [0, 1]$. We will assume that we wish to fit this model under some monotonicity constraints – so that the probability of $y = 1$ is known to either increase or decrease with certain covariates. This is modeled through the constraint $x_i > 0$ and $x_i < 0$ respectively. An example where such restrictions occur naturally is in logistic regression for questionnaires, see [24]. In following we consider the case $x_j \geq 0$ for $j = 1, \dots, p$ along with the additional linear constraint $\sum_j x_j \leq K$ where $K = 10$.

For simplicity we use a flat prior over the space of parameters values consistent with our constraints. By Bayes' rule the posterior π satisfies

$$\pi(x) \propto \prod_{i=1}^N f(y_i x^T \xi_i) \text{ for } x \in \mathcal{O},$$

where \mathcal{O} is the space of parameter values consistent with our constraints. We implement the BPS with subsampling. As explained in the introduction, subsampling is a key benefit of using piecewise deterministic sampling methods; see Section 3. We use reflection at the boundary

i.e. $Q_b(s, v, du) = \delta_{(I-2P_{n(s)})v}(du)$ for $s \in \partial\mathcal{O}$. We can bound the switching intensity by a linear function of time, even when we use the subsampling estimator for the switching rate. See the supplementary material, Section 2, for details on the application of subsampling in this example. We use $n = 10,000$ and $p = 20$ and generate artificial data based on ξ and x^* whose components are a realization i.i.d. of uniformly distributed random variables satisfying the imposed constraints.

We compare the performance of BPS to standard MALA and HMC schemes, in terms of effective sample size (ESS) per epoch of data evaluation. For each scheme we obtain the distribution of ESS based on 10 independent realisations of each chain. In Figure 1(a) we plot for each scheme, the distribution of ESS per epoch with respect to the function $f_1(x) = \frac{1}{p}(x_1 + \dots + x_p)$. Similarly, In Figure 1(b) we plot the ESS per epoch for each chain with respect to the function $f_2(x) = \log \pi(x)$. The performance of MALA and HMC appears commensurate and the BPS demonstrates a clear advantage over both in terms of ESS per epoch.

The HMC and MALA schemes were tuned by minimising the ESS with respect to the step-size, calculated from exploratory runs. For HMC we use 5 leap-frog steps. We find that we must tune both HMC and MALA to have a small step size due to proposals being rejected at the boundary. The ESS is estimated based on asymptotic variance using the batch means method; see [4, 5] for details.

For specific types of constraints more efficient implementations of HMC and MALA are possible, either by introducing an appropriate transformation of the restricted state space, or by reflecting the posterior distribution along the constraint boundaries. Moreover, we note that there exists a version of HMC which can sample from truncated Gaussian distributions [19]. However, to our knowledge there is no efficient HMC or MALA scheme able to handle generally restricted domains.

The Bouncy Particle Sampler for this model was implemented using `PDSampler.jl` while the corresponding HMC and MALA samplers implemented with `Klara.jl`. The code for this numerical experiment along with results are carefully presented in github.com/tlienart/ConstrainedPDMP/.

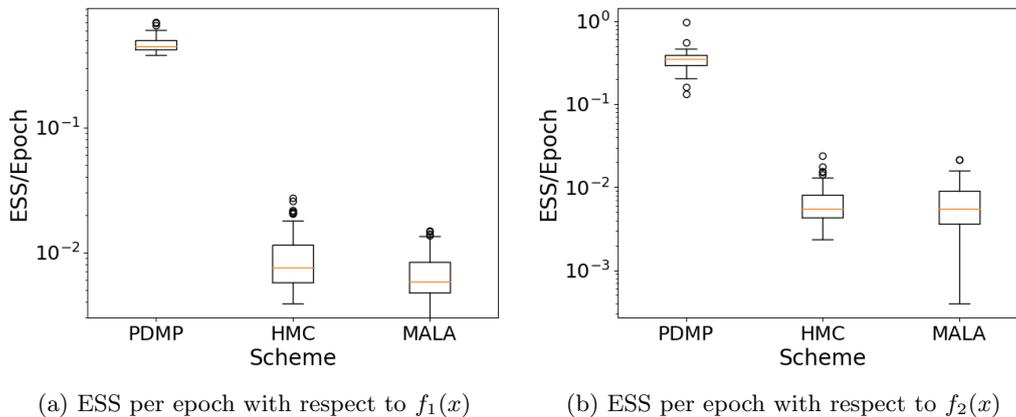


Figure 1: Average ESS per epochs of data evaluation for MALA, HMC and PDMP (BPS) applied to logistic regression with $p = 20$ and $n = 10000$ and parameter x constrained to be nonnegative and satisfy $\sum_j x_j \leq 10$. The graphic is based on 10 independent runs for each HMC, MALA and BPS for each choice of number of epochs.

5. Discussion

This work provides a framework for describing a general class of PDMC methods which are ergodic with respect to a given target probability distribution. Open questions remain on how the choice of intensity function, velocity transition kernel as well as other parameters of the system influence the overall performance of the scheme. The problem of understanding the true computational cost of such PDMC schemes is more subtle than for classical discrete time MCMC schemes: often one needs to find a balance between fast mixing of the continuous time Markov process and having a switching rate that is relatively cheap to simulate. For example, when using subsampling the mixing of the Markov process is slower than without subsampling, but the computational cost per simulated switch is significantly smaller. Further investigation is required to understand this delicate balance.

Acknowledgements

All authors thank the Alan Turing Institute and Lloyds registry foundation for support. S.J.V. gratefully acknowledges funding through EPSRC EP/N000188/1. J.B., P.F. and G.R. gratefully acknowledge EPSRC ilike grant EP/K014463/1. A.B.D. acknowledges grant EP/L020564/1 and A.D. acknowledges grant EP/K000276/1.

References

- [1] J. H. Albert and S. Chib. Bayesian analysis of binary and polychotomous response data. *Journal of the American Statistical Association*, 88(422):669–679, 1993.
- [2] R. Bardenet, A. Doucet, and C. Holmes. On Markov chain Monte Carlo methods for tall data. *Journal of Machine Learning Research*, 18(18):1–43, 2017.
- [3] S. Bellavia, M. Macconi, and B. Morini. An interior point Newton-like method for non-negative least-squares problems with degenerate solution. *Numerical Linear Algebra with Applications*, 13(10):825–846, 2006.
- [4] J. Bierkens and A. Duncan. Limit theorems for the Zig-Zag process. *Advances in Applied Probability*, 49(3), jul 2017.
- [5] J. Bierkens, P. Fearnhead, and G. Roberts. The Zig-Zag process and super-efficient sampling for Bayesian analysis of big data. *arXiv preprint arXiv:1607.03188*, 2016.
- [6] J. Bierkens and G. Roberts. A piecewise deterministic scaling limit of lifted Metropolis-Hastings in the Curie- Weiss model. *Annals of Applied Probability*, 27(2):846–882, 2017.
- [7] A. Bouchard-Côté, S. J. Vollmer, and A. Doucet. The bouncy particle sampler: A non-reversible rejection-free Markov chain Monte Carlo method. *arXiv preprint arXiv:1510.02451*, 2015.
- [8] M. H. A. Davis. Piecewise-deterministic Markov processes: A general class of non-diffusion stochastic models. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 353–388, 1984.
- [9] P. Diaconis, S. Holmes, and R. Neal. Analysis of a nonreversible Markov chain sampler. *Annals of Applied Probability*, 10(3):726–752, 2000.
- [10] L. Fahrmeir and G. Tutz. *Multivariate Statistical Modelling based on Generalized Linear Models*. Springer, New York, 2001.
- [11] P. Fearnhead, J. Bierkens, M. Pollock, and G. O. Roberts. Piecewise deterministic Markov processes for continuous-time Monte Carlo. *arXiv preprint arXiv:1611.07873*, 2016.

- [12] N. Galbraith. On Event-Chain Monte Carlo Methods. Master’s thesis, Department of Statistics, Oxford University, 2016.
- [13] J. Geweke. Exact inference in the inequality constrained normal linear regression model. *Journal of Applied Econometrics*, 1(2):127–141, 1986.
- [14] Y. Guo and M. Berman. A comparison between subset selection and l1 regularisation with an application in spectroscopy. *Chemometrics and Intelligent Laboratory Systems*, 118:127–138, 2012.
- [15] H. Kim and H. Park. Sparse non-negative matrix factorizations via alternating non-negativity-constrained least squares for microarray data analysis. *Bioinformatics*, 23(12):1495–1502, 2007.
- [16] P. A. Lewis and G. S. Shedler. Simulation of nonhomogeneous Poisson processes by thinning. *Naval Research Logistics Quarterly*, 26(3):403–413, 1979.
- [17] Y. Nishikawa, M. Michel, W. Krauth, and K. Hukushima. Event-chain algorithm for the Heisenberg model. *Physical Review E*, 92(6):63306, 2015.
- [18] A. Pakman, D. Gilboa, D. Carlson, and L. Paninski. Stochastic bouncy particle sampler. *arXiv preprint arXiv:1609.00770*, 2016.
- [19] A. Pakman and L. Paninski. Exact Hamiltonian Monte Carlo for truncated multivariate Gaussians. *Journal of Computational and Graphical Statistics*, 23(2):518–542, 2014.
- [20] S. Patterson and Y. W. Teh. Stochastic gradient Riemannian Langevin dynamics on the probability simplex. In *Advances in Neural Information Processing Systems*, pages 3102–3110, 2013.
- [21] E. A. J. F. Peters and G. De With. Rejection-free Monte Carlo sampling for general potentials. *Physical Review E*, 85(2):1–5, 2012.
- [22] L. E. Train. *Discrete Choice Methods with Simulation*. Cambridge university press, 2009.
- [23] K. S. Turitsyn, M. Chertkov, and M. Vucelja. Irreversible Monte Carlo algorithms for efficient sampling. *Physica D: Nonlinear Phenomena*, 240(4-5):410–414, 2011.

- [24] G. Tutz and J. Gertheiss. Rating scales as predictors—the old question of scale level and some answers. *Psychometrika*, 79(3):357–376, 2014.
- [25] M. Welling and Y. W. Teh. Bayesian learning via stochastic gradient Langevin dynamics. In *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, pages 681–688, 2011.

Supplement to *Piecewise Deterministic Markov Processes for Scalable Monte Carlo on Restricted Domains*

Joris Bierkens, Alexandre Bouchard-Côté, Arnaud Doucet, Andrew B. Duncan,
Paul Fearnhead, Thibaut Lienart, Gareth Roberts, Sebastian J. Vollmer

January 7, 2018

In the supplementary material we provide a theoretical background for the framework for restricted domain PDMC (Section 1, which includes the Zig-Zag sampler as further example). The detailed application of subsampling to the logistic regression example may be found in Section 2.

1 Stationary distribution for PDMPs on restricted domains

From [2, Section 5] the process Z_t will have infinitesimal generator given by the closure of the operator

$$\mathcal{L}f(x, v) = v \cdot \nabla_x f(x, v) + \sum_{i=1}^N \lambda_i(x, v) \int_{\mathcal{V}} (f(x, u) - f(x, v)) Q_i(x, v, du), \quad (x, v) \in E, \quad (1)$$

where $\mathcal{D}(\mathcal{L})$ is the set of functions which are continuously differentiable with respect to x on \mathcal{O} , which is decaying to infinity as $\|x\| \rightarrow \infty$ and such that

$$f(x, v) = \int_{\mathcal{V}} f(x, u) Q_b(x, v, du), \quad (2)$$

for all $(x, v) \in \partial\mathcal{O} \times \mathcal{V}$. Based on this identification of the infinitesimal generator we can now provide a formal proof that the conditions of Proposition 1 of the paper are sufficient to ensure invariance of $\pi \otimes \rho$.

Sketch Proof of Proposition 1. Without loss of generality we take $\pi(x) = \exp(-U(x))$, i.e. the proportionality factor in $\pi(x) \propto \exp(-U(x))$ is assumed to be 1. We shall only provide a formal proof of this result, by demonstrating that

$$\int_{\mathcal{O} \times \mathcal{V}} \mathcal{L}f(x, v) \pi(x) dx \rho(dv) = 0, \quad \text{for all } f \in \mathcal{D}(\mathcal{L}),$$

so that \mathcal{L} is infinitesimally invariant. A rigorous proof would require establishing that $\mathcal{D}(\mathcal{L})$ as defined above is a core for the extended generator. This is a technical result which we defer for future work.

For $f \in \mathcal{D}(\mathcal{L})$,

$$\begin{aligned}
& \int_{\mathcal{O}} \int_{\mathcal{V}} \int_{\mathcal{V}} \sum_{i=1}^N \lambda_i(x, v) [f(x, u) - f(x, v)] Q_i(x, v, du) \pi(x) dx \rho(dv) \\
&= \int_{\mathcal{O}} \int_{\mathcal{V}} \sum_{i=1}^N f(x, u) \left[\int_{\mathcal{V}} \lambda_i(x, v) Q_i(x, u, dv) - \lambda_i(x, u) \right] \rho(du) \pi(x) dx \\
&= - \int_{\mathbb{R}^d} \int_{\mathcal{V}} f(x, u) u \cdot \nabla U(x) \rho(du) e^{-U(x)} dx \\
&= \int_{\mathcal{O}} \int_{\mathcal{V}} f(x, u) u \cdot \nabla e^{-U(x)} \rho(du) dx \\
&= - \int_{\mathcal{O}} \int_{\mathcal{V}} u \cdot \nabla_x f(x, u) e^{-U(x)} \rho(du) dx + \int_{\partial\mathcal{O}} \int_{\mathcal{V}} f(\sigma, u) (u \cdot n(\sigma)) e^{-U(\sigma)} \rho(du) d\sigma,
\end{aligned}$$

where the boundary term arises from integration by parts with respect to x . Considering the boundary integral, by applying (4) (in the paper) which is assumed to hold on $\partial\mathcal{O}$ and (2) (above) we obtain

$$\begin{aligned}
& \int_{\partial\mathcal{O}} \int_{\mathcal{V}} f(\sigma, u) (u \cdot n(\sigma)) e^{-U(\sigma)} \rho(du) d\sigma \\
&= \int_{\partial\mathcal{O}} \int_{\mathcal{V}} \int_{\mathcal{V}} f(\sigma, v) Q_b(\sigma, u, dv) (u \cdot n(\sigma)) e^{-U(\sigma)} \rho(du) d\sigma \\
&= \int_{\partial\mathcal{O}} \int_{\mathcal{V}} \int_{\mathcal{V}} f(\sigma, v) Q_b(\sigma, v, du) (u \cdot n(\sigma)) e^{-U(\sigma)} \rho(dv) d\sigma \\
&= - \int_{\partial\mathcal{O}} \int_{\mathcal{V}} f(\sigma, v) (v \cdot n(\sigma)) e^{-U(\sigma)} \rho(dv) d\sigma,
\end{aligned}$$

so that the boundary term evaluates to zero.

It follows that

$$\int_{\mathcal{O}} \int_{\mathcal{V}} \mathcal{L}f(x, v) \pi(x) dx \rho(dv) = \int_{\mathcal{O}} \int_{\mathcal{V}} (u \cdot \nabla_x f(x, u) - u \cdot \nabla_x f(x, u)) \pi(x) dx \rho(dv) = 0,$$

so that $\pi(x) dx \otimes \rho(dv)$ is infinitesimally invariant with respect to Z_t . \square

Another possible behaviour at the boundary is to generate the new reflected direction independently of the angle of incidence. This will also preserve the invariant distribution provided that ρ is isotropic.

Proposition 1. *Consider the process Z_t as in the previous proposition, such that conditions (2) and (3) (of the paper) hold and the distribution ρ has mean zero. Then $\pi(x) dx \otimes \rho(dv)$ will be an invariant distribution for the process Z_t if $Q_b(x, v, du)$ is independent of v for all $x \in \partial\mathcal{O}$.*

Sketch Proof of Proposition 1. Let $f \in \mathcal{D}(\mathcal{L})$, so that f satisfies (2). By the assumptions on Q_b in Proposition 1 (of the paper), this implies that $f(x, v) = f(x)$ for all $x \in \partial\mathcal{O}$. Following the proof of Proposition 1 above, the boundary integral term becomes

$$\int_{\partial\mathcal{O}} \int_{\mathcal{V}} f(s, u) (u \cdot n(s)) e^{-U(s)} \rho(du) ds = \int_{\partial\mathcal{O}} f(s) e^{-U(s)} n(s) ds \cdot \int_{\mathcal{V}} u \rho(du),$$

which is zero if ρ has mean zero, as required. \square

1.1 The Zig-Zag sampler

The Zig-Zag sampler [1] can be recovered by choosing $N = d$ and picking as velocity space $\mathcal{V} = \{-1, +1\}^d$ equipped with discrete uniform distribution ρ , defining switching rates $\lambda_i(x, v) = \max(v_i \partial_{x_i} U(x), 0)$. The corresponding switching kernels over new directions are given by

$$Q_i(x, v, dv') = \delta_{F_i v}(dv'),$$

where $F_i : \mathcal{V} \rightarrow \mathcal{V}$ denotes the operation of flipping the i -th component, i.e. $(F_i v)(i) = -v(i)$, and $(F_i v)(j) = v(j)$ for $j \neq i$.

2 Derivation of dominating intensity for logistic regression example

A valid choice of L can be derived as follows: Notice that $(\log f(z))' = f(-z)$ and $f'(z) = f(z)(1 - f(z))$ so that we obtain

$$\begin{aligned} \frac{\partial}{\partial x} \log f(y_i | x) &= f(-y_i x^\top \xi_i) y_i \xi_i \\ \frac{\partial}{\partial^2 x} \log f(y_i | x) &= -f(-y_i x^\top \xi_i) (1 - f(-y_i x^\top \xi_i)) \xi_i \xi_i^\top \end{aligned}$$

Using $p(1 - p) \leq \frac{1}{4}$ for $p \in [0, 1]$

$$\sup_{\|w\| \leq 1} \left| w^\top \frac{\partial}{\partial^2 x} \log f(y_i | x) w \right| \leq \frac{1}{4} \|\xi_i\|^2$$

So Equation (7) (of the paper) holds with

$$L := \frac{1}{4} \max_{i=1, \dots, n} \|\xi_i\|$$

as defined above.

This is a linear upper bound on the intensity which can be used to sample according to (8) (of the paper) and then used for thinning as introduced in Section 1 of the paper.

References

- [1] J. Bierkens, P. Fearnhead, and G. Roberts. The Zig-Zag process and super-efficient sampling for Bayesian analysis of big data. *arXiv preprint arXiv:1607.03188*, 2016.
- [2] M. H. A. Davis. Piecewise-deterministic Markov processes: A general class of non-diffusion stochastic models. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 353–388, 1984.