

Finite-Time Analysis of Projected Langevin Monte Carlo

Authored by:

Sebastien Bubeck
Ronen Eldan
Joseph Lehec

Abstract

We analyze the projected Langevin Monte Carlo (LMC) algorithm, a close cousin of projected Stochastic Gradient Descent (SGD). We show that LMC allows to sample in polynomial time from a posterior distribution restricted to a convex body and with concave log-likelihood. This gives the first Markov chain to sample from a log-concave distribution with a first-order oracle, as the existing chains with provable guarantees (lattice walk, ball walk and hit-and-run) require a zeroth-order oracle. Our proof uses elementary concepts from stochastic calculus which could be useful more generally to understand SGD and its variants.

1 Paper Body

A fundamental primitive in Bayesian learning is the ability to sample from the posterior distribution. Similarly to the situation in optimization, convexity is a key property to obtain algorithms with provable guarantees for this task. Indeed several Markov Chain Monte Carlo methods have been analyzed for the case where the posterior distribution is supported on a convex set, and the negative log-likelihood is convex. This is usually referred to as the problem of sampling from a log-concave distribution. In this paper we propose and analyze a new Markov chain for this problem which could have several advantages over existing chains for machine learning applications. We describe formally our contribution in Section 1.1. Then in Section 1.2 we explain how this contribution relates to various line of work in different fields such as theoretical computer science, statistics, stochastic approximation, and machine learning.

Main result

Let $K \subset \mathbb{R}^n$ be a convex set such that $0 \in K$, K contains a Euclidean ball of radius $r > 0$ and is contained in a Euclidean ball of radius R . Denote PK the Euclidean projection on K (i.e., $PK(x) = \arg\min_{y \in K} \|x - y\|$ where $\|\cdot\|$ denotes the Euclidean norm in \mathbb{R}^n), and $k \in K$ the gauge of K defined

by $\|x\|_K = \inf\{t \geq 0; x \in tK\}$, $x \in \mathbb{R}^n$. Let $f: K \rightarrow \mathbb{R}$ be a L -Lipschitz and μ -smooth convex function, that is f is differentiable and satisfies $\langle x, y \rangle \in K, -\langle f(x) - f(y), x - y \rangle \leq L \|x - y\|^2$, and $\|f(x) - f(y)\| \leq L \|x - y\|$. We are interested in the problem of sampling from the probability measure μ on \mathbb{R}^n whose density with respect to the Lebesgue measure is given by: $Z d\mu = \exp(-f(x)) 1_{\{x \in K\}}$, where $Z = \int_K \exp(-f(y)) dy$.

We denote $m = \mathbb{E}[\|X\|]$, and $M = \mathbb{E}[\|k\|_K]$, where μ is uniform on the sphere $S^{n-1} = \{x \in \mathbb{R}^n : \|x\| = 1\}$. In this paper we study the following Markov chain, which depends on a parameter $\gamma \in (0, 1]$, and where μ_1, μ_2, \dots is an i.i.d. sequence of standard Gaussian random variables in \mathbb{R}^n , and $X_0 = 0$,

$X_{k+1} = \gamma X_k + (1-\gamma) \mu_k$. (1) We call the chain (1) projected Langevin Monte Carlo (LMC). Recall that the total variation distance between two measures μ, ν is defined as $TV(\mu, \nu) = \sup_A |\mu(A) - \nu(A)|$ where the supremum is over all measurable sets A . With a slight abuse of notation we sometimes write $TV(X, \mu)$ where X is a random variable distributed according to μ (respectively ν). e^{-c} means that there exists $c \in \mathbb{R}, C \in (0, 1]$ such that $\mu \leq C e^{-c}$. The notation $\nu_n = O(u \nu_n)$ (respectively $\nu_n \leq C u \nu_n$) (respectively ν_n). Our main result shows that for an appropriately chosen step-size and number of iterations, one has convergence in total variation distance of the iterates (X_k) to the target distribution μ . Theorem 1 Let $\gamma \in (0, 1]$. One has $TV(X_N, \mu) \leq \gamma$ provided that $\gamma \leq 1/(N+1)$.

6

$\gamma \leq 1/(N+1)$ Note that by viewing μ, L, r as numerical constants, using $M \leq 1/r$, and assuming $R \leq n$ and $m \leq n^{3/4}$, the bound reads $\gamma \leq 1/(N+1)$. Observe also that if f is constant, that is μ is the uniform measure on K , then $L = 0$, $m \leq n$, and μ one can show that $M = O(1/\sqrt{n})$, which yields the bound: $\gamma \leq 1/(N+1)$.

Context and related works

There is a long line of works in theoretical computer science proving results similar to Theorem 1, starting with the breakthrough result of Dyer et al. [1991] who showed that the lattice walk mixes in e^{23} steps. The current record for the mixing time is obtained by Lovász and Vempala [2007], $O(n \log n)$ for the hit-and-run walk. These chains (as well as other popular chains who show a bound of $O(n \log n)$ such as the ball walk or the Dikin walk, see e.g. Kannan and Narayanan [2012] and references therein) all require a zeroth-order oracle for the potential f , that is given x one can calculate the value $f(x)$. On the other hand our proposed chain (1) works with a first-order oracle, that is given x one can calculate the value of $\nabla f(x)$. The difference between zeroth-order oracle and firstorder oracle has been extensively studied in the optimization literature (e.g., Nemirovski and Yudin [1983]), but it has been largely ignored in the literature on polynomial-time sampling algorithms. We also note that hit-and-run and LMC are the only chains which are rapidly mixing from any starting point (see Lovász and Vempala [2006]), though they have this property for seemingly very different reasons. When initialized in a corner of the convex body, hit-and-run might take a long time to take a step, but once it moves it escapes very

far (while a chain such as the ball walk would only do a small step). On the other hand LMC keeps moving at every step, even when initialized in a corner, thanks for the projection part of (1). Our main motivation to study the chain (1) stems from its connection with the ubiquitous stochastic gradient descent (SGD) algorithm. In general this algorithm takes the form $x_{k+1} = \text{PK}(x_k - \eta \nabla f(x_k) + \xi_k)$ where ξ_1, ξ_2, \dots is a centered i.i.d. sequence. Standard results in approximation theory, such as Robbins and Monro [1951], show that if the variance of the noise $\text{Var}(\xi_1)$ is of smaller order than the step-size η then the iterates (x_k) converge to the minimum of f on K (for a step-size decreasing sufficiently fast as a function of the number of iterations). For the specific noise

sequence that we study in (1), the variance is exactly equal to the step-size, which is why the chain deviates from its standard and well-understood behavior. We also note that other regimes where SGD does not converge to the minimum of f have been studied in the optimization literature, such as the constant step-size case investigated in Pflug [1986], Bach and Moulines [2013]. The chain (1) is also closely related to a line of works in Bayesian statistics on Langevin Monte Carlo algorithms, starting essentially with Tweedie and Roberts [1996]. The focus there is on the unconstrained case, that is $K = \mathbb{R}^n$. In this simpler situation, a variant of Theorem 1 was proven in the recent paper Dalalyan [2014]. The latter result is the starting point of our work. A straightforward way to extend the analysis of Dalalyan to the constrained case is to run the unconstrained chain with an additional potential that diverges quickly as the distance from x to K increases. However it seems much more natural to study directly the chain (1). Unfortunately the techniques used in Dalalyan [2014] cannot deal with the singularities in the diffusion process which are introduced by the projection. As we explain in Section 1.3 our main contribution is to develop the appropriate machinery to study (1). In the machine learning literature it was recently observed that Langevin Monte Carlo algorithms are particularly well-suited for large-scale applications because of the close connection to SGD. For instance Welling and Teh [2011] suggest to use mini-batch to compute approximate gradients instead of exact gradients in (1), and they call the resulting algorithm SGLD (Stochastic Gradient Langevin Dynamics). It is conceivable that the techniques developed in this paper could be used to analyze SGLD and its refinements introduced in Ahn et al. [2012]. We leave this as an open problem for future work. Another interesting direction for future work is to improve the polynomial dependency on the dimension and the inverse accuracy in Theorem 1 (our main goal here was to provide the simplest polynomial-time analysis).

1.3 Contribution and paper organization

As we pointed out above, Dalalyan [2014] proves the equivalent of Theorem 1 in the unconstrained case. His elegant approach is based on viewing LMC as a discretization of the diffusion process $dX_t = -\nabla f(X_t)dt + \sqrt{2\eta}dW_t$, where (W_t) is a Brownian motion. The analysis then proceeds in two steps, by deriving first the mixing time of the diffusion process, and then showing that the discretized process is “close” to its continuous version. In Dalalyan [2014] the first step

is particularly transparent as he assumes λ -strong convexity for the potential f , which in turns directly gives a mixing time of order $1/\lambda$. The second step is also simple once one realizes that LMC (without projection) can be viewed as the diffusion process $dX_t = dW_t - \lambda \nabla f(X_t) dt$. Using Pinsker's inequality and Girsanov's formula it is then a short calculation to show that the total variation distance between X_t and X_t is small. The constrained case presents several challenges, arising from the reflection of the diffusion process on the boundary of K , and from the lack of curvature in the potential (indeed the constant potential case is particularly important for us as it corresponds to λ being the uniform distribution on K). Rather than a simple Brownian motion with drift, LMC with projection can be viewed as the discretization of reflected Brownian motion with drift, which is a process of the form $dX_t = dW_t - \lambda \nabla f(X_t) dt + \lambda L(dt)$, where $X_t \in K$, $t \geq 0$, L is a measure supported on $\{t \geq 0 : X_t \in \partial K\}$, and ν_t is an outer normal unit vector of K at X_t . The term $\lambda L(dt)$ is referred to as the Tanaka drift. Following Dalalyan [2014] the analysis is again decomposed in two steps. We study the mixing time of the continuous process via a simple coupling argument, which crucially uses the convexity of K and of the potential f . The main difficulty is in showing that the discretized process (X_t) is close to the continuous version (X_t) , as the Tanaka drift prevents us from a straightforward application of Girsanov's formula. Our approach around this issue is to first use a geometric argument to prove that the two processes are close in Wasserstein distance, and then to show that in fact for a reflected Brownian motion with drift one can deduce a total variation bound from a Wasserstein bound. In this extended abstract we focus on the special case where f is a constant function, that is λ is uniform on the convex body K . The generalization to an arbitrary smooth potential can be found in the supplementary material. The rest of the paper is organized as follows. Section 2 contains the main technical arguments. We first remind the reader of Tanaka's construction (Tanaka [1979]) of reflected Brownian motion in Section 2.1. We present our geometric argument to bound the

Wasserstein distance between (X_t) and (X_t) in Section 2.2, and we use our coupling argument to bound the mixing time of (X_t) in Section 2.3. The derivation of a total variation bound from the Wasserstein bound is discussed in Section 2.4. Finally we conclude the paper in Section 3 with some preliminary experimental comparison between LMC and hit-and-run.

2

The constant potential case

In this section we derive the main arguments to prove Theorem 1 when f is a constant function, that is $\lambda f = 0$. For a point $x \in \partial K$ we say that ν is an outer unit normal vector at x if $\langle \nu, x \rangle = 1$ and $\langle \nu, y \rangle \leq 0$, $\forall y \in K$.

For $x \in \partial K$ we say that ν is an outer unit normal at x . We define the support function h_K of K by $h_K(y) = \sup \{\langle x, y \rangle : x \in K\}$, $y \in \mathbb{R}^n$.

Note that h_K is also the gauge function of the polar body of K . 2.1

The Skorokhod problem

Let $T \in \mathbb{R}_+ \cup \{+\infty\}$ and $w : [0, T) \rightarrow \mathbb{R}^n$ be a piecewise continuous path with $w(0) \in K$. We say that $x : [0, T) \rightarrow \mathbb{R}^n$ and $\gamma : [0, T) \rightarrow \mathbb{R}^n$ solve the Skorokhod problem for w if one has $x(t) \in K$, $\forall t \in [0, T)$, $x(t) = w(t) + \gamma(t)$, $\gamma(t) \in [0, T)$, and furthermore γ is of the form $\int_0^t \gamma_s L(ds)$, $\gamma_t \in [0, T)$, 0

where γ_s is an outer unit normal at $x(s)$, and L is a measure on $[0, T]$ supported on the set $\{t \in [0, T) : x(t) \in \partial K\}$. The path x is called the reflection of w at the boundary of K , and the measure L is called the local time of x at the boundary of K . Skorokhod showed the existence of such a pair (x, γ) in dimension 1 in Skorokhod [1961], and Tanaka extended this result to convex sets in higher dimensions in Tanaka [1979]. Furthermore Tanaka also showed that the solution is unique, and if w is continuous then so is x and γ . In particular the reflected Brownian motion in K , denoted (X_t) , is defined as the reflection of the standard Brownian motion (W_t) at the boundary of K (existence follows by continuity of W_t). Observe that by Itô's formula, for any smooth function g on \mathbb{R}^n , $\int_0^t \nabla g(X_s) \cdot dW_s = g(X_t) - g(X_0) - \frac{1}{2} \int_0^t \Delta g(X_s) ds$. To get a sense of what a solution typically looks like, let us work out the case where w is piecewise constant (this will also be useful to realize that LMC can be viewed as the solution to a Skorokhod problem). For a sequence $g_1, \dots, g_N \in \mathbb{R}^n$, and for $t \geq 0$, we consider the path: $w(t) =$

$$\sum_{k=1}^N \mathbf{1}_{\{t \in [k, (k+1))\}} g_k,$$

Define $(x_k)_{k=0, \dots, N}$ inductively by $x_0 = 0$ and $x_{k+1} = PK(x_k + g_k)$. It is easy to verify that the solution to the Skorokhod problem for w is given by $x(t) = x_{\lfloor t \rfloor}$ and $\gamma(t) = \sum_{k=1}^{\lfloor t \rfloor} L(ds)$, where the measure L is defined by (denoting δ_s for a dirac at s) $L =$

$$\sum_{k=1}^N \delta_{x_k + g_k} \otimes PK(x_k + g_k) - \delta_{x_k},$$

and for $s = x_k$, $\gamma_s =$

$$x_k + g_k - PK(x_k + g_k) - x_k + g_k - PK(x_k + g_k) - 4$$

2.2

Discretization of reflected Brownian motion

Given the discussion above, it is clear that when f is a constant function, the chain (1) can be viewed as the reflection (X_t) of a discretized Brownian motion $W_t := W_{\lfloor t \rfloor}$ at the boundary of K (more precisely the value of $X_{\lfloor t \rfloor}$ coincides with the value of X_k as defined by (1)). It is rather clear that the discretized Brownian motion (W_t) is "close" to the path (W_t) , and we would like to carry this to the reflected paths (X_t) and (\tilde{X}_t) . The following lemma extracted from Tanaka [1979] allows to do exactly that. Lemma 1 Let w and \tilde{w} be piecewise continuous path and assume that (x, γ) and $(\tilde{x}, \tilde{\gamma})$ solve the Skorokhod problems for w and \tilde{w} , respectively. Then for all time t we have

$$-x(t) + x(t) - 2 \int_0^t w(s) ds - 2 \int_0^t w(s) ds + 2 \int_0^t w(s) ds + 2 \int_0^t w(s) ds$$

Applying the above lemma to the processes (W_t) and (W_t) at time $T = N$ yields (note that $WT = W_T$) $Z_T Z_T - XT + X_T - \int_0^T hW_t + W_t$, $\int_0^T iL(dt) + 2 \int_0^T hW_t + W_t$, $\int_0^T iL(dt) = 0$

We claim that the second integral is equal to 0. Indeed, since the discretized process is constant on the intervals $[k, (k+1))$ the local time L is a positive combination of Dirac point masses at $1, 2, \dots, N$. On the other hand $W_k = W_k$ for all integer k , hence the claim. Therefore $Z_T hW_t + W_t$, $\int_0^T iL(dt) = XT + X_T - \int_0^T hW_t + W_t$

Using the inequality $\langle x, y \rangle \leq \|x\| \|y\|$ we get Z

$$-XT + X_T - \int_0^T \sup_{t \in [0, T]} \|W_t - W_T\| dt \leq \int_0^T \|W_t - W_T\| dt$$

$$\int_0^T \|W_t - W_T\| dt \leq \int_0^T \|W_t - W_T\| dt$$

Taking the square root, expectation and using Cauchy-Schwarz we get $\mathbb{E} \int_0^T \|W_t - W_T\| dt \leq \sqrt{\mathbb{E} \int_0^T \|W_t - W_T\|^2 dt}$

$$\mathbb{E} \int_0^T \|W_t - W_T\| dt \leq \sqrt{\mathbb{E} \int_0^T \|W_t - W_T\|^2 dt} \leq \sqrt{\mathbb{E} \int_0^T \|W_t - W_T\|^2 dt}$$

(3)

The next two lemmas deal with each term in the right hand side of the above equation, and they will show that there exists a universal constant C such that $\frac{1}{4} \mathbb{E} \int_0^T \|W_t - W_T\|^2 dt \leq C (\log(T/\epsilon))^{3/4} T^{1/2} M^{1/2}$. (4) We discuss why the above bound implies a total variation bound in Section 2.4. Lemma 2 We have, for all $t \geq 0$

$$\int_0^t \|W_s - W_T\|^2 ds \leq \int_0^t \|W_s - W_T\|^2 ds$$

Proof By Itô's formula $d\|W_t - W_T\|^2 = 2\langle W_t - W_T, dW_t \rangle + n dt - 2\langle W_t - W_T, dW_t \rangle + n dt$. Now observe that by definition of the reflection, if t is in the support of L then $\langle W_t - W_T, dW_t \rangle = \langle W_t - W_T, dW_t \rangle + \langle W_t - W_T, dW_t \rangle$. In other words $\langle W_t - W_T, dW_t \rangle = \langle W_t - W_T, dW_t \rangle$. Therefore $\int_0^t \|W_s - W_T\|^2 ds \leq \int_0^t \|W_s - W_T\|^2 ds + \int_0^t \|W_s - W_T\|^2 ds$

0

The first term of the right-hand side is a martingale, so using that $X_0 = 0$ and taking expectation we get the result.

5

Lemma 3 There exists a universal constant C such that $\mathbb{E} \int_0^T \sup_{t \in [0, T]} \|W_t - W_T\| dt \leq C M$

P

$$\mathbb{E} \int_0^T \sup_{t \in [0, T]} \|W_t - W_T\| dt \leq C M$$

$$[0, T]$$

Proof Note that

"

#

$$\mathbb{E} \int_0^T \sup_{t \in [0, T]} \|W_t - W_T\| dt = \mathbb{E} \int_0^T \sup_{t \in [0, T]} \|W_t - W_T\| dt$$

$$\max_{0 \leq i \leq N-1} Y_i$$

where $Y_i =$

$$k W_t - W_t - k K \cdot \sup_{t \in [i, (i+1)]} W_t$$

Observe that the variables (Y_i) are identically distributed, let $p \geq 1$ and write $\frac{1}{p} \sum_{i=0}^{N-1} Y_i$

$$\frac{1}{N} \sum_{i=0}^{N-1} Y_i \leq \frac{1}{N} \sum_{i=0}^{N-1} (k W_t - W_t - k K \cdot \sup_{t \in [i, (i+1)]} W_t)$$

We claim that

(5) $k W_t - W_t - k K \cdot \sup_{t \in [i, (i+1)]} W_t \leq C p^{-1} M$ for some constant C , and for all $p \geq 2$. Taking this for granted and choosing $p = \log(N)$ in the previous inequality yields the result (recall that $N = T/\epsilon$). So it is enough to prove (5). Observe that since (W_t) is a martingale, the process $M_t = k W_t - W_t - k K \cdot \sup_{t \in [0, t]} W_t$ is a submartingale. By Doob's maximal inequality $k W_0 - k K \cdot \sup_{t \in [0, T]} W_t \leq k M_T \leq k p^{-1} M$, [0, T]

for every $p \geq 2$. Letting μ_n be the standard Gaussian measure on \mathbb{R}^n and using Khintchine's inequality we get $\int_{\mathbb{R}^n} |Z|^{1/p} d\mu_n \leq C p^{1/2} \int_{\mathbb{R}^n} |Z| d\mu_n$ (dx) $\int_{\mathbb{R}^n} |Z| d\mu_n \leq C p^{1/2} \int_{\mathbb{R}^n} |Z| d\mu_n$

\mathbb{R}^n

Lastly, integrating in polar coordinate, it is easily seen that $\int_{\mathbb{R}^n} |Z| d\mu_n \leq C n^{1/2} M$. \mathbb{R}^n

2.3

A mixing time estimate for the reflected Brownian motion

Given a probability measure μ supported on K , we let μ_t be the law of X_t when X_0 has law μ . The following lemma is the key result to estimate the mixing time of the process (X_t) . Lemma 4 Let $x, x_0 \in K$

$$\|x - x_0\| - TV(\mu_x, \mu_{x_0}) \leq C \sqrt{t}$$

The above result clearly implies that for a probability measure μ on K , $TV(\mu_t, \mu) \leq C \sqrt{t}$. Since μ (the uniform measure on K) is stationary for reflected Brownian motion, we obtain $TV(\mu_t, \mu) \leq C \sqrt{t}$. (6) In other words, starting from 0, the mixing time of (X_t) is of order t^2 .

We now turn to the proof of the above lemma. Proof The proof is based on a coupling argument. Let (W_t) be a Brownian motion starting from 0 and let (X_t) be a reflected Brownian motion starting from x :

$$dX_t = dx + dW_t + \frac{1}{2} L_t dt$$

where (L_t) and L satisfy the appropriate conditions. We construct a reflected Brownian motion (X_t^0) starting from x_0 as follows. Let $\tau = \inf\{t \geq 0; X_t = X_t^0\}$, and for $t \leq \tau$ let S_t be the orthogonal reflection with respect to the hyperplane $(X_t - X_t^0)^\perp$. Then up to time τ , the process (X_t^0) is defined by $(0 \leq t \leq \tau) X_t^0 = x_0 + dW_t^0 + \frac{1}{2} L_0 dt$ where L_0 is a measure supported on $\{t \geq 0; X_t^0 \in K\}$, and L_0 is an outer unit normal at X_t^0 for all such t . After time τ we just set $X_t^0 = X_t$. Since S_t is an orthogonal map (W_t^0) is a Brownian motion and thus (X_t^0) is a reflected Brownian motion starting from x_0 . Therefore $TV(\mu_x, \mu_{x_0}) \leq P(X_\tau \neq X_t^0) = P(\tau \leq t)$.

Observe that on $[0, \tau] \cap \tau_0$ $dW_t = dW_{t_0} = (I - S_t)(dW_t) = 2hV_t, dW_t \perp V_t$, where $V_t =$

$$X_t - X_{t_0} - X_t - X_{t_0} = \dots$$

So

$$d(X_t - X_{t_0}) = 2hV_t, dW_t \perp V_t \implies \int_{t_0}^t L(dt) + \int_{t_0}^t L_0(dt) = 2(dB_t) \cdot V_t \implies \int_{t_0}^t L(dt) + \int_{t_0}^t L_0(dt), \text{ where } Z$$

t

$$hV_s, dW_s \perp,$$

$$B_t =$$

on $[0, \tau]$.

0

Observe that (B_t) is a one-dimensional Brownian motion. Itô's formula then gives $dg(X_t - X_{t_0}) = 2h \cdot g(X_t - X_{t_0}), V_t \perp dB_t \implies h \cdot g(X_t - X_{t_0}), \int_{t_0}^t L(dt) + h \cdot g(X_t - X_{t_0}), \int_{t_0}^t L_0(dt) + 2 \int_{t_0}^t g(X_t - X_{t_0})(V_t, V_t) dt$, for every smooth function g on \mathbb{R}^n . Now if $g(x) = -x$ then $g(X_t - X_{t_0}) = V_t$ so $h \cdot g(X_t - X_{t_0}), V_t \perp = 1, h \cdot g(X_t - X_{t_0}), \int_{t_0}^t \int_{t_0}^t 0$ on the support of L , and $h \cdot g(X_t - X_{t_0}), \int_{t_0}^t \int_{t_0}^t 0$ on the support of L_0 . Moreover $\int_{t_0}^t g(X_t - X_{t_0}) = -X_t - Y \cdot P(X_t - Y_t)$ where P_x denotes the $t - t_0$ orthogonal projection on x . In particular $\int_{t_0}^t g(X_t - Y_t)(V_t) = 0$. We obtain $-X_t - X_{t_0} - \int_{t_0}^t -x - x_0 - 2B_t$, on $[0, \tau]$. Therefore $P(\tau \leq t) = P(\tau \leq t)$ where $\tau \leq 0$ is the first time the Brownian motion (B_t) hits the value $-x - x_0 - 2$. Now by the reflection principle $-x - x_0 - P(\tau \leq t) = 2P(0 \leq 2B_t \leq -x - x_0 - 2) = 2 \cdot 2 \cdot t$

2.4

From Wasserstein distance to total variation

To conclude it remains to derive a total variation bound between X_T and $X_{T+\tau}$ using (4). The details of this step are deferred to the supplementary material where we consider the case of a general logconcave distribution. The intuition goes as follows: the processes $(X_T + s)_{s \geq 0}$ and $(X_{T+\tau} + s)_{s \geq 0}$ both evolve according to a Brownian motion until the first time s that one process undergoes a reflection. But if T is large enough and τ is small enough then one can easily get from (4) (and the fact that the uniform measure does not put too much mass close to the boundary) that X_T and $X_{T+\tau}$ are much closer to each other than they are to the boundary of K . This implies that one can couple them (just as in Section 2.3) so that they meet before one of them hits the boundary.

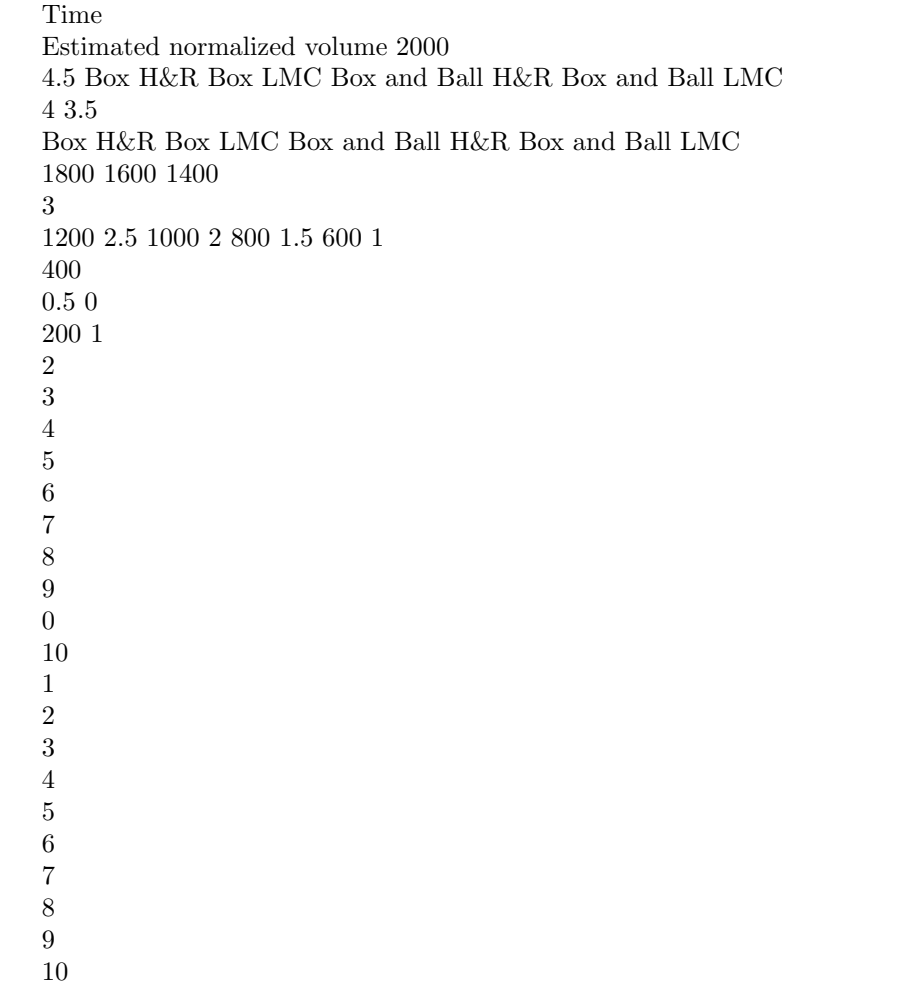
3

Experiments

Comparing different Markov Chain Monte Carlo algorithms is a challenging problem in and of itself. Here we choose the following simple comparison procedure based on the volume algorithm 7

developed in Cousins and Vempala [2014]. This algorithm, whose objective is to compute the volume of a given convex set K , proceeds in phases. In each phase ‘ it estimates the mean of a certain function under a multivariate Gaussian restricted to K with (unrestricted) covariance Σ . In . Cousins and Vempala provide a Matlab implementation of the entire algorithm, where in each phase

the target mean is estimated by sampling from the truncated Gaussian using the hit-and-run (H&R) chain. We implemented the same procedure with LMC instead of H&R, and we choose the step-size $\gamma = 1/(\ell n^2)$, where ℓ is the smoothness parameter of the underlying log-concave distribution (in particular here $\ell = 1/\epsilon^2$). The intuition for the choice of the step-size is as follows: the scaling in inverse smoothness comes from the optimization literature, while the scaling in inverse dimension squared comes from the analysis in the unconstrained case in Dalalyan [2014].



We ran the volume algorithm with both H&R and LMC on the following set of convex bodies: $K = [1, 1]^n$ (referred to as the "Box") and $K = [1, 1]^n \cap 2n B^n$ (referred to as the "Box and Ball"), where $n = 10 - k$, $k = 1, \dots, 10$. The computed volume (normalized by 2^n for the "Box" and by $0.2 \cdot 2^n$ for the "Box and Ball") as well as the clock time (in seconds) to terminate are reported in the figure above. From these experiments it seems that LMC and H&R roughly compute similar values for the volume (with H&R

being slightly more accurate), and LMC is almost always a bit faster. These results are encouraging, but much more extensive experiments are needed to decide if LMC is indeed a competitor to H&R in practice.

8

2 References

S. Ahn, A. Korattikara, and M. Welling. Bayesian posterior sampling via stochastic gradient fisher scoring. In ICML 2012, 2012. F. Bach and E. Moulines. Non-strongly-convex smooth stochastic approximation with convergence rate $\mathcal{O}(1/n)$. In Advances in Neural Information Processing Systems 26 (NIPS), pages 773–781, 2013. B. Cousins and S. Vempala. Bypassing kls: Gaussian cooling and an $\mathcal{O}(n^3)$ volume algorithm. Arxiv preprint arXiv:1409.6011, 2014. A. Dalalyan. Theoretical guarantees for approximate sampling from smooth and log-concave densities. Arxiv preprint arXiv:1412.7392, 2014. M. Dyer, A. Frieze, and R. Kannan. A random polynomial-time algorithm for approximating the volume of convex bodies. Journal of the ACM (JACM), 38(1):1–17, 1991. R. Kannan and H. Narayanan. Random walks on polytopes and an affine interior point method for linear programming. Mathematics of Operations Research, 37:1–20, 2012. L. Lovász and S. Vempala. Hit-and-run from a corner. SIAM J. Comput., 35(4):985–1005, 2006. L. Lovász and S. Vempala. The geometry of logconcave functions and sampling algorithms. Random Structures & Algorithms, 30(3):307–358, 2007. A. Nemirovski and D. Yudin. Problem Complexity and Method Efficiency in Optimization. Wiley Interscience, 1983. G. Pflug. Stochastic minimization with constant step-size: asymptotic laws. SIAM J. Control and Optimization, 24(4):655–666, 1986. H. Robbins and S. Monro. A stochastic approximation method. Annals of Mathematical Statistics, 22:400–407, 1951. A. Skorokhod. Stochastic equations for diffusion processes in a bounded region. Theory of Probability & Its Applications, 6(3):264–274, 1961. H. Tanaka. Stochastic differential equations with reflecting boundary condition in convex regions. Hiroshima Mathematical Journal, 9(1):163–177, 1979. L. Tweedie and G. Roberts. Exponential convergence of langevin distributions and their discrete approximations. Bernoulli, 2(4):341–363, 1996. M. Welling and Y.W. Teh. Bayesian learning via stochastic gradient langevin dynamics. In ICML 2011, 2011.

9