

# Some aspects of robustness in modern Markov Chain Monte Carlo

Sam Power<sup>1</sup> and Giorgos Vasdekis<sup>2</sup>

<sup>1</sup>School of Mathematics, University of Bristol

<sup>2</sup>School of Mathematics, Statistics and Physics, Newcastle University

November 27, 2025

## Abstract

Markov Chain Monte Carlo (MCMC) is a flexible approach to approximate sampling from intractable probability distributions, with a rich theoretical foundation and comprising a wealth of exemplar algorithms. While the qualitative correctness of MCMC algorithms is often easy to ensure, their practical efficiency is contingent on the ‘target’ distribution being reasonably well-behaved.

In this work, we concern ourselves with the scenario in which this good behaviour is called into question, reviewing an emerging line of work on ‘robust’ MCMC algorithms which can perform acceptably even in the face of certain pathologies.

We focus on two particular pathologies which, while simple, can already have dramatic effects on standard ‘local’ algorithms. The first is *roughness*, whereby the target distribution varies so rapidly that the numerical stability of the algorithm is tenuous. The second is *flatness*, whereby the landscape of the target distribution is instead so barren and uninformative that one becomes lost in uninteresting parts of the state space. In each case, we formulate the pathology in concrete terms, review a range of proposed algorithmic remedies to the pathology, and outline promising directions for future research.

## 1 Introduction

Markov Chain Monte Carlo (MCMC) is a computational approach to sampling from a probability measure which is specified in terms of an unnormalised density function, a task which is ubiquitous in various facets of Statistics, Machine Learning, Signal Processing, and Computational Physics. Given a ‘target’ measure  $\pi$  living on a reasonable state space  $\mathcal{X}$ , the modern MCMC toolbox by now contains a range of ‘standard’ methods for solving the sampling problem (i.e. producing approximate samples from  $\pi$ ), with some illustrative methods including Random Walk Metropolis [127, 5], Gibbs Sampling [57, 7], Langevin Monte Carlo [124, 43, 49, 50], Hamiltonian Monte Carlo [111, 20], Hit-and-Run [10, 126], Slice Sampling [110, 109, 120], and more. A general trend within the area is to seek general-purpose methods, which can be applied to a wide range of targets, without needing to focus on low-level details of the problem, and still expect reasonable (if not necessarily optimal) performance.

While basic MCMC often ‘works out of the box’ for well-behaved target distributions – indeed, this is somehow the success story of MCMC at large – there is an increasing realisation that problems of contemporary interest are prone to deviate from this good behaviour in various ways, leading the reliability of these default algorithms to be called into question somewhat. The present paper thus seeks to review current practice in MCMC, with a focus on modern developments in ‘robust’ algorithms which are able to perform well, even in the face of somewhat ‘pathological’ target distributions.

The format of the paper is then as follows: in Section 1.1 we present the basic notation we will be using, in Section 2, we recall the construction of some prominent MCMC methods, and the associated principles

for designing such methods. We then give a brief interlude describing some pathologies which can cause these prominent methods to break down. In Section 3, we focus on the pathology of ‘roughness’ of the target distribution, illustrating how it can arise, and reviewing a number of strategies for resolving it in practice. In Section 4, we focus on the sibling pathology of ‘heavy-tailedness’ of the target distribution, similarly treating how and why it arises, and how and why it can be addressed. We then conclude with some review of the area as a whole, pointing to some meaningful open problems and suggestions for future research.

## 1.1 Notation

We will be using  $\pi$  to denote a probability measure on  $\mathbf{R}^d$  and when the distinction is clear, we will be abusing notation, using  $\pi$  to denote the density of the measure with respect to the Lebesgue measure (assuming it exists). With that in mind, we will be writing

$$\pi(x) = \frac{1}{Z} \exp(-U(x)),$$

where  $U$  is defined up to an additive constant, and  $Z$  is the normalising constant, which in most settings will be assumed to be unknown, or at least difficult to calculate. To indicate this, we will sometimes be writing

$$\pi(x) \propto \exp(-U(x))$$

to indicate that  $\pi$  is known up to a multiplicative constant.  $U$  will sometimes be called *the potential*.  $\mu$  will typically denote a measure on  $\mathbf{R}^{2d}$ , having  $\pi$  as its first marginal. Typically,  $\mu$  will be constructed by augmenting  $\mathbf{R}^d$  with an auxiliary variable  $v \in \mathbf{R}^d$ ; in that setting, we will write  $\psi$  for the law of  $v$ , then writing  $K = -\log \psi$  for the associated potential.  $P, Q$  will generally denote Markov kernels on  $\mathbf{R}^d$ , which we will again occasionally conflate with their densities with respect to the Lebesgue measure. We write  $\text{TV}(\nu, \nu')$  to denote the total variation distance between probability measures  $\nu$  and  $\nu'$ .  $\mathcal{N}(m, C)$  denotes the Gaussian distribution with mean vector  $m$  and covariance matrix  $C$ .

$\|\cdot\|$  will denote the Euclidean norm in  $\mathbf{R}^d$ , and  $\|\cdot\|_1$  will denote the 1-norm.

For a continuous-time path  $(x_t)_{t \geq 0} \in \mathbf{R}^d$ , we write  $\dot{x}_t$  or  $\dot{x}$  to denote the derivative with respect to time  $t$ .  $\nabla$  will be used to denote the usual Euclidean gradient, with  $D$  sometimes being used for the same notion in the univariate case.

## 2 Design of MCMC Procedures

### 2.1 MCMC Basics

In concrete terms, the setup and goals of MCMC are as follows, focusing on the (dominant) setting of sampling on  $\mathcal{X} = \mathbf{R}^d$ : one is interested in approximate sampling from the probability measure  $\pi \in \mathcal{P}(\mathbf{R}^d)$ , whose density with respect to Lebesgue measure is known, at least up to a possibly-unknown normalising constant. As a standard abuse of notation, we will also use  $\pi$  for this density.

Given  $\pi$ , the MCMC architecture should design a Markov kernel  $P$  such that the Markov chain on  $\mathbf{R}^d$  which is driven by  $P$  will converge in distribution to this  $\pi$ . Typically, this is achieved by at least imposing that the kernel  $P$  leaves  $\pi$  *invariant*, i.e.

$$x \sim \pi, \quad y \sim P(x, \cdot) \implies y \sim \pi \quad \text{marginally},$$

and in practice, it is common to even impose that the kernel  $P$  be *reversible* with respect to  $\pi$ , in the sense that

$$x \sim \pi, \quad y \sim P(x, \cdot) \implies (x, y) \stackrel{d}{=} (y, x),$$

or at least that some similarly ‘local’ invariance condition is satisfied. Of course, each of these properties is satisfied by the trivial kernel  $P(x, dy) = \delta(x, dy)$  (which is certainly not fit for purpose as an MCMC kernel!), and so practical kernels must generally also satisfy some non-degeneracy conditions, to the tune of irreducibility, aperiodicity, and a meaningful form of *ergodicity*. For the latter point, a minimal requirement is that for a suitably large set of  $x \in \mathbf{R}^d$ , it holds that

$$\text{as } n \rightarrow \infty, \quad \text{TV}(P^n(x, \cdot), \pi) \rightarrow 0,$$

i.e. that for more-or-less arbitrary initialisations, the chain will indeed converge towards  $\pi$  in total variation distance. Beyond this qualitative ergodicity, a more demanding condition which should be satisfied by ‘good’ MCMC kernels is *exponential ergodicity*, whereby for a similarly large set of  $x \in \mathbf{R}^d$ , one can even write that

$$\text{as } n \rightarrow \infty, \quad \text{TV}(P^n(x, \cdot), \pi) \leq \rho^n \cdot V(x), \quad (1)$$

where the implied rate constant  $\rho$  is taken uniform in  $x$ , but the prefactor  $V(x)$  is allowed to depend on  $x$ .

Given an exponentially-ergodic Markov kernel  $P$  which is  $\pi$ -invariant, a number of desirable consequences hold, e.g. for ergodic averages along the MCMC trace, one can obtain Laws of Large Numbers, Central Limit Theorems, Concentration Inequalities for suitable expectands, and so on. In this respect, exponential ergodicity reassures the MCMC user that their method is fit for purpose, at least in principle (of course, in practice, if the constant  $\rho$  is unacceptably close to 1, then these reassurances are of little use). In this regard, we view exponential ergodicity as a ‘gold standard’ of sorts for separating ‘bad’ and ‘good’ kernels, coarse as the distinction may be.

While there are many possible ways in which to construct such a process, a prevailing strategy for constructing such processes in practice is to identify some continuum process which is known to be precisely  $\pi$ -invariant, and then find some way to simulate it in practice. Of course, other principles for deriving MCMC algorithms do exist, but cannot be our focus here.

Upon identifying such a continuum process, a number of options arise.

In certain situations, the continuum process is not just an idealised object, but can actually be realised algorithmically. A simple instance of this comes with Markov jump processes (MJPs), which can typically be simulated exactly by use of the Doob-Gillespie algorithm (initially introduced as a practical method in [61] for biochemical simulations; see also [1]). For Piecewise-Deterministic Markov Processes (PDMPs), exact simulation is possible in some cases, provided that the dynamics are tractable and that certain inhomogeneous Poisson point processes are feasible to simulate (following e.g. [90]; see also discussion in [30, 22] and similar works). For Stochastic Differential Equations (SDEs) of Itô type, exact simulation tends to only be possible for rather stylised models, either those with substantial analytic tractability (e.g. Gauss-Markov processes), or satisfying strong a priori bounds; see e.g. [19, 17, 24]. For SDEs with more general driving noises, the situation is typically even worse.

Moving beyond the somewhat restrictive class of processes for which exact simulation is feasible, the world of numerical analysis comes to our rescue, and provides us with an ample set of tools for approximate-yet-accurate simulation of a much wider class of processes. While these approximations typically incur some bias in the invariant measure of the numerical process, there are various strategies for quantifying and correcting for this bias, at least to the point that the impact of the bias falls below that of the Monte Carlo variance. Among such strategies, we highlight Multilevel Monte Carlo [60] as an archetypal example with wide and relatively straightforward applicability. For large MJPs, Gillespie’s algorithm can become costly to implement exactly, and so approximation strategies such as tau-leaping [62] become appealing (see also [75] for other ways of approximating a certain relevant class of MJPs). For PDMPs, circumstances can dictate that an approximate treatment is necessitated for either the dynamics, the jumps, or both of them. When the dynamics are tractable, it is natural to consider e.g. piecewise-constant or piecewise-linear approximations to the event rates (as in [114, 41, 3]). When the dynamics are also intractable, then the use of numerical methods for ODEs like splitting schemes becomes appealing; this was pursued in [14]. For SDEs, the literature on numerical approximation is yet more vast (see e.g. [119, 72, 107, 51, 71]), though it

is fair to say that in the context of MCMC simulation, the set of popular strategies is plainly dominated by the Euler-Maruyama method and various splitting schemes.

Within certain communities (the Bayesian statistical community perhaps chief among them), there is a cultural preference to avoid questions of asymptotic bias entirely by constructing implementable Markov kernels whose algorithmic invariant measure coincides exactly with the intended target measure. A generic approach to this is furnished by the ‘propose-accept/reject’ paradigm, as exemplified by [105, 68, 132, 4]. In this paradigm, the user ‘proposes’ a move according to some exactly-implementable Markov kernel, evaluates the favourability of the move according to some carefully-designed quantitative criterion, and then either ‘accepts’ the move to this proposed point, or ‘rejects’ it, staying in place. This strategy (which is known as the *Metropolis-Hastings filter* or simply “*Metropolisation*”) provides a rather generic approach to converting approximate numerical dynamics into dynamics which are genuinely ergodic with respect to  $\pi$ , eliminating concerns of asymptotic bias, albeit at the cost of slowing down the original dynamics with this conservatism; we term this the *Monte Carlo-Exact* paradigm for MCMC.

## 2.2 From Continuum Processes to Practical Algorithms

Having outlined these three high-level approaches to the design of MCMC kernels, in this sub-section, we review some continuum processes of interest, and detail some ways in which these processes have been translated into practical algorithms, following this rubric.

The aim is not to cover all possible processes which could be considered, nor to review all algorithms which have been derived in this manner. Our goal is to instead focus on processes and algorithms which are of clear popular interest, have been adopted as ‘default’ methods to some extent or another, and whose merits and limitations are relatively well-understood at a rigorous level.

### 2.2.1 The Overdamped Langevin Diffusion

Given a continuously-differentiable distribution  $\pi$  on a Euclidean space  $\mathbf{R}^d$ , the *overdamped Langevin diffusion* targeting  $\pi$  is the Itô SDE given by

$$dX_t = \nabla \log \pi(X_t) dt + \sqrt{2} dW_t \quad (2)$$

where  $W_t$  is a usual Wiener process on  $\mathbf{R}^d$ . Under rather mild conditions on  $\pi$ , this process is both reversible and ergodic with respect to  $\pi$ , and exponential ergodicity is assured under various possible conditions, the most restrictive (and informative) of which roughly asks that  $\pi$  be asymptotically approximately log-concave at infinity, or at least similarly confining.

Outside of a handful of rather simple cases, exact simulation of this process is challenging. The most general treatments tend to impose some rather strong conditions on the growth of  $\nabla \log \pi$ , and tend to only apply in settings of quite low (effective) dimension; these are not really seen to be viable approaches in the MCMC context (though have found value in other contexts of statistical inference; see e.g. [18, 55, 115]).

As concerns approximate numerical simulation, the options become somewhat more appealing. Application of a straightforward Euler-Maruyama discretisation to the Langevin process leads to the ‘Unadjusted Langevin Algorithm’ (ULA), also ‘Langevin Monte Carlo’ (LMC), which takes the form

$$\text{ULA: } X_n = X_{n-1} + h \cdot \nabla \log \pi(X_{n-1}) + \sqrt{2 \cdot h} \cdot \xi_n, \quad \xi_n \sim \mathcal{N}(0, \mathbf{I}_d) \quad (3)$$

(with  $h > 0$  a step-size parameter). This has been studied extensively (see e.g. [43, 49, 50, 39, 47]), and represents a reliable baseline method in various contexts where high accuracy is not a first priority.

Finally, to take the Langevin diffusion into the Monte Carlo-Exact paradigm, application of the Metropolis-Hastings device to the aforementioned Euler-Maruyama discretisation yields the so-called *Metropolis-Adjusted*

*Langevin Algorithm (MALA)* (introduced variously in [125, 16]), well-studied in works including [124, 26, 40, 138]. In this method, writing  $Q$  for the transition kernel corresponding to the ULA update, one simulates a ‘proposal’ move  $Y_n \sim Q(X_{n-1}, \cdot)$ , evaluates the *Metropolis-Hastings* acceptance probability  $\alpha(X_{n-1}, Y_n)$ , where

$$\alpha(x, y) = \min \left\{ 1, \frac{\pi(y) \cdot Q(y, x)}{\pi(x) \cdot Q(x, y)} \right\}, \quad (4)$$

and then moves to  $Y_n$  with this probability, otherwise remaining at  $X_{n-1}$  (that is  $X_n$  is either set to be  $Y_n$  or  $X_{n-1}$ , with probabilities  $\alpha(X_{n-1}, Y_n)$  and  $1 - \alpha(X_{n-1}, Y_n)$  respectively).

To showcase the performance of the MALA algorithm, let us consider a two-dimensional Gaussian distribution with mean  $m = [0, 0]^\top$  and covariance matrix  $C = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$ , meaning that the density is of the form

$$\pi(x) = \frac{1}{2\pi\sqrt{\det(C)}} \cdot \exp\left(-\frac{1}{2}x^\top C^{-1}x\right), \quad x \in \mathbf{R}^2. \quad (5)$$

We ran MALA with step-size  $h = 0.35$  on this target for  $N = 10^4$  iterations, starting from  $x_0 = (4, 5)$ , rather far from the mode. The step-size was chosen so that the algorithm accepts between 50 – 60% of the proposed jumps, as recommended by e.g. [122]. Our results are summarised in Figure 1, where we present four plots. The first is the scatter plot, which shows the path of the process in  $\mathbf{R}^2$ . The second plot is the traceplot of the first coordinate against time, with the marginal density plotted on the right for reference. From these two plots, we can see that the process captures the shape of the target very efficiently. The third plot is the autocorrelation plot, showing how correlated the algorithm’s samples are as a function of the time increment between them. The plot shows a fairly rapid decay of this correlation to zero as time-distance between samples increases, which indicates good algorithmic performance (noting that one would ideally have i.i.d. samples from the target, with no autocorrelations whatsoever). Finally, the fourth plot is the histogram of the samples of the first coordinate, with the true marginal distribution overlaid with red colour for reference.

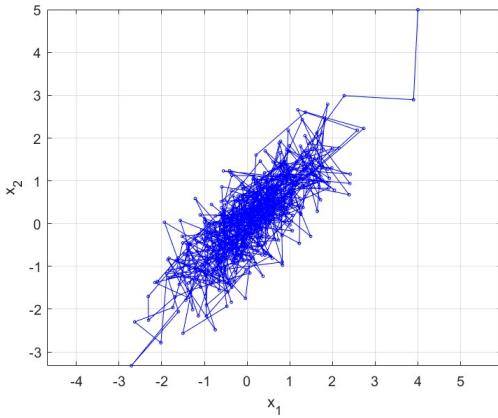
Finally, in terms of using numerical approximations of diffusions to motivate MCMC algorithms, we note in passing that while the seminal Random Walk Metropolis Algorithm (see e.g. [127, 5] is seldom presented as a numerical method for simulating the Langevin diffusion *per se*, there exist a range of theoretical results (e.g. [58, 59]) which reveal that it can be reasonably interpreted in this way.

### 2.2.2 Piecewise-Deterministic Markov Processes

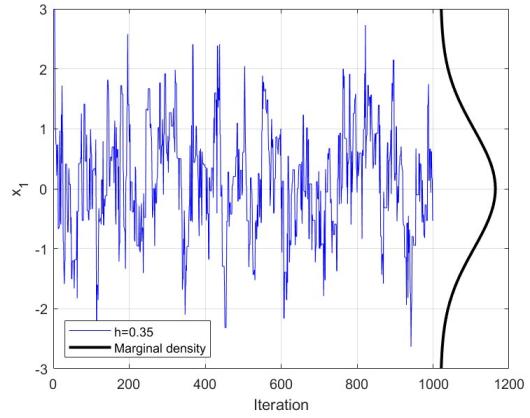
In informal terms, a Piecewise-Deterministic Markov Process (PDMP) is a continuous-time Markov process evolving in a continuous space whose dynamics comprise deterministic evolution along the flow of an ordinary differential equation, punctuated by instantaneous Markovian jumps. Although initially introduced by [44] as a model for problems in Operations Research and Control, this class of processes has undergone a sort of renaissance in the last decade, as their application to MCMC simulation has been examined in greater depth.

For the purposes of this review, we focus our attention on three prototypical PDMPs for Monte Carlo applications: *Refreshed Hamiltonian Dynamics*, the *Bouncy Particle Sampler*, and the *Zig-Zag Process*. In each of these, the state undergoing dynamics is a kinetic particle, i.e. a position accompanied by a velocity, i.e.  $z = (x, v) \in \mathbf{R}^d \times \mathbf{R}^d$ , designed so that at equilibrium, the particle  $z$  is distributed according to  $\mu(dz) = \pi(dx) \cdot \psi(dv)$ , where  $\psi$  is some simple, known distribution over a suitable ‘velocity space’. Interpreting  $v$  as the velocity of the system, we refer to  $K(v) := -\log \psi(v)$  as the *Kinetic energy*, with  $U$  being called the *potential energy*.

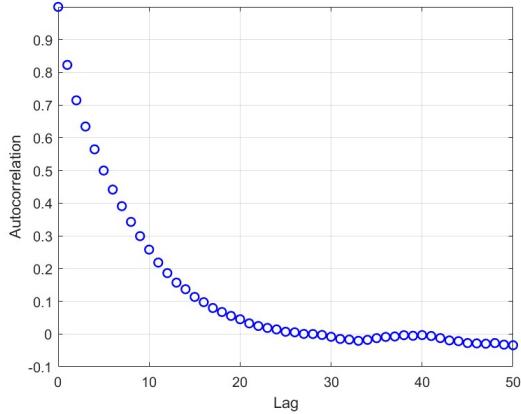
#### Refreshed Hamiltonian Dynamics



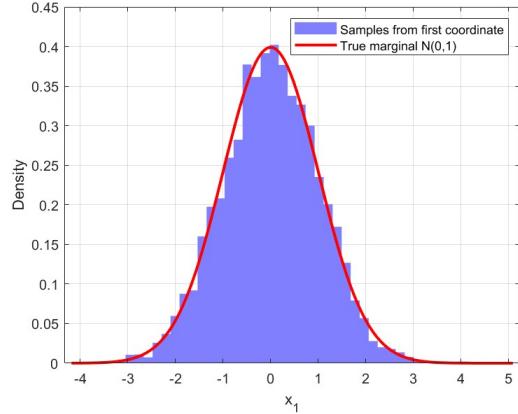
(a) Scatter plot



(b) Traceplot of first coordinate samples



(c) Autocorrelation of first coordinate samples



(d) Histogram of first coordinate samples

Figure 1: MALA algorithm on a two-dimensional correlated Gaussian target (5). Step-size  $h = 0.35$ . Number of iterations  $N = 10^4$ . Starting point  $x_0 = (4, 5)$ .

For Refreshed Hamiltonian Dynamics (introduced by [27] as ‘Randomised Hamiltonian Monte Carlo (RHMC)’<sup>1</sup>), one takes  $\psi(dv) = \mathcal{N}(dv; 0, \mathbf{I}_d)$ , meaning that  $K(v) = \frac{1}{2}v^2$ , and the particle moves according to Hamilton’s equations with ‘Hamiltonian’  $\mathcal{H}(z) = -\log \mu(z)$ , i.e.

$$\dot{x} = v, \quad \dot{v} = \nabla \log \pi(x), \quad (6)$$

punctuated by jumps in the velocity of the form

$$v' \sim \mathcal{N}(dv'; \rho \cdot v, (1 - \rho^2) \cdot \mathbf{I}_d)$$

for some correlation parameter  $\rho \in [-1, 1]$ , occurring at some time-homogeneous Poisson rate  $\lambda > 0$ .

Outside of very simple cases, exact simulation of the Hamiltonian dynamics is intractable, and so this is rarely treated in practice. Fortunately, numerical methods for Hamiltonian dynamics are rather well-developed, and have been especially well-served by the field of Geometric Numerical Integration [65]. The dominant approach to approximate simulation of Refreshed Hamiltonian Dynamics involves use of the splitting integrator of

<sup>1</sup>We prefer here to distinguish the idealised physical process with its application to simulation, cf. the use of ‘Langevin diffusion’ and ‘Langevin Monte Carlo’.

(Størmer-Verlet, Leapfrog, Strang, etc.) to solve the dynamics, with the jump times and events handled exactly by simulating exponential random variables.

Monte Carlo-Exact resolutions of Refreshed Hamiltonian Dynamics are ultimately rather similar in character; one again simulates the dynamics approximately using the Verlet integrator (with step-size  $h$ ) for a randomly-chosen number of steps (e.g.  $L \sim \text{Poisson}(\lambda \cdot h^{-1})$ ) and then accepting or rejecting the final state  $z_L$  with a suitable probability, following the Metropolis-Hastings prescription (as in (4)); this is typically what is meant by (Metropolised) ‘Hamiltonian Monte Carlo’. A number of variations of Hamiltonian Monte Carlo exist, with many intriguing idiosyncrasies, but for the cohesion of presentation, we focus on this one.

Below we present a numerical study targeting the two-dimensional Gaussian distribution (5) of the previous sub-section. In the special case of the Gaussian distribution, the solution to (6) is given in closed form, and one can directly implement the Randomised Hamiltonian Monte Carlo (RHMC) algorithm without any need for numerical discretisation or Metropolis-adjustment. As in the previous sub-section, we ran the algorithm for  $N = 10^4$  iterations, starting from  $x_0 = (4, 5)$ . The correlation parameter and refreshment rate were chosen as  $\rho = 0$  and  $\lambda = 0.2$  respectively. We present our results in Figure 2, which include the scatter plot and the traceplot of the first coordinate, along with the autocorrelation plot and the histogram of the first coordinate samples. Interestingly, the correlations between samples appear to decrease significantly faster than in MALA; at the intuitive level, one could attribute this to the ‘momentum-driven’ dynamics of RHMC, which can allow the process to explore the space faster.

### Bouncy Particle Sampler and the Zig-Zag Process

By contrast to Refreshed Hamiltonian Dynamics, the Bouncy Particle Sampler (BPS) instead works with ODE dynamics which are agnostic to the target distribution  $\pi$ , but introduces jump rates and jump types which depend intimately on the details of  $\pi$ . To be precise, the dynamics in question are simple free transport, i.e.

$$\dot{x} = v, \quad \dot{v} = 0.$$

The jump rates are no longer time-homogeneous, but they depend on the instantaneous position and velocity of the process, influenced by the rate at which the log-density of the target is changing along these dynamics. In particular, the rate at  $(x, v)$  is given by

$$\lambda(x, v) = \max\{0, \langle v, -\nabla \log \pi(x) \rangle\}.$$

Upon the occurrence of a jump event, the position  $x$  stays in place, but the velocity  $v$  undergoes a specular reflection against the level sets of the log-density, i.e.

$$v' = \left( \mathbf{I}_d - 2 \cdot \frac{(\nabla \log \pi(x)) (\nabla \log \pi(x))^\top}{(\nabla \log \pi(x))^\top (\nabla \log \pi(x))} \right) \cdot v,$$

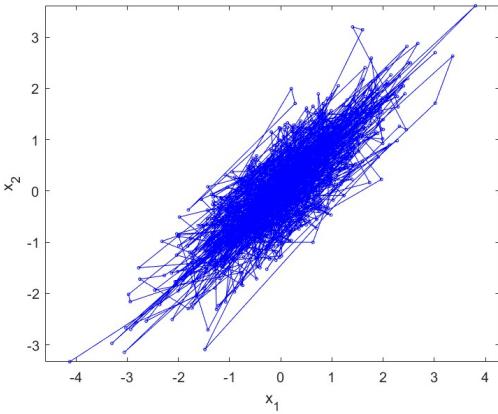
which forms the eponymous ‘bounce’ and preserves the uniform measure on the sphere  $\mathbf{S}^{d-1}$  according to which the velocities are typically distributed.

The Zig-Zag Process has a similar character, using the same deterministic dynamics, but differing in the domain of the velocities and the nature of the jumps. In particular, the velocity of the Zig-Zag Process takes values in  $\{\pm 1\}^d$ , and each jump event simply ‘flips’ one of the  $d$  components. More precisely, for  $1 \leq i \leq d$ , at rate

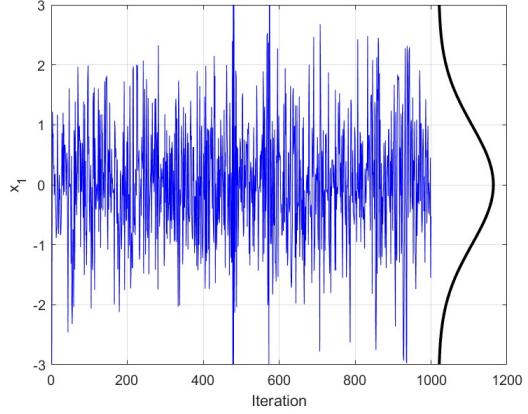
$$\lambda_i(x, v_i) = \max\{0, -v_i \cdot \partial_i \log \pi(x)\},$$

one observes a ‘jump event of type  $i$ ’, at which the position  $x$  stays in place, the values  $\{v_j : 1 \leq j \leq d, j \neq i\}$  stay in place, and the  $i$ th component  $v_i$  negates its own value, jumping to  $-v_i$ . This leads to the eponymous ‘zig-zagging’ trajectories from which the process takes its name.

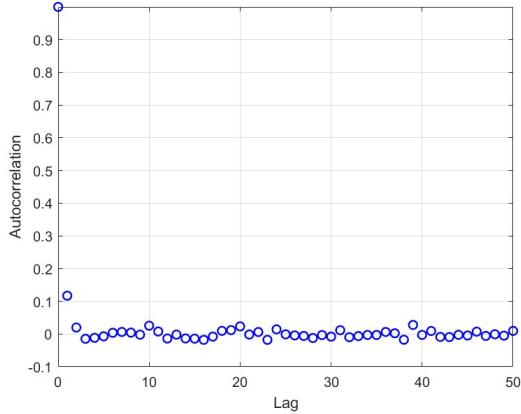
The two processes occupy rather similar locations vis-à-vis their practical implementability. Due to the simplicity of the deterministic dynamics, exact simulation of both the BPS and the ZZP is rather more



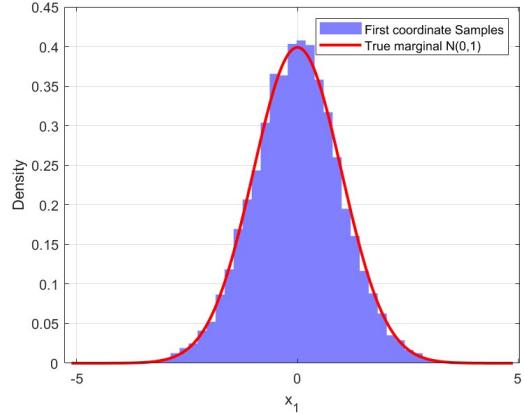
(a) Scatter plot



(b) Traceplot of first coordinate samples



(c) Autocorrelation of first coordinate samples



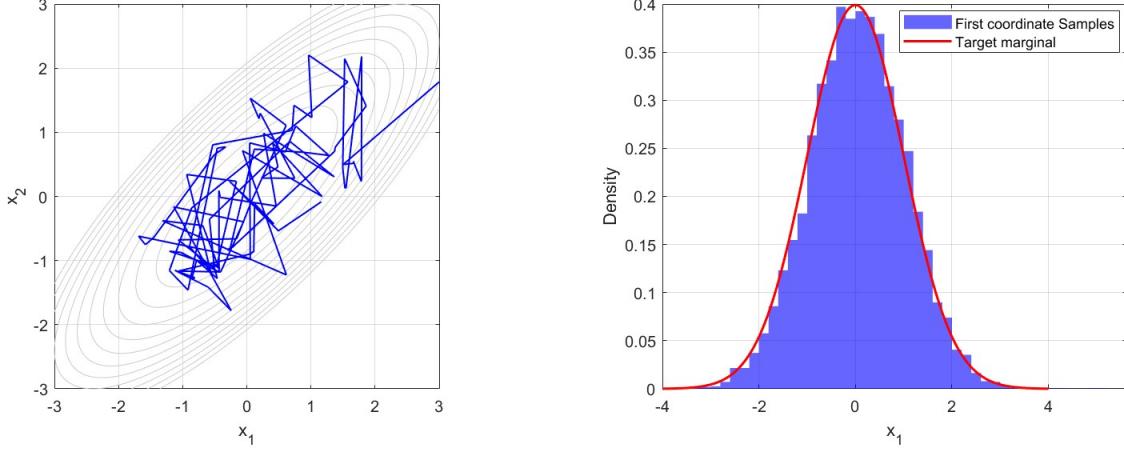
(d) Histogram of first coordinate samples

Figure 2: Randomized Hamiltonian Monte Carlo (RHMC) algorithm on a two-dimensional correlated Gaussian target (5). Correlation parameter  $\rho = 0$ . Rate  $\lambda = 0.2$ . Number of iterations  $N = 10^4$ . Starting point  $x_0 = (4, 5)$ .

feasible than for Refreshed Hamiltonian Dynamics. The practical difficulty arises instead in the simulation of these jump events, which requires the simulation of an inhomogeneous-in-time Poisson process. While there are a range of target distributions for which this task is achievable, in general, one requires some a priori understanding of the shape of  $\log \pi$  to make this efficient. A number of numerical discretisation strategies have been proposed (see e.g. [114, 41, 3]) which are more readily applicable to generic problems, and can be observed to incur an acceptable discretisation bias in many cases. The recent work of [38] proposes a Metropolis-type wrapper which allows to ‘sanitise’ some of these discretisations for an audience who demand Monte Carlo-Exactness, removing the bias in the invariant measure with a suitable accept-reject mechanism.

We consider the two-dimensional correlated Gaussian (5) of the previous sub-sections to showcase the performance of the Bouncy Particle Sampler (BPS) and the Zig-Zag Sampler (ZZS) via a numerical simulation. We ran both algorithms for  $N = 10^4$  direction switches, starting from  $x_0 = (4, 5)$ . For the BPS we used refresh rate  $\lambda = 0.66$  (using methodological guidance from [30]), while for ZZS we used  $\lambda = 0$ , which is conjectured to be optimal in terms of minimising the asymptotic variance (e.g. [21]). In Figure 3 we present the scatter plot of the BPS for the first 100 direction switches, and the histogram of the first coordinate samples with the marginal target density overlaid. For the scatter plot, we also add the level sets of the target on the background to indicate how the direction switches work as a result of bouncing on the level

sets. We present the same plots for the Zig-Zag process in Figure 4.



(a) Scatter plot of the first 100 direction switches. Level sets on the background.

(b) Histogram of the first coordinate

Figure 3: Bouncy Particle Sampler (BPS) on a bivariate correlated Gaussian target (5). Number of direction switches:  $N = 10^4$ . Starting point:  $x_0 = (4, 5)$ . Refresh rate  $\lambda = 0.66$ .

### 2.2.3 Omissions

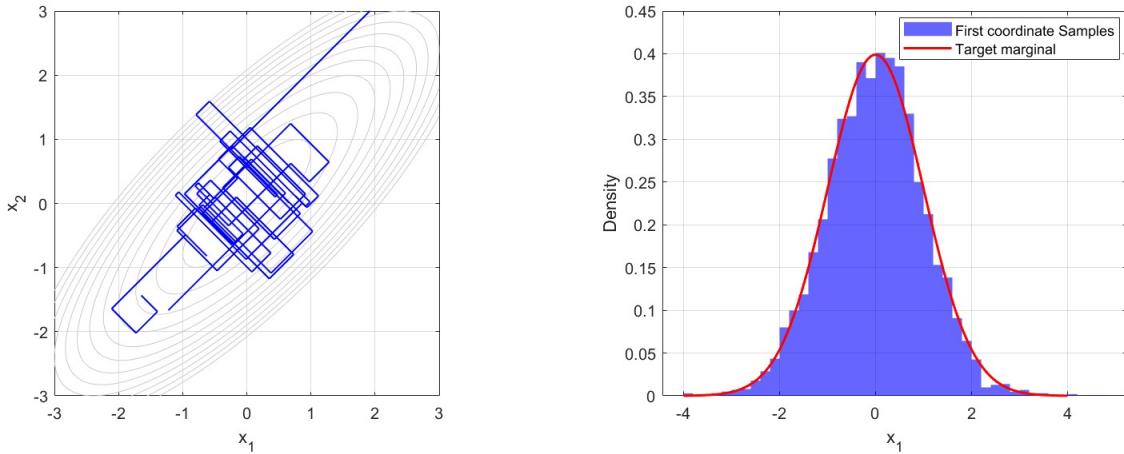
We pause for a minute here to discuss some other prominent MCMC algorithms which will not feature as explicitly in our subsequent discussion.

The first omission is the Random Walk Metropolis (RWM) algorithm [105, 127, 5]. Due to the aforementioned interpretation of the RWM as a zeroth-order discretisation of the Langevin diffusion, one might anticipate that aside from generally moving more slowly than gradient-based discretisations, it will not really encounter any additional difficulties, and all available evidence suggests that this is indeed the case. For similar reasons, while Monte Carlo methods based around the Underdamped (also ‘Kinetic’, ‘Second-Order’, etc.) Langevin Diffusion [37] are of practical interest, their qualitative features are arguably subsumed by a combination of the Overdamped Langevin Diffusion and the Refreshed Hamiltonian Dynamics, and so neither do they present unique challenges in terms of robustness. The same comments apply equally to other methods based upon similar ODEs and SDEs.

The other class of notable omissions will be MCMC methods whose dynamics are ‘non-local’ in some sense. The chief example of this class is the Gibbs sampler [57, 141], and its many descendants based on principles of iterative conditional simulation, including Hit-and-Run [10, 126], Slice Sampling [110, 126], and various other ‘auxiliary variable’ methods. We choose to spare these methods from discussion on the grounds that the qualitative picture is rather different here, to the effect that non-local samplers face different challenges to their robustness, thus requiring rather different solutions. As such, rather than introducing further notational and conceptual overhead to this chapter, we prefer to confine our discussion to methods based on local dynamics derived from the Langevin diffusion and the Refreshed Hamiltonian Dynamics.

## 2.3 Interlude

These methods are general-purpose, widely-used, and theoretically well-understood to work well under reasonable conditions. Conventional results tend to operate under the rather strong assumptions that the target distribution  $\pi$  be strongly log-concave, and that  $\log \pi$  have a Lipschitz-continuous gradient. More refined



(a) Scatter plot of the first 100 direction switches. Level sets on the background.

(b) Histogram of the first coordinate

Figure 4: Zig-Zag Sampler (ZZS) on a two-dimensional correlated Gaussian target (5). Number of direction switches:  $N = 10^4$ . Starting point:  $x_0 = (4, 5)$ .

results have been obtained under various weakenings of these assumptions, usually amounting to some sort of light-tailedness property for  $\pi$  combined with some relaxed notion of quantitative continuity for  $\nabla \log \pi$ . While it is difficult to obtain a tight characterisation of conditions on  $\pi$  which ensure quantitatively efficient sampling, experience dictates that conditions of this form are qualitatively appropriate.

Taking this as given, it bears mentioning that these assumptions can fail in contemporary statistical problems, even those which are relatively conventional. Sometimes, this is rather benign, and the practical performance of the methods remains acceptable, even in the absence of theoretical guarantees. Sometimes, however, the failure of these assumptions leads to a genuine breakdown of the methods, which fail to produce an acceptable output in a reasonable time-frame. This can come in various forms: the cost of each iteration of the method might blow up, the numerical stability of the method might fail, the range of stable step-sizes for a numerical discretisation might collapse, the rate of convergence as a function of the number of iterations might slow to a crawl, and so on.

This has prompted an emerging line of MCMC research that seeks to develop sampling algorithms that retain the general-purpose character of well-established methods, but have an additional character of ‘robustness’. In particular, a key aim is for the algorithm to perform effectively when the classical assumptions on  $\pi$  are satisfied, and to degrade gracefully when these assumptions fail. This should be understood less as a task of ‘acceleration’, but more of devising simple, lightweight strategies for avoiding potentially-catastrophic pitfalls.

In the next two sections, we will focus on two different types of pathologies that often arise in various practical statistical settings. The first one, ‘roughness’, relates to settings in which the target distribution has reduced regularity (e.g. when  $\nabla \log \pi$  is non-Lipschitz, non-differentiable, non-continuous, or similar). The second one, ‘heavy-tailedness’ concerns settings in which the target fails to concentrate its probability mass well in the bulk of the space, with substantial mass at a large distance from e.g. the mode. One can characterise the former as referring to bad ‘local’ behaviour of  $\pi$ , whereas the latter refers to bad ‘global’ behaviour of  $\pi$ . To indicate differences between these two types of distributions, in Figure 5 we show the growth of the negative log-density (also known as potential) for various light (typically tails lighter than Gaussian) and heavy (tails heavier than Laplace) densities.

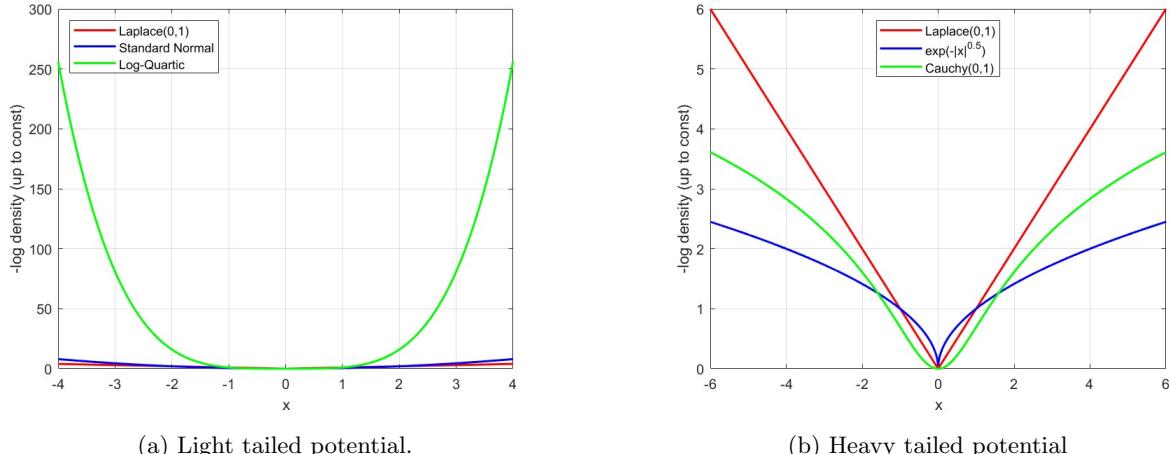


Figure 5: Plots of the negative log-densities of various densities. The left plot shows the growth of Laplace potential and of other densities with lighter tails. The right plot shows the growth for heavier tails.

### 3 Roughness

#### 3.1 Formulation of Pathology, Basic Examples

Our starting point in this section is that for target distributions  $\pi$  with good smoothness properties, if one starts with some  $\pi$ -invariant continuum process, then numerical discretisation essentially ‘works’. Broadly speaking, this entails that the numerical process inherits the pertinent stability properties of the continuum process, and that by taking sufficiently fine discretisations, the discretisation error can be reduced arbitrarily. In the context of Monte Carlo-Exact procedures, the consequences are similar, in that for sufficiently small step-sizes, acceptance probabilities can be increased to an appropriate level, and the resulting Markov chain is not too ‘sticky’, i.e. the chain ‘actually moves’ on most steps. This is somehow a minimal requirement for a discrete-time Markov chain to be performant; various elementary arguments can quantify the intuition that a chain which moves only rarely can only converge slowly.

Concretely, ‘good smoothness properties’ tend to enforce that the target distribution varies locally in predictable ways. A usual assumption to this effect is ‘ $L$ -Smoothness’ of the potential  $U = -\log \pi$ , or equivalently, that the ‘force’  $\nabla U$  is  $L$ -Lipschitz, i.e.

$$x, y \in \mathbf{R}^d \implies \|\nabla U(x) - \nabla U(y)\| \leq L \cdot \|x - y\|$$

for some  $L \in (0, \infty)$ . This is sometimes weakened to a more general modulus of continuity assumption, i.e. that for some suitable  $\psi : \mathbf{R}_+ \rightarrow \mathbf{R}_+$ , it holds that

$$x, y \in \mathbf{R}^d \implies \|\nabla U(x) - \nabla U(y)\| \leq \psi(\|x - y\|).$$

Under an assumption of this form, various quantitative consequences can be derived. As an example, under these assumptions, [6] shows that in the Random Walk Metropolis algorithm, by tuning the step-size based on  $\psi$  and the ambient dimension  $d$ , one can ensure that uniformly over the state space, the acceptance rates out of any state can be uniformly bounded from below. Scaling analysis of a more specific model problem in [135] suggests that these estimates do not need to be overly conservative. Along similar lines, [39] show that for Langevin Monte Carlo, for polynomial  $\psi$  (corresponding to smoothness of Hölder type), one can design step-sizes so that the asymptotic bias of the numerical invariant measure is controlled well. In each case, one sees that ‘rougher’ target distributions are subject to more stringent assumptions on viable step-sizes, often in a dimension-dependent way.

Although this ‘modulus of continuity’ approach to smoothness is rather flexible, it does preclude certain interesting classes of target distribution. To this point, observe that although the modulus  $\psi$  is a priori quite general, there are some implied constraints on both its form and growth. In particular, note that any nontrivial control of the form

$$\|x - y\| \leq r \implies \|\nabla U(x) - \nabla U(y)\| \leq \Lambda$$

can be inductively bootstrapped into the global estimate

$$\|\nabla U(x) - \nabla U(y)\| \leq \Lambda \cdot \left\lceil \frac{\|x - y\|}{r} \right\rceil \leq \frac{\Lambda}{r} \cdot \{r + \|x - y\|\}.$$

As such, any  $\nabla U$  satisfying such an estimate is necessarily of at most *linear growth* at infinity, and hence corresponds to potentials of at most quadratic growth, i.e. with tails no lighter than Gaussian. One is then led to wonder what this implies about the performance of such algorithms when applied to such light-tailed targets: shall we observe reasonable performance which happens to be unsupported by the existing theory, or shall we observe genuine problems?

To this end, we introduce an illustrative toy example which will help demonstrate some of these features.

**Example 3.1** (‘Polynomially-Steep’ Potential). *Noting that any potential  $U$  for which  $\nabla U$  admits a non-trivial modulus of continuity is of at most quadratic growth at infinity, it is natural to seek counter-examples by constructing potentials of super-quadratic (but still polynomial) growth at infinity, e.g.  $U(x) = \|x\|^q$  for some  $q > 2$ . For concreteness (and to avoid issues about the existence and regularity of higher-order derivatives), we will generally focus on the quartic potential  $U(x) = \|x\|^4$ , i.e. the choice  $q = 4$ , and the induced target distribution  $\pi(x) \propto \exp(-\|x\|^4)$ . Note that this distribution remains log-concave and very well-confined, with tails much lighter than those of the Gaussian distribution.*

To this end, we perform some illustrative simulations on the target distribution associated to the quartic potential in  $d = 1$  (see Figure 6). We observe that when the step-size is not appropriately small ( $h = 0.1$ , magenta line), the Metropolis-Adjusted Langevin Algorithm fails to explore the state space, rapidly becoming ‘stuck’, rejecting every proposed move in our simulation run. Use of smaller step-sizes can delay this pathology to some extent, but any eventual excursion of the process to a larger value of  $x$  will lead to similar issues. At the same time, use of very small step-sizes ( $h = 10^{-4}$  red line,  $h = 10^{-3}$  green line) will lead to very slow space exploration, as indicated in Figure 6. This serves as strong empirical evidence that ‘steep’ gradients can be problematic for standard gradient-based MCMC algorithms, which is validated theoretically in other works; see e.g. [124, 26].

At the other end of the spectrum, one might ask what happens when the force does admit a bounded modulus of continuity, but one which does not vanish continuously at 0.

**Example 3.2** (‘Locally-Sharp’ Potential). *When the obstruction to smoothness is ‘local’ roughness rather than growth at infinity, an archetypal example is the ‘Laplace-type’ potential associated to the  $\ell^1$  norm, i.e.  $U(x) = \|x\|_1$  for  $x \in \mathbf{R}^d$ . In this setting, one computes that away from the coordinate axes,  $\nabla U(x) = \text{sign}(x)$ . By considering points in a neighbourhood of the origin, one sees that the best possible  $\psi$  is  $\psi(r) = 2d$  for  $r > 0$ . While this is favourable in terms of uniform boundedness, it is problematic in that one has no immediate guarantees to generate acceptable moves, even at a small step-size. Moreover, one sees that high dimensionality of  $x$  apparently makes matters even worse.*

To demonstrate some of these features, we consider the application of MALA to the associated Laplace target in dimension  $d = 10^4$ , i.e.

$$\pi(x) \propto \prod_{i=1}^{10^4} \exp(-|x_i|). \tag{7}$$

Figure 7 documents the results, showing the traceplot and histogram of the first coordinate of the algorithm run for  $10^5$  iterations. We considered two different step-sizes. The first row presents the results with step-size  $h = 0.01$ , chosen so that the algorithm accepts between 50 – 60% of the proposed jumps, as suggested

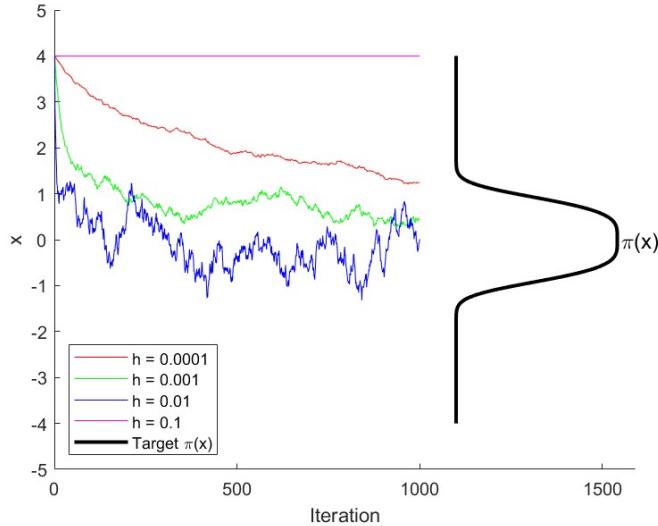


Figure 6: MALA trace plots with various step-sizes. The target density ( $\pi(x) \propto \exp(-x^4)$ ) is overlaid on the right hand side of the graph.

by the literature (see e.g. [118]). From the trace plot, the algorithm seems to be performing sufficiently well. However, when one changes the step-size, the algorithmic behaviour quickly deteriorates. To showcase this, we ran the algorithm for different step-sizes  $h \in \{0.01, 0.02, \dots, 0.1\}$ . For each step-size we ran 100 i.i.d. copies of the algorithm for  $N = 10^5$  iterations. We report the average (over iterations and algorithmic configurations) Acceptance Probability and Mean Squared Error when estimating the target expectation (in this case the true quantity is zero) as functions of step-sizes. It is evident that different step-sizes lead to drastically different behaviour, with any step-size larger than  $h = 0.01$  leading to poor algorithmic performance.

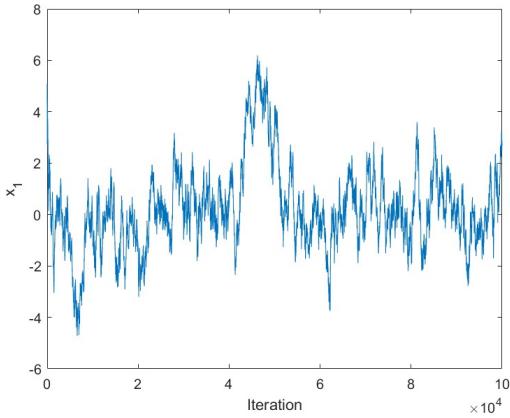
Finally, there are various other intermediate model problems which exhibit qualitatively similar pathologies, with densities which blow up at some natural boundary, as one observes in Beta-type distributions. In a similar category, there are targets which are subject to hard boundaries, e.g. Gaussian distributions subject to non-negativity constraints. These boundaries typically then create computational issues for MCMC algorithms such as MALA. Most algorithms fail to behave appropriately when close to boundary, either by not visiting it enough, or by proposing moves outside the bounded state space, which are ultimately rejected. To illustrate the problem, we introduce a stylised model problem.

**Example 3.3** (Divergent Potential with Natural Boundary). *Fix the domain  $\mathcal{B} = \{x \in \mathbf{R}^d : \|x\| < 1\}$ , and consider the potential  $U : \mathcal{B} \rightarrow \mathbf{R}$  which is given by  $U(x) = \log\left(\frac{1}{1-\|x\|^2}\right)$ . One sees immediately that  $U$  blows up near to the boundary of  $\mathcal{B}$ , as do all derivatives of  $U$ . When considering the associated probability measure ‘at positive temperature’, this blow-up corresponds to the vanishing of the target density; at ‘negative temperature’, it instead reflects a concentration of mass near this boundary (cf. Beta distributions with parameters in  $(0, 1)$ ). In either case, the blow-up of  $U$  and its derivatives can be expected to cause difficulties when discretising  $\pi$ -invariant dynamics.*

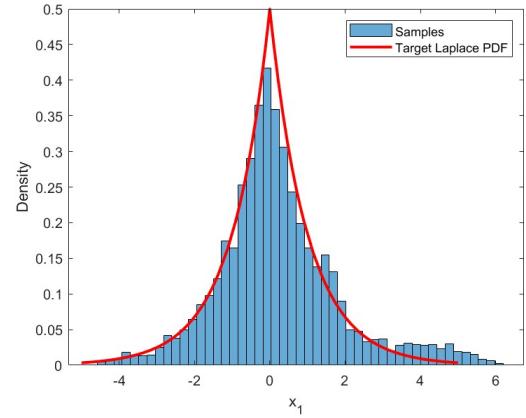
We present now a numerical example for this divergent potential in the univariate setting, at ‘temperature’  $\tau = -2$ , yielding the specific density

$$\pi(x) \propto (1-x^2)^{-\frac{1}{2}}, \quad x \in (-1, 1), \tag{8}$$

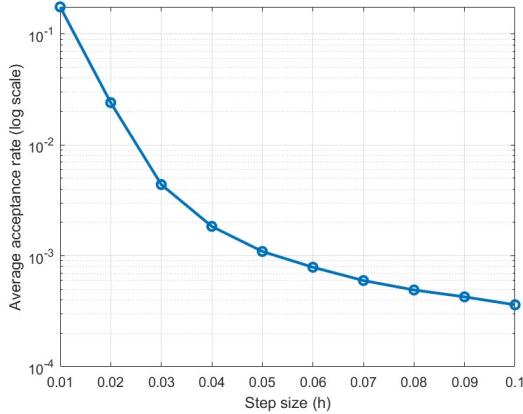
which blows up as one approaches  $x = \pm 1$ . The algorithm runs for  $N = 1000$  iterations, with step-size  $h = 0.1$  (again chosen so that around 50–60% of the proposed jumps are accepted), starting from  $x_0 = 0.5$ .



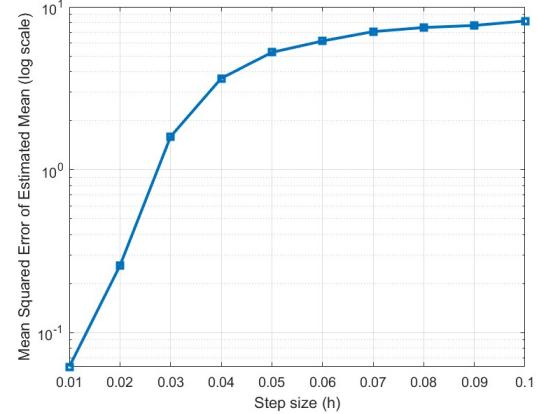
(a) Traceplot of first coordinate. Step-size  $h = 0.01$ .



(b) Histogram of first coordinate samples. Step-size  $h = 0.01$ .



(c) Average acceptance probability over step-sizes.



(d) Mean square error of target expectation over step-sizes

Figure 7: MALA on the  $10^4$ -dimensional Laplace target (7). First row: traceplot and histogram with optimally tuned step-size  $h = 0.01$ . Second row: Average acceptance probability and Mean squared error (estimating the target expectation) over step-size.

We present the traceplot and histogram of the MALA algorithm in Figure 8. It is evident that the process has not captured the target’s behaviour sufficiently well, and there are areas close to the boundary that are not explored at all. The traceplot also indicates that the process can get quite stuck close to the boundary. Closer examination of the algorithm output reveals that this results from frequent proposals outside of the domain  $(-1, 1)$ , which result in immediate rejections.

A closely-related scenario, not uncommon in applications, is when the potential itself is ‘nice’ *per se*, but the domain is constrained for other reasons.

**Example 3.4** (Nice Potential, Artificial Boundary). *A simple example of this setting would be a Gaussian distribution, subject to a collection of affine constraints, i.e. the distribution of  $X \sim \mathcal{N}(0, \mathbf{C})$ , subject to the constraints  $a_j^\top X \leq b_j$  for  $j$  in some finite index set  $\mathcal{J}$ , with  $a_j$  and  $b_j$  a suitable collection of vectors and scalars respectively. Such distributions arise naturally in the study of shape-constrained Gaussian processes, as well as for algorithmic reasons in data augmentation strategies for certain generalised linear models (GLMs). A particularly tractable instance (which nevertheless exhibits a number of relevant pathologies) is to consider a spherical Gaussian random variable  $X \sim \mathcal{N}(0, \sigma^2 \cdot \mathbf{I}_d)$ , constrained to the  $\ell^\infty$  box  $B =$*

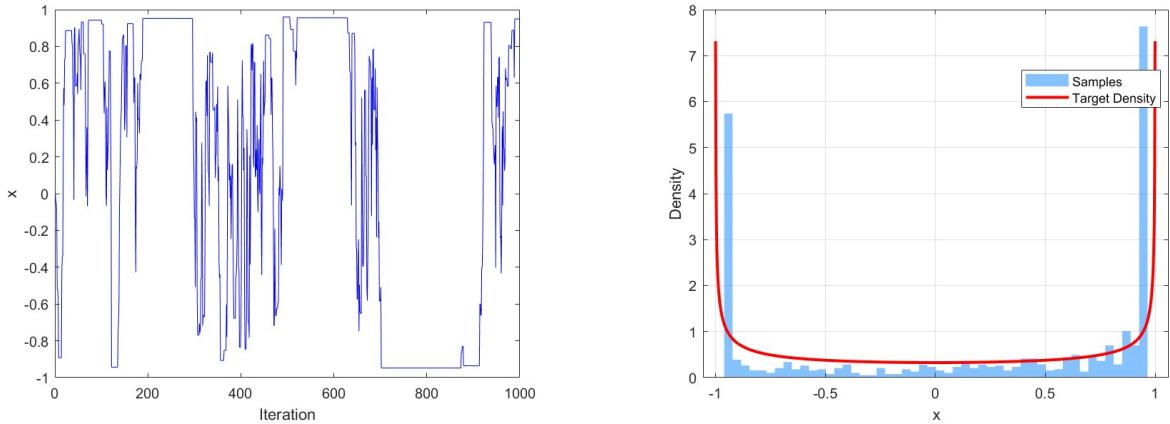


Figure 8: MALA traceplot and histogram on the one-dimensional target with exploding boundary at  $-1$  and  $1$ , i.e.  $\pi(x) \propto (1 - x^2)^{-\frac{1}{2}}$ . Step-size  $h = 0.1$ .

$\{x \in \mathbf{R}^d : -1 \leq x_i \leq 1 \text{ for all } i\}$ . We refer to this as the ‘box-constrained Gaussian’ target.

To illustrate, we run some simulations on a box-constrained Gaussian target in  $d = 10^4$  dimensions with  $\sigma^2 = 1$ , i.e.

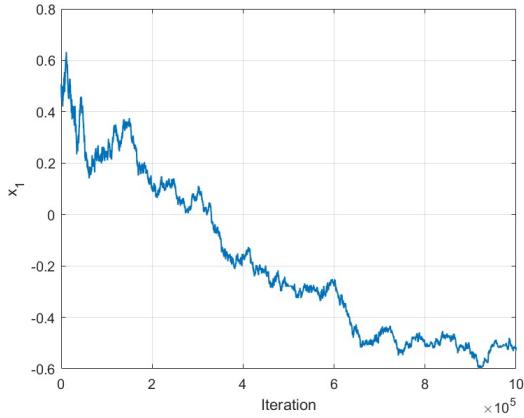
$$\pi(x) \propto \exp\left(-\frac{1}{2}(x_1^2 + x_2^2 + \dots + x_d^2)\right), \quad x \in [-1, 1]^{10^4}. \quad (9)$$

Note this is a significant constraint; due to the high dimensionality of the target, a sample from the *unconstrained* Gaussian distribution will very rarely fall inside this box. We ran the MALA algorithm on this target for  $N = 10^5$  iterations, starting from the point  $(0.5, 0.5, \dots, 0.5)$  and we present our result in Figure 9. In the first row we present the traceplot of the first coordinate and its histogram with a step-size chosen to be  $h = 5 \cdot 10^{-7}$  so that around  $50 - 60\%$  of the proposed jumps are accepted. It is evident from the histogram that the algorithm needs more time to converge and has not fully captured the shape of the target, while from the traceplot, we see that the algorithm moves around very slowly. Furthermore, in the second row, we present the results of an analysis similar to Figure 7. We ran the process for different step-sizes  $h \in \{10^{-7}, 5 \cdot 10^{-7}, 10^{-6}, 5 \cdot 10^{-6}, 10^{-5}, 5 \cdot 10^{-5}\}$ , we ran 100 copies of each process for  $= 10^5$  iterations, and we present the average acceptance probability and the Mean Squared Error (when estimating the target expectation, which in this case is zero) as a function of step-sizes. It is evident from the plots that the Mean Squared Errors are significantly large, given that the range of values each coordinate can take is limited to  $[-1, 1]$ . The acceptance probability heavily depends on the chosen step-size, and the algorithm’s performance rapidly deteriorates at larger step-sizes.

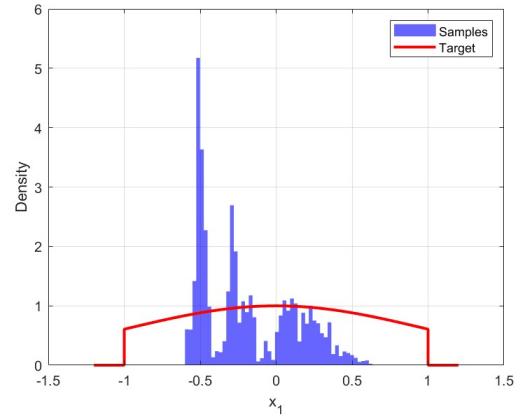
Given the problems that state-of-the-art algorithms may experience in the presence of roughness, we now proceed to describe some resolutions to these pathologies.

### 3.2 Proposed Solutions

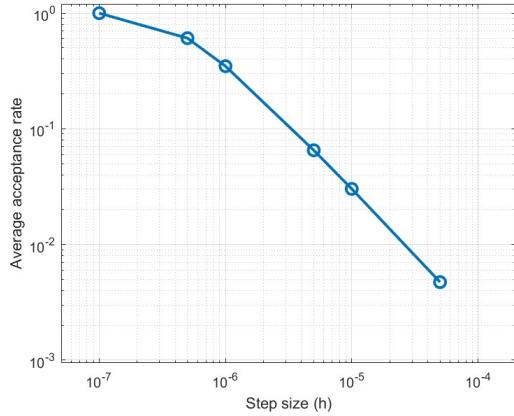
Returning to the original approach of designing MCMC algorithms by taking inspiration from a suitable continuum process, we describe two popular strategies for ‘robustifying’ MCMC in the face of roughness. A first strategy is to focus on discretising the process more carefully, perhaps taking inspiration from numerical-analytic approaches to stiff dynamics. A second strategy is to concoct entirely new continuum processes which admit more stable discretisations ‘by design’, or for which discretisation can somehow be avoided entirely. In



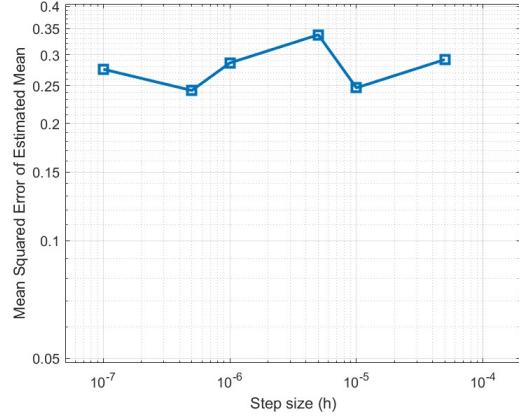
(a) Traceplot of first coordinate. Step-size  $h = 5 \cdot 10^{-7}$ .



(b) Histogram of first coordinate samples. Step-size  $h = 5 \cdot 10^{-7}$ .



(c) Average acceptance probability over step-sizes.



(d) Mean square error of target expectation over step-sizes

Figure 9: MALA on the  $10^4$ -dimensional box-constrained Gaussian target (9), ran for  $N = 10^5$  iterations. First row: traceplot and histogram with optimally tuned step-size  $h = 5 \cdot 10^{-7}$ . Second row: Average acceptance probability and mean squared error (estimating the target expectation) over step-size.

this section, we will review some individual proposals which tackle roughness through some implementation of these principles, focusing more on their conceptual genesis than upon their application to specific examples of practical relevance and scale.

### 3.2.1 Truncated Langevin Monte Carlo

In the setting of ‘steep’ gradients, whereby  $\nabla U$  is perhaps locally-Lipschitz without being globally-Lipschitz, a natural strategy is to perform some sort of truncation. Indeed, in many settings, steepness of gradients is associated with light tails for the invariant measure (it is relatively straightforward to concoct such examples in e.g. the case of convex  $U$ ), and so given  $\delta \in (0, 1)$ , one can often find some rather moderately-sized  $R = R_\delta \in (0, \infty)$  so that

$$x \sim \pi \implies \text{with probability } \geq 1 - \delta, \quad \|\nabla U(x)\| \leq R.$$

As such, in e.g. a Langevin-type algorithm, one could envision replacing all instances of the ‘raw’ drift  $\nabla U(x)$  with the  $R$ -truncated drift  $\tau_R \circ \nabla U$ , where

$$\tau_R(g) := \begin{cases} g & \text{if } \|g\| \leq R \\ R \cdot \frac{g}{\|g\|} & \text{if } \|g\| > R \end{cases}$$

and obtain a process with desirable long-time behaviour (in terms of converging to a sensible invariant measure, even without Metropolis-adjustment) while enjoying good numerical stability properties. In particular, for most ‘reasonable’ targets, the truncated drift will be uniformly Lipschitz, and so ‘standard’ analyses will apply quite directly.

In the context of Monte Carlo-Exact methods, this strategy appears to have first been introduced by [124] as the ‘Metropolis Adjusted Langevin Truncated Algorithm’ (MALTA). [8] later studied this approach in the context of Adaptive MCMC, observing that the intrinsic stability of the truncated algorithm allows for well-behaved adaptation of algorithmic hyperparameters. Additionally, [28] study the numerical-analytic properties of truncated and Metropolis-adjusted procedures when viewed as numerical integrators for SDEs, again observing favourable properties.

We note in passing that such strategies might be characterised more precisely as ‘drift-truncated’, to contrast with other methods which explicitly truncate the state space; see e.g. [108]. Such methods have apparently seen less widespread use in the MCMC world.

In Figure 10, we present a numerical simulation of the Truncated MALA algorithm applied on the one-dimensional log-quartic target from Example 3.1 for various step-sizes. Comparing the trace plots with those of MALA in Figure 6, we see that the algorithmic performance is critically improved for large step-sizes. Of particular note is the step-size  $h = 0.1$  where the Truncated algorithm is exploring the space in an efficient manner, whereas the same step-size for MALA was inducing an algorithm that did not move at all. For smaller step-sizes the performance seems similar to MALA; this is to be expected, as at small step-size, both schemes converge to the Langevin diffusion. Similar improvement in algorithmic performance was observed in simulations for higher-dimensional targets.

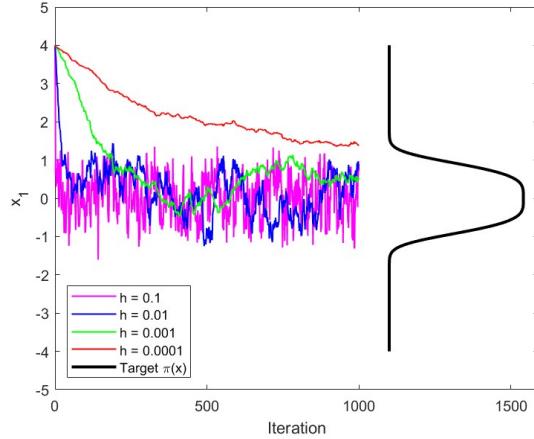


Figure 10: Truncated MALA traceplots with various step-sizes. The target density ( $\pi(x) \propto \exp(-x^4)$ ) is overlaid on the right hand side of the graph.

### 3.2.2 Tamed Langevin Monte Carlo

The ‘tamed’ approach to MCMC has its genesis in the literature on the numerical analysis of stochastic differential equations. A result of [79] established that for a variety of SDEs with non-Lipschitz drift, the

standard Euler-Maruyama scheme is almost-surely explosive, and hence fails to be fit for purpose as a discretisation method. Subsequent work by the same authors [80] introduced the ‘Tamed Euler’ scheme, which ‘tames’ the original drift to be better-behaved, performing a sort of ‘soft truncation’ by analogy with the previous section. We follow here a simplified presentation of a perspective on taming given in [131], focusing for simplicity on the case of SDEs with additive noise.

In general, given the SDE

$$dX_t = b(X_t) dt + \sigma dW_t,$$

‘taming’ is accomplished by introducing some step-size-dependent function  $G^h(x)$  which vanishes as  $h \rightarrow 0^+$  (e.g.  $G^h(x) = h^\alpha \cdot \|b(x)\|$  for some  $\alpha > 0$ ), defining

$$b^h(x) = (1 + G^h(x))^{-1} \cdot b(x),$$

and solving the ‘tamed’ SDE

$$dX_t = b^h(X_t) dt + \sigma dW_t, \quad (10)$$

with a conventional Euler-Maruyama scheme of step-size  $h$ . For a well-chosen  $G^h$ , the drift  $b^h$  can have much milder growth at infinity than the original  $b$ , and so the stability properties of this composite method can be obtained more directly. Moreover, since  $G^h(x)$  vanishes as  $h \rightarrow 0^+$ , for small  $h$ , the two SDEs should follow one another reasonably closely. Indeed, [80] establish strong convergence guarantees for the ‘Tamed Euler’ scheme in some generality, allowing for results to be obtained in settings for which the basic Euler scheme is known to fail, and for which implicit schemes can be rather costly to implement.

In the context of Monte Carlo simulation, the approach of taming has largely been used to stabilise numerical approximations to the Overdamped Langevin Diffusion (2). [33] introduced the ‘Tamed ULA’ approach (including a variant with ‘coordinate-wise’ taming, natural for high-dimensional applications with meaningful coordinate structure), proving some theoretical guarantees, and performing some exploration of a Metropolis-adjusted variant. Various follow-up works ([99, 101, 100]) have continued to complete the theoretical picture for these approaches, focusing predominantly on the unadjusted setting.

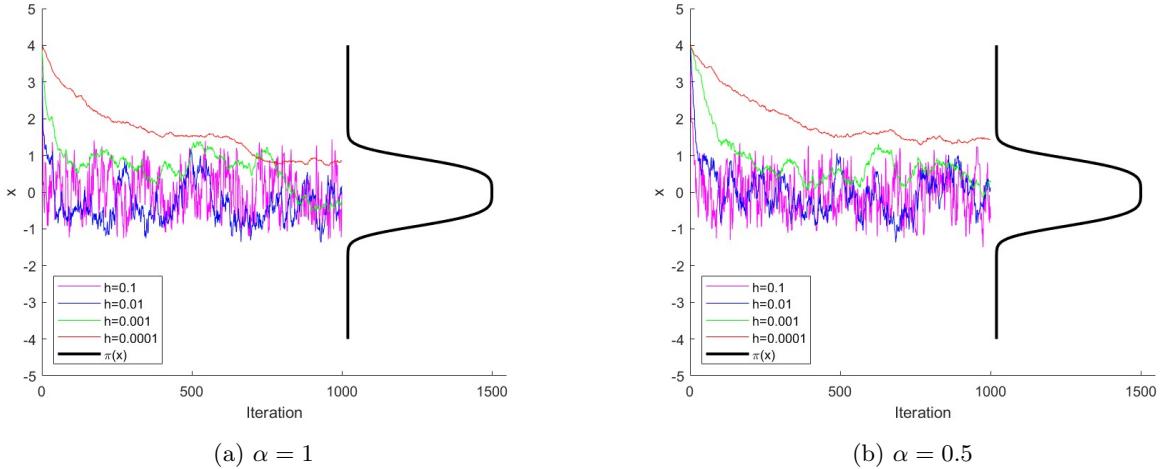


Figure 11: Tamed Langevin Monte Carlo with various step-sizes. The target density  $\pi(x) \propto \exp(-x^4)$  is overlaid on the right-hand side of each graph.

In Figure 11 we present a numerical simulation of the Tamed Metropolis Adjusted Langevin algorithm. We consider again the one-dimensional log-quartic target  $\pi(x) \propto \exp(-x^4)$ , and use the Tamed scheme that

aims to discretise the Langevin equation

$$\begin{aligned} dX_t &= \nabla \log \pi(X_t) dt + \sqrt{2} dW_t \\ &= -X_t^3 dt + \sqrt{2} dW_t \end{aligned}$$

using an update of the form (10), with

$$G^h(x) = h^\alpha \cdot \|\nabla \log \pi(x)\| \quad (11)$$

$$= h^\alpha \cdot \|x\|^3 \quad (12)$$

The left plot considers the case where  $\alpha = 1$ , while the right one has  $\alpha = 0.5$ . Both plots consider different step-sizes, just as in the previous section. Comparing the trace plots with those of MALA in Figure 6, we once again see that the algorithmic performance is critically improved for large step-sizes. The algorithmic performance seems similar to the Truncated algorithm, discussed in Section 3.2.1.

Finally, to further showcase the differences between algorithms, in Figure 12 we consider the two-dimensional target of the form

$$\pi(x_1, x_2) \propto \exp(-x_1^4 - x_2^4 - 5 \cdot x_1 \cdot x_2) \quad (13)$$

and for each of the algorithms MALA, Truncated MALA and Tamed MALA, we plot the vector field that indicates the expected proposed jump of the algorithm from each current position. The colour indicates the log density of the target at the current point and the magnitude of the arrow indicates how large the expected jump will be. By inspecting the figure, we see that away from the mode, the drift which is used in MALA is rather large in magnitude, which runs the risk of inducing numerical instability. On the other hand, for the other two algorithms, the drift tends to stay uniform in magnitude, while it still guides the process towards the mode, leading to more stable numerical behaviour.

### 3.2.3 Proximal Langevin Monte Carlo

A benefit of the taming approach as a numerical method is that it induces stability while retaining an explicit implementation. From a numerical-analytic perspective, another conventional approach to stability is to instead take an *implicit* approach. Such approaches often observe similar stability benefits and improved accuracy, albeit at the cost of a more involved implementation.

As applied to the Langevin diffusion, this approach has been implemented in a couple of different ways. The original work of [116] starts from the perspective of sampling from a target whose potential  $U$  is convex but non-smooth. Potentials of this form are historically common in optimisation-based approaches to image processing, and gained some popularity in sampling-based Bayesian approaches to the same task. With this non-smoothness in mind, the author proposes to mollify the target potential  $U$  according to the ‘Moreau-Yosida convolution’, i.e. constructing the ‘Moreau-Yosida envelope’ as

$$U^\lambda(x) := \inf \left\{ U(y) + \frac{\|x - y\|^2}{2 \cdot \lambda} : y \in \mathbf{R}^d \right\}$$

to obtain an approximate target with improved regularity properties. In particular (see Theorem 4.1.4 and Proposition 4.1.5 of [76]), assuming that  $U$  is convex with a closed graph, then

$$U^\lambda(x) \xrightarrow{\lambda \rightarrow 0} U(x)$$

point-wise, while the mollified potential  $U^\lambda$  is automatically differentiable, and its gradient is automatically Lipschitz-continuous, with Lipschitz constant upper-bounded by  $\lambda^{-1}$ . Moreover, provided that the original  $U$  is convex, so too will  $U^\lambda$  be. Combining these properties, one sees that usual numerical discretisations of the gradient flow with respect to  $U^\lambda$  will behave stably, even in the presence of stochastic noise and other perturbations. Bearing all of this in mind,  $U^\lambda$  presents itself favourably as a candidate for more naive

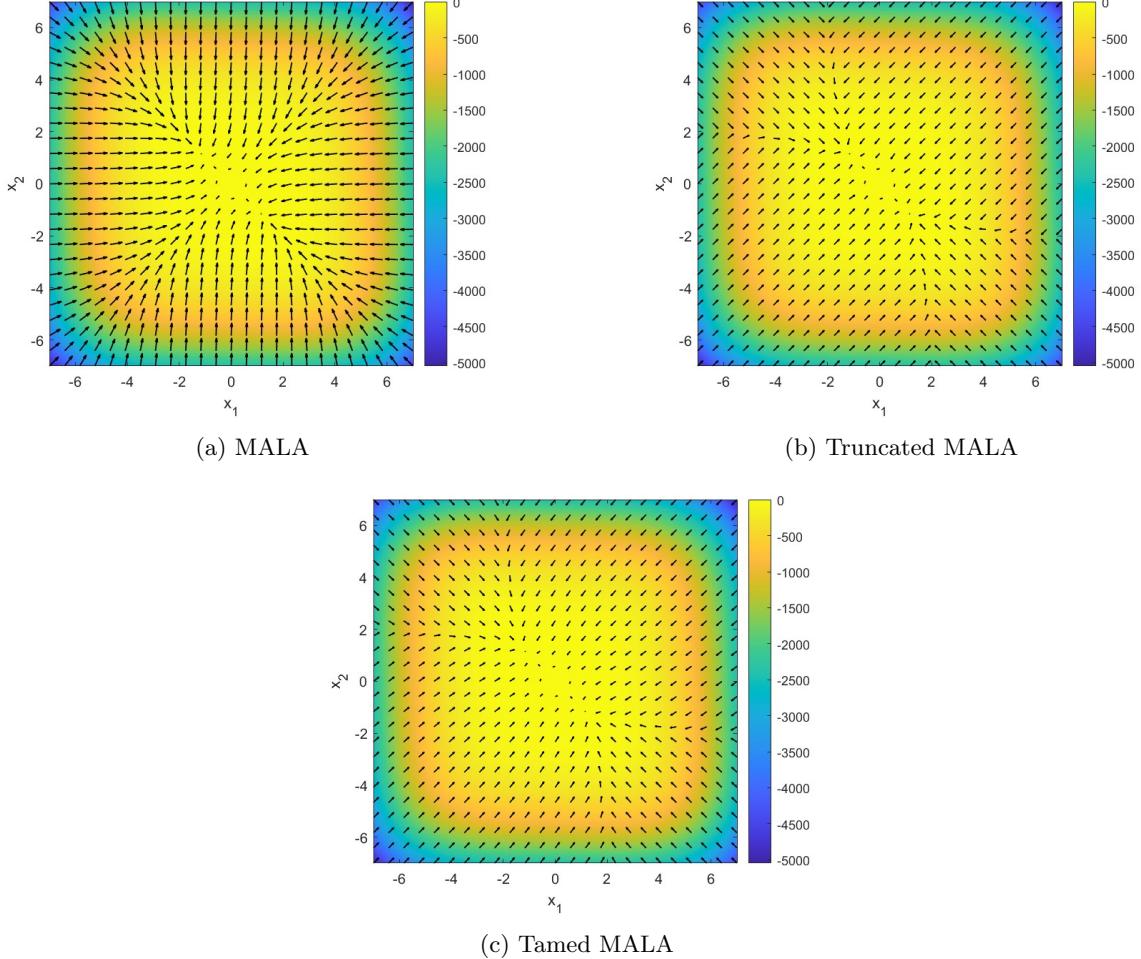


Figure 12: Vector fields of the expected proposed jump from the current position for MALA, Truncated MALA ( $R = 10$ ) and Tamed MALA. Target  $\pi$  as in (13). Step-size  $h = 0.1$ . Colours indicate the level sets of log-density.

numerical discretisation, whereby explosivity and instability of the process are essentially ruled out by design. For a graphical representation of the Moreau-Yosida envelope, see Figure 13, where it is applied with  $\lambda = 1$  on the (non-differentiable at zero) Laplace potential.

Following in this direction, using Euler discretisation on the Overdamped Langevin Diffusion with target measure proportional to  $\exp\{-U^\lambda(x)\}$ , one obtains the ‘Proximal ULA’ (P-ULA) proposal as

$$X_n = X_{n-1} - h \cdot \nabla U^\lambda(X_{n-1}) + \sqrt{2 \cdot h} \cdot \xi_n, \quad \xi_n \sim \mathcal{N}(0, \mathbf{I}_d).$$

This framing highlights an interesting point of comparison with some of the previous solutions: while the Truncated and Tamed Langevin strategies can be viewed as employing a modified numerical scheme to sample from the original target distribution (up to discretisation error), the Proximal Langevin approach instead fixes the numerical scheme, and systematically changes the target distribution. We will see in later discussion that this conceptual distinction allows for some specific functionalities which are not available to all approaches.

A priori, one could expect that computing  $U^\lambda$  and  $\nabla U^\lambda$  at each step might add substantial complexity to the algorithm. However, in various settings, the additional burden is relatively manageable. In particular, for small enough  $\lambda$ , the optimisation formulation which defines  $U^\lambda$  is a strongly-convex minimisation problem,

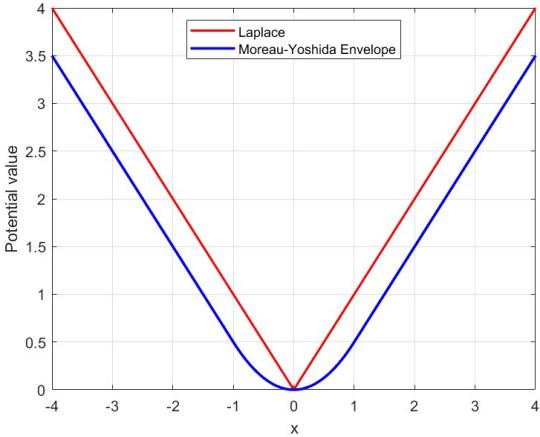


Figure 13: Laplace potential and its associated Moreau-Yosida envelope with  $\lambda = 1$

which can be rapidly solved using the methods of modern convex optimisation (see e.g. [31, 34]). Moreover, the gradient  $\nabla U^\lambda$  can be obtained by examining the value of  $y$  which solves this minimisation problem, that is

$$\begin{aligned} \nabla U^\lambda(x) &= \lambda^{-1} \cdot (x - \text{prox}_{\lambda \cdot U}(x)) \\ \text{prox}_{\lambda \cdot U}(x) &:= \arg \min \left\{ U(y) + \frac{\|x - y\|^2}{2 \cdot \lambda} : y \in \mathbf{R}^d \right\}, \end{aligned}$$

where  $\text{prox}_{\lambda \cdot U}$  is the so-called ‘proximal operator’ associated to  $U$ . As such, both  $U^\lambda$  and its gradient can be obtained ‘in a single pass’, so to speak. Observe also that upon taking  $\lambda = h$ , the P-ULA proposal simplifies to

$$X_n = \text{prox}_{h \cdot U}(X_{n-1}) + \sqrt{2 \cdot h} \cdot \xi_n, \quad \xi_n \sim \mathcal{N}(0, \mathbf{I}_d).$$

As an aside for readers unfamiliar with convex analysis, it may be useful for one’s intuition to verify that when  $U$  corresponds to the characteristic function of a closed, convex set  $\mathcal{K}$  (i.e. 0 for  $x \in \mathcal{K}$  and infinite otherwise), the Moreau envelope takes the form  $U^\lambda(x) = \frac{1}{2\lambda} \cdot \text{dist}^2(x, \mathcal{K})$ , and the proximal operator is the projection onto  $\mathcal{K}$ . In this sense, the Moreau-Yosida smoothing can be viewed as a way of generalising these notions of distance and projection from sets to arbitrary convex functions.

Having established these preliminaries, [116] identifies a range of practical problems in high-dimensional Bayesian image analysis for which basic Langevin algorithms perform poorly, the new proposals can be computed efficiently, and the resulting proximal Langevin algorithms – whether unadjusted (P-ULA) or Metropolis-adjusted (P-MALA) – perform favourably.

To showcase the performance of P-MALA, we return to the example of a  $10^4$ -dimensional product-of-Laplace target distribution (7), following Example 3.2. We report our results in Figure 14. On the first row we report the traceplot and histogram of the first coordinate of a P-MALA algorithm, ran for  $N = 10^5$ , using step-size  $h = 0.01$ . The step-size was chosen so that the average acceptance of proposed step is close to 36% (see e.g. [42]). The chain was initialised at  $(5, 5, \dots, 5)$ , a significant distance from the mode. On the second row, we present a similar analysis to Figure 7, testing the behaviour of the algorithm over different step-sizes and reporting the average (over 100 independent runs) acceptance probability and Mean Squared Error (MSE) when estimating the target expectation (in this case zero). While the variations in acceptance probability and MSE seem similar to Figure 7, it can be observed that for all step-sizes, the average acceptance probability for P-MALA is larger than it is for MALA, while the MSE is always lower (albeit still quite large) for most step-sizes, indicating better algorithmic performance. Comparing with Figure 7, it seems that the PMALA has captured the shape of the distribution in a more efficient manner.

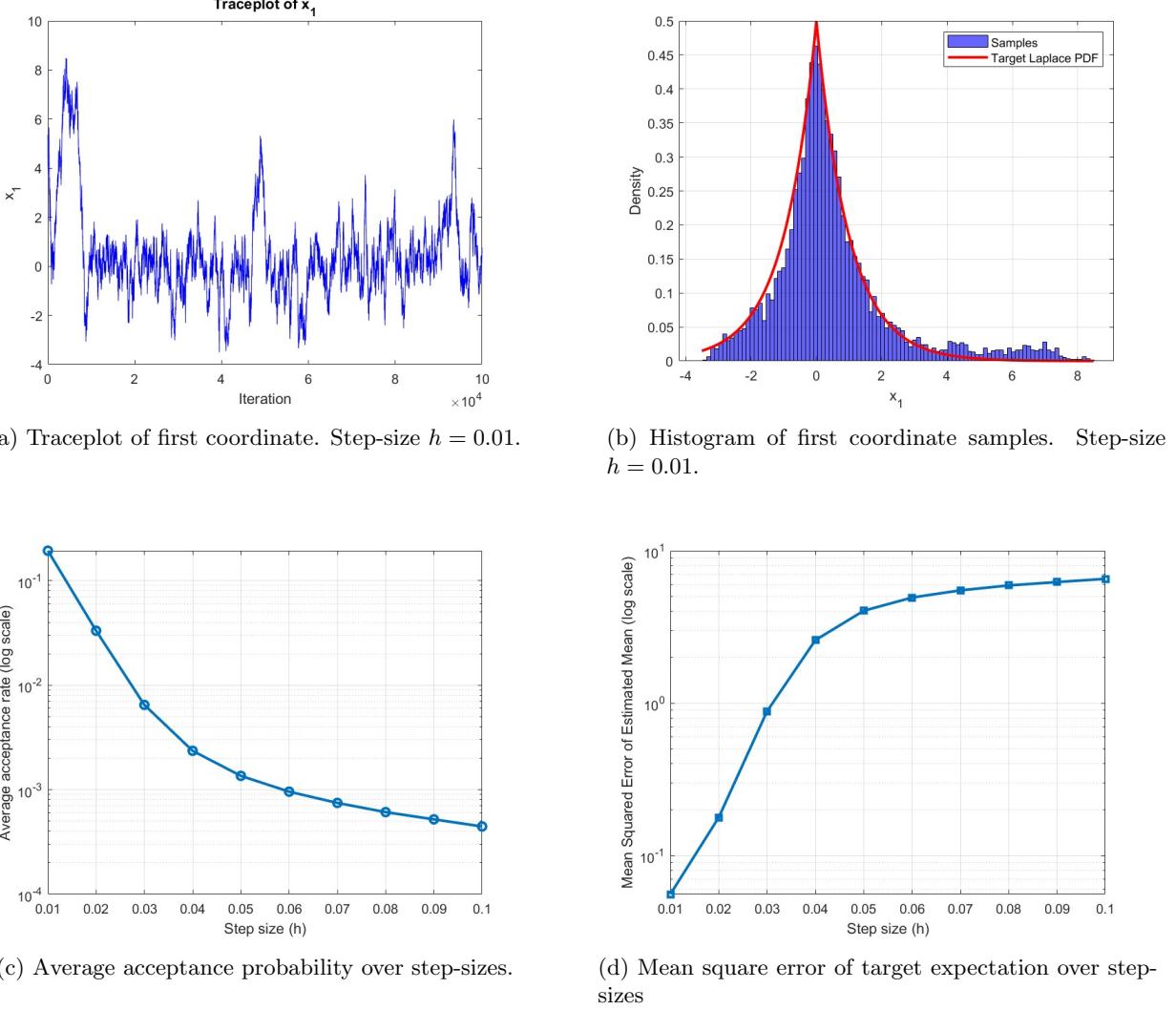


Figure 14: Proximal MALA on the  $10^4$ -dimensional Laplace target (7). First row: traceplot and histogram with optimally tuned step-size  $h = 0.01$ . Second row: Average acceptance probability and Mean squared error (estimating the target expectation) over step-size.

As a brief aside, we note that while the Moreau-Yosida convolution always delivers a potential whose gradient is Lipschitz, for small regularisation parameters  $\lambda$ , this Lipschitz constant can also be small, leading to residual stiffness in the dynamics of the Langevin diffusion. In this setting, it has proven fruitful [117] to make use of ‘explicit, stabilised’ numerical integrators which are able to resolve this stiffness while retaining a large nominal step-size, at an increased but worthwhile per-iteration computational cost.

In terms of other algorithms that are based on the idea of Moreau-Yosida convolution, [137] studies a related algorithm which transposes the order of computing  $\text{prox}_{\lambda \cdot U}$  and injecting additive Gaussian noise, and [12] study a variant which accommodates the possibility of computing the proximal operator inexactly, each obtaining convergence results along the way.

Subsequent work of [77] revisits these ideas from a more numerical-analytic perspective, instead framing the P-ULA proposal as a (non-standard) implicit Euler-Maruyama discretisation of the Langevin diffusion. Building on this perspective, the authors propose to generate moves by the ‘theta method’, whereby one

fixes a  $\vartheta \in [0, 1]$  and transitions from  $X_{n-1}$  to  $X_n$  by solving the nonlinear equation

$$X_n = X_{n-1} - h \cdot \{\vartheta \cdot \nabla U(X_{n-1}) + \bar{\vartheta} \cdot \nabla U(X_n)\} + \sqrt{2 \cdot h} \cdot \xi_n, \quad \xi_n \sim \mathcal{N}(0, \mathbf{I}_d),$$

where  $\vartheta + \bar{\vartheta} = 1$ . The resulting algorithm is then termed the ‘Implicit Langevin Algorithm’ (ILA). For  $\vartheta = 0$ , this is the conventional explicit ULA proposal, which can be expected to suffer from stability issues as detailed earlier. For  $\vartheta > 0$ , this is now an *implicit* proposal, which requires the solution of an auxiliary optimisation problem, but is expected to perform favourably with a view to numerical and dynamical stability. While this formulation is superficially different from that of P-ULA, some rearrangement shows that  $X_n$  can be obtained from  $X_{n-1}$  as the solution to

$$X_n = \arg \min \left\{ U(y) + \frac{\|y - \bar{X}_{n,\bar{\vartheta}}\|^2}{2 \cdot \vartheta \cdot h} : y \in \mathbf{R}^d \right\}$$

$$\text{where } \bar{X}_{n,\bar{\vartheta}} = X_{n-1} - h \cdot \bar{\vartheta} \cdot \nabla U(X_{n-1}) + \sqrt{2 \cdot h} \cdot \xi_n,$$

i.e. by minimising a similar quadratic penalisation of the potential  $U$ . As such, the implementation complexity of ILA should be comparable to that of P-ULA. Note that ILA is superficially closer to the method proposed in [137] than to P-ULA, in that the proximal operator is applied *after* adding noise to the current iterate.

One might expect this methodology to perform reasonably in similar situations to P-ULA, although the work [77] does not really examine the non-smooth case empirically, focusing instead on taking larger step-sizes in the setting of strongly-log-concave targets with Lipschitz-continuous forces, albeit with a large condition number.

An interesting hybrid approach was presented recently in [128]. This approach proposes to run MCMC which genuinely targets the Moreau-smoothing  $\pi^\lambda$  of the target distribution  $\pi$ , and then use these samples as the basis for an importance sampling correction. In practice, this amounts to using the *proposal* kernel corresponding to P-ULA (or similar), but performing the Metropolis-adjustment with respect to  $\pi^\lambda$  rather than  $\pi$ . One thus benefits from the good stability properties of the MCMC kernels, allowing for a relatively friendly implementation. Moreover, while the invariant measure of the implemented MCMC kernel is not the target distribution  $\pi$ , for a well-chosen  $\lambda$ , one expects to sample from a rather good approximation to  $\pi$ , and so the importance sampling adjustment should be rather gentle, leading to only a mild inflation of variance. Observe that implementing the same conceptual program with the Truncated or Tamed Langevin approach would not be entirely straightforward, as the invariant measure of the unadjusted methods is both different from  $\pi$  and difficult to characterise explicitly.

### 3.2.4 Barker-Langevin Monte Carlo

The preceding methods largely have their genesis in viewing existing MCMC algorithms as discretisations of diffusion processes. The so-called ‘Barker proposal’, as introduced in [95], originates from a different perspective of ‘local balancing’, initially put forward as a general principle in [141]. In intuitive terms, the local balancing framework takes as given some target-agnostic ‘local exploration kernel’  $Q(x, dy)$ , and seeks to correct it gently so as to obtain an improved proposal kernel. In particular, the authors considered ‘tilted’ kernels of the form

$$\bar{Q}(x, dy) \propto Q(x, dy) \cdot \beta \left( \frac{\pi(y)}{\pi(x)} \right)$$

for some ‘balancing function’  $\beta : \mathbf{R}_+ \rightarrow \mathbf{R}_+$ . By taking this function to satisfy the symmetry condition

$$\beta(r) = r \cdot \beta(r^{-1}),$$

and assuming that  $Q$  is symmetric (in that the law  $Q(x, \cdot)$  has a density which is invariant under swapping its arguments), one obtains a kernel which is approximately  $\pi$ -reversible, and is hence expected to interface well with Metropolis-adjustment. In general, simulation from the kernel  $\bar{Q}$  is not directly tractable, and so some additional approximations are required in order to obtain a practical kernel.

**Example 3.5** (Connections with ULA). *Focus on the case in which  $Q_\sigma$  is a simple Gaussian random walk with proposal standard deviation equal to  $\sigma$ , and consider  $\beta(r) = r^{\frac{1}{2}}$ . One can then compute that the resulting  $\bar{Q}$  takes the form*

$$\begin{aligned}\bar{Q}_\sigma(x, dy) &\propto Q_\sigma(x, dy) \cdot \sqrt{\frac{\pi(y)}{\pi(x)}} \\ &= Q_\sigma(x, dy) \cdot \exp\left\{\frac{1}{2}(\log \pi(y) - \log \pi(x))\right\} \\ &\approx Q_\sigma(x, dy) \cdot \exp\left(\frac{1}{2}\langle \nabla \log \pi(x), y - x \rangle\right).\end{aligned}$$

*It can be verified that this final law coincides with the ULA kernel. As per our earlier discussion, while this method can perform well for regular targets, one anticipates various instabilities in the face of more serious steepness.*

An insight of [95] is that one can consider the same construction, but replacing the square root with the so-called ‘Barker’ balancing function  $\beta(r) = \frac{r}{1+r}$  (named due to a connection with an early work of [9]). Similarly to Example 3.5, a first order Taylor approximation of the form  $\log \pi(y) - \log \pi(x) \approx \langle \nabla \log \pi(x), y - x \rangle$  yields

$$\begin{aligned}\bar{Q}_\sigma(x, dy) &\propto Q_\sigma(x, dy) \cdot \frac{\pi(y)}{\pi(x) + \pi(y)} \\ &\approx Q_\sigma(x, dy) \cdot s(\langle \nabla \log \pi(x), y - x \rangle),\end{aligned}$$

where  $s : t \mapsto (1 + \exp(-t))^{-1}$  is a standard sigmoid function, which also coincides with the distribution function of the standard Logistic distribution. The resulting kernel

$$\hat{Q}(x, dy) = 2 \cdot \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}|y - x|^2\right) \cdot s(\langle \nabla \log \pi(x), y - x \rangle) dy \quad (14)$$

is then viable as a proposal kernel which can be corrected via a Metropolis acceptance step. Interestingly, sampling from the proposal  $\hat{Q}_\sigma$  has a simple algorithmic interpretation: first, ‘pre-propose’ a uninformed move by generating some standard Gaussian noise  $\xi$ , and by looking at the alignment of  $\xi$  with  $\nabla \log \pi(x)$ , then decide whether it will be more fruitful to propose a ‘real’ move to  $y_+ = x + \sigma \cdot \xi$ , or to its ‘reflection’,  $y_- = x - \sigma \cdot \xi$ . In particular, given  $\xi$ , one samples  $y_+$  with probability  $p(x, \xi) = s(\langle \nabla \log \pi(x), \sigma \cdot \xi \rangle)$ , and  $y_-$  with the complementary probability  $1 - p(x, \xi)$ .

A key property of this proposal is that while it is gradient-informed (and hence makes active use of the local geometry of the target distribution), it also enjoys a distinctive ‘bounded influence’ property: even in the face of massive gradients, one does not make massive moves, as the original size of the jump  $\sigma \cdot \xi$  comes from a Gaussian distribution. Instead, one simply becomes overwhelmingly likely to move in a specific direction. In this regard, the proposal circumvents one of the challenging effects of steep gradients.

In [95], however, it was shown that the spectral gap of Metropolised  $\hat{Q}_\sigma$  is no better than twice the spectral gap of the Metropolised Gaussian Random Walk  $Q_\sigma$ , which indicates that the computational gains obtained by choosing such a method are rather limited. An important practical variation of the algorithm was thus suggested. Assuming that the target lives in  $\mathbf{R}^d$ , instead of choosing whether to add or subtract the noise vector  $\xi$  ‘all at once’, one can instead look at each of the  $d$  coordinates of the noise *separately*, and consider all  $2^d$  possible combinations in which this noise can be added to the current state  $x$ . More specifically, one first samples a Gaussian random variable  $\xi \sim \mathcal{N}(0, \mathbf{I}_d)$ . For any coordinate  $i \in \{1, \dots, d\}$ , one then samples a random variable  $b_i$  such that

$$b_i = \begin{cases} 1 & \text{with probability } p_i(x, \xi), \\ -1 & \text{with probability } 1 - p_i(x, \xi), \end{cases}$$

where

$$p_i(x, \xi) = s \left( \frac{\partial}{\partial x_i} \log \pi(x) \cdot \sigma \cdot \xi_i \right)$$

and for  $b = (b_1, \dots, b_d) \in \{\pm 1\}^d$ , one proposes to move to  $y_b = x + \sigma \cdot b \odot \xi$ , where  $\odot$  denotes the element-wise product. In other words, the algorithm decides whether to accept or reflect in the direction of each coordinate *separately*, conditionally on  $x$  and  $\xi$ . This modified procedure has become known as the “(coordinate-wise) Barker proposal”. It preserves many of the desirable stability properties of the original  $\hat{Q}$ , but allows for extra flexibility and adaptivity to the information provided by the gradient, leading to improved sampling efficiency, particularly in high dimension. While the proposal is not exact in terms of exactly preserving  $\pi$ , it can similarly be corrected by Metropolis-adjustment without losing most of its computational efficiency; see [136] for a scaling analysis in support of this observation.

Working backwards from this construction, recent work of [94] proposes to use the same ‘propose-reflect’ mechanism as a tool for numerical discretisation of general stochastic differential equations, naturally focusing on applications with similarly ‘steep’ drifts. Interestingly, one can interpret the Barker proposal as a special case of this SDE discretisation, when applied to the Overdamped Langevin Diffusion (2), which further motivates the use of the Barker proposal. It also presents a number of opportunities for future developments in which Metropolis-adjustment of the dynamics is not a key priority, as is often the case for e.g. stochastic gradient methods; see e.g. [104].

In order to showcase the performance of the Barker proposal, we consider two numerical examples. The first target is our running univariate log-quartic model problem. In Figure 15, we present the traceplots of the algorithm run for  $N = 1000$  iterations, starting from  $x_0 = 4$  and for four different step-sizes ( $h = 0.1, 0.01, 0.001, 0.0001$ ). Contrary to MALA, Barker shows a remarkable robustness in the behaviour across the different step-sizes. Interestingly, for this particular example, it seems that for small step-sizes, the method is able to find the mode of the distribution faster than for the previously mentioned algorithms. We conjecture that this is due to the reflection mechanism, which allows for better-informed proposal jumps, and hence to fewer rejected moves.

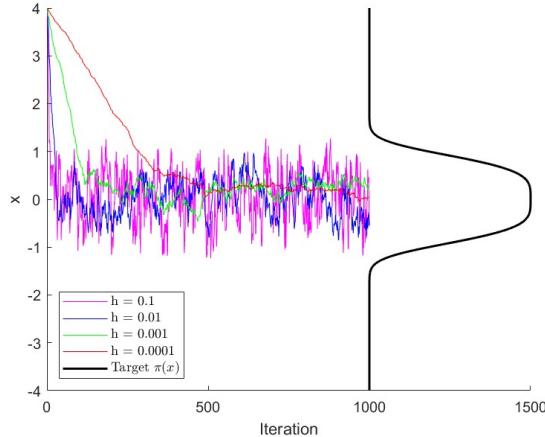


Figure 15: Barker trace plots with various step-sizes. The target density ( $\pi(x) \propto \exp(-x^4)$ ) is overlaid on the right hand side of the graph.

The second target is of the form (8), supported on  $(-1, 1)$ , and exploding at the boundary. The algorithm is run for  $N = 1000$  iterations, with step-size  $h = 0.1$ , starting from 0.5. We present the traceplot and histogram of the MALA algorithm in Figure 16. Compared to Figure 8, the Barker traceplot seems to get stuck less often at the boundary, while the histogram looks much more symmetric and tends to capture the shape of the target more efficiently. We should emphasise here, that over independent realisations of the algorithm, the Barker chains showcased much more consistent behaviour when compared to MALA,

which has frequent realisations that completely failed to explore the area away from the boundaries, and would presumably require drastically more iterations in order to converge reasonably.

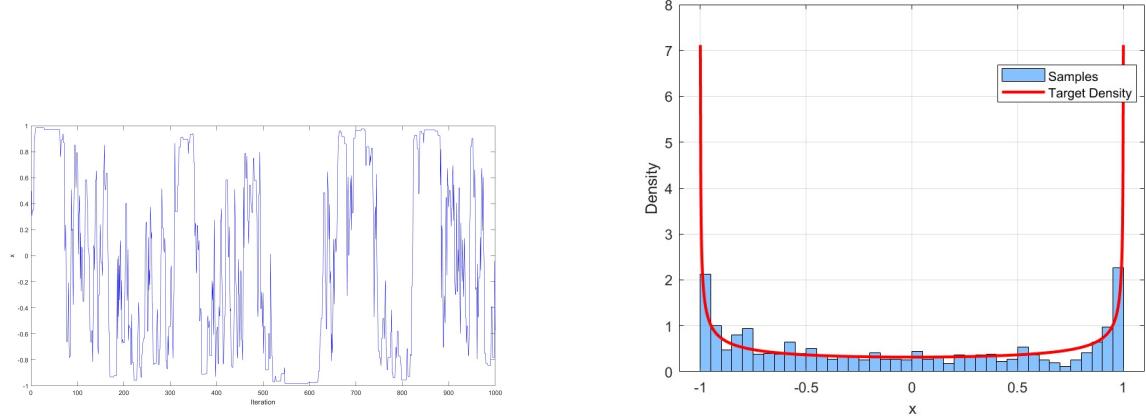


Figure 16: Barker traceplot and histogram on the one-dimensional target with exploding boundary at  $-1$  and  $1$  ( $\pi(x) \propto (1 - x^2)^{-\frac{1}{2}}$ ). Step-size  $h = 0.1$ .

In order to give a visual representation of the Barker proposal, in Figure 17 we plot the densities of the Barker, MALA and Random Walk proposal (Gaussian centered around the current position), when targeting the log-quartic target  $\pi(x) \propto \exp(-x^4)$ . We use step-size  $h = 0.5$  and the current position for all three proposals is  $x = -2$ . It is evident that the Gaussian Random walk does not capture the shape of the target well, assigning most of its mass close to the current point  $x = -2$ . MALA on the other hand, failed to an “overshooting” phenomenon, jumping massively to the other side of the tails of the target (centered around 6). The Barker seems to be the middle ground, centering its mass closer towards zero (the target mean).

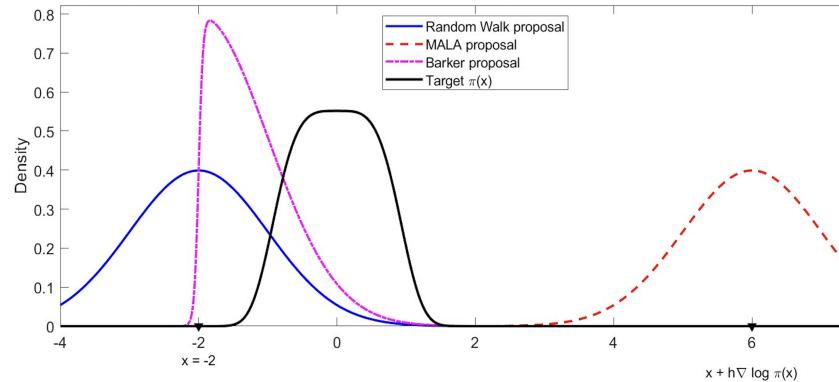


Figure 17: Density plots of Barker, MALA and Random Walk proposals. Target  $\pi(x)$  is the log-quartic. Step-size  $h = 0.5$ . Current position is  $x = -2$ .

### 3.2.5 Hamiltonian Monte Carlo with Non-Quadratic Kinetic Energy

Beyond Langevin-type MCMC methods, a particularly popular class of modern gradient-based MCMC algorithms are developed around *Hamiltonian Monte Carlo* (HMC), a broad term which we take to encompass various implementable discretisations of the aforementioned Refreshed Hamiltonian Dynamics process of [27]. Despite various conceptual differences to Langevin Monte Carlo (which include moving from stochastic dynamics to deterministic dynamics, from ‘one-step’ proposal mechanisms to ‘trajectorial’ proposals, and so on), in practice, Hamiltonian Monte Carlo suffers from the same qualitative shortcomings in the face of

steep gradients; as soon as  $\nabla \log \pi$  is of superlinear growth at infinity, the numerical stability of conventional numerical integrators can collapse. By analogy with Proximal Langevin algorithms, [36] proposed to run HMC on a Moreau-Yosida-smoothed version of the target distribution; this appears not to have been taken up much in practice, though the recent appearance of the preprint [129] suggests that the door has not been closed entirely.

An interesting alternative approach was proposed by [93], who leave the potential  $U$  unchanged, and instead seek to modify the *kinetic* energy associated to the underlying (fictitious) Hamiltonian system. More precisely, where the usual HMC samples from a joint target distribution  $\mu$  (with states given by  $z = (x, v)$ ) which satisfies

$$-\log \mu(z) = U(x) + \frac{1}{2} \|v\|^2 + \text{const.},$$

one can instead seek to modify the marginal distribution of the momentum<sup>2</sup>  $v$  and sample according to

$$-\log \mu(z) = U(x) + K(v) + \text{const.},$$

for some other well-chosen function  $K : \mathbf{R}^d \rightarrow \mathbf{R}$ . In this setting, the Hamiltonian dynamics instead take the form

$$\dot{x} = \nabla K(v), \quad \dot{v} = -\nabla U(x),$$

which exposes that steepness in  $\nabla U$  might be alleviated by suitable combination with  $\nabla K$ . Indeed, the authors argue in favour of designing  $K$  such that  $\nabla K \circ \nabla U$  is asymptotically of linear growth as  $\|x\| \rightarrow \infty$ . This facilitates the numerical stability of ‘usual’ discretisations of Hamiltonian dynamics, while also ensuring that the state  $x$  ultimately has a linear drift back towards the ‘bulk’ of the state space, i.e. that the dynamics are not flattened to the point of inefficiency. They then advocate for a ‘relativistic’ kinetic energy choice (following [98]) of the form

$$K(v) = \left(1 + \|v\|^2\right)^{a/2} - 1, \quad a \geq 1$$

(modulo various multiplicative scaling factors), which has the desired quadratic-like behaviour for small  $\|v\|$ , but grows like  $\|v\|^a$  for large  $\|v\|$ , ensuring a gentle perturbation to conventional Hamiltonian dynamics in the bulk of the space, with a flexible regularising influence out in the tails. In Figure 18, we present the contour plots of the Hamiltonian dynamics for the one-dimensional log-quartic target (that is with potential  $U(x) = x^4$ ) and two choices of kinetic energy. The horizontal axis represents the  $x$ -space, with the vertical axis as the velocity space. These plots show the curves along which the Hamiltonian energy  $-\log \mu(x, v)$  remains constant, along which the dynamics move. In particular, if we start from  $(x, v) = (1, 0)$ , the path of the dynamics is indicated with bold.

A specific construction with some fascinating properties involves taking  $K(v) = \|v\|_1$  in the above construction, corresponding to momenta following the multivariate Laplace distribution. In this setting, it holds that  $\nabla K(v) = \text{sign}(v)$ , so that (ignoring singularities at which the coordinates of  $v$  change sign) one always has  $\dot{x} \in \{\pm 1\}^d$ , i.e. the velocity of the process lives on the discrete hypercube, and hence the  $x$  component has a uniformly bounded (in fact, constant) speed, enabling stable discretisation in some generality.

This peculiarity of the velocities also hints to an intriguing connection with Piecewise-Deterministic Markov Processes, and in particular, to the Zig-Zag Process of [22]. Indeed, work of [112] explores this connection in greater depth, witnessing this ‘Hamiltonian Monte Carlo with Laplace-Distributed Momentum’ as a ‘Non-Markovian Zig-Zag Process’. In this analogy, the times at which some momentum coordinate  $v_i$  crosses 0 correspond qualitatively to the ‘velocity coordinate flip’ events of the Zig-Zag Process. While the system follows deterministic dynamics, randomness is induced via random refreshments of the momentum  $v$ , ensuring that appropriate stochasticity enters the system. Empirically, the authors observe that in the Hamiltonian

---

<sup>2</sup>We hope that the expert reader will forgive us for sticking to ‘velocity-type’ notation here, in the interests of cohesion with our presentation of PDMPs.

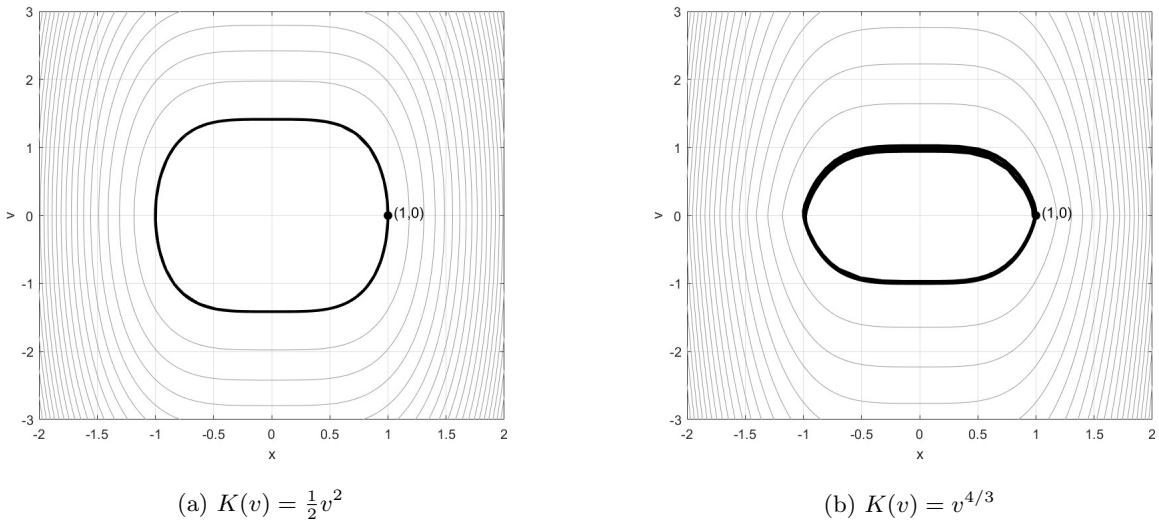


Figure 18: Contour plots of the Hamiltonian dynamics with the quartic potential  $U(x) = x^4$  and different kinetic energies.  $x$ - axis is the  $x$ -space,  $y$ -axis is the velocity space. With bold, the path of the process starting from  $x = 1, v = 0$ .

formulation, the process is able to travel further in space for an equivalent number of ‘events’, following the intuition that with these ‘more deterministic’ dynamics, one can suppress unnecessary stochasticity and promote efficient spatial exploration.

We conclude with a numerical study of the light tailed one-dimensional log-quartic target of the form  $\pi(x) \propto \exp(-x^4)$ , introduced in Example 3.1. We consider the Hamiltonian dynamics with kinetic energy of the form  $K(v) = |v|^{4/3}$ , in accordance with the suggestion of [93]. We use the Leapfrog algorithm to discretise the Hamiltonian dynamics, using step-size  $h$  and  $L$  leapfrog steps per move. This leads to discretisation errors which we correct via the usual Metropolis approach, ensuring that the chain targets the correct distribution. We ran the algorithm for  $N = 1000$  iterations and starting position  $x_0 = 5$ . Our results are summarised in Figure 19, where we present the traceplots of four algorithms with various choices of step-size and leapfrog steps. On the legend, the various choices of  $h$  and  $L$  are described, along with the average acceptance probability of the proposed jumps. It is evident that, compared to MALA, or other algorithms previously mentioned, this HMC algorithm has the potential to reach the mode of the target and move around the space faster, an attribute of having momentum due to the Hamiltonian dynamics. On the other hand, there are now two tuning parameters,  $h$  and  $L$ , instead of only  $h$ , and a wrong choice can lead to very inefficient algorithms (see the magenta plot).

### 3.2.6 Piecewise-Deterministic Monte Carlo

For the preceding examples, the impact of working with ‘rough’ potentials  $U$  is felt most keenly through the difficulty of stably and accurately discretising the stiff dynamics of some associated stochastic process. Accordingly, the solutions focus on using more refined discretisation strategies, or seeking a modified stochastic process whose discretisation suffers less from these concerns of stiffness. A somewhat radical proposal is then to seek a new stochastic process for which discretisation is not only stable, but can even be avoided entirely. Other things being equal, the appeal of such an approach is reasonably clear; on the other hand, the existence of such an approach is much less clear.

In this regard, Piecewise-Deterministic Markov Processes play a rather distinguished role. While the Langevin and Hamiltonian approaches to Monte Carlo simulation are built around the numerical approximation of a dynamical system which is rich with information about the target distribution, but challenging

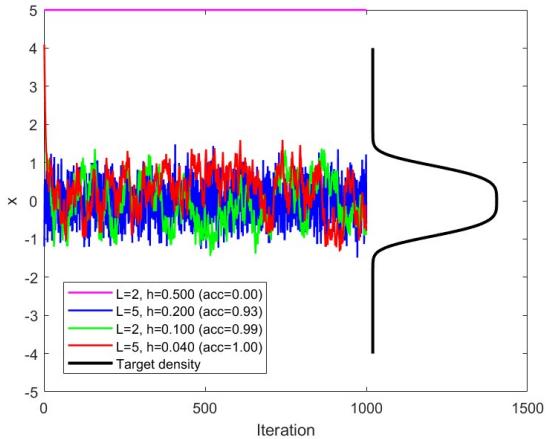


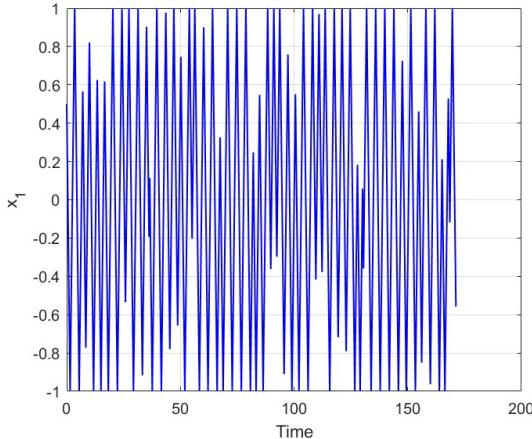
Figure 19: Hamiltonian Monte Carlo with kinetic energy  $K(v) = |v|^{\frac{4}{3}}$ . The target density ( $\pi(x) \propto \exp(-x^4)$ ) is overlaid on the right hand side of the graph. The legend indicates the various step-sizes ( $h$ ) and leapfrog steps ( $L$ ) along with the average acceptance probability.

to construct exactly, Piecewise-Deterministic Monte Carlo methods like the Bouncy Particle Sampler (BPS) and Zig-Zag Process (ZZ) instead explore the target distribution using rather rudimentary dynamics which are largely uninformed by the geometry of the target distribution, and ensure their long-term stability by intermittently punctuating these continuous dynamics with discontinuous jumps in velocity space. The usual story is a simple and compelling one: pick a direction, run in that direction until it starts to look like a bad idea; once that happens, pick a new direction ('bounce') in a reasonable way, and carry on running in that new direction. The details of how to select directions well are more involved, but are by now reasonably well-understood.

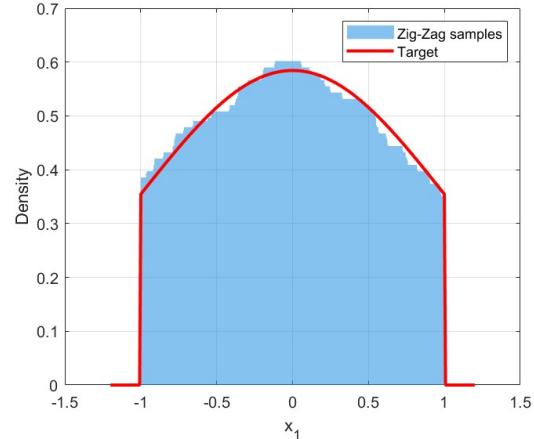
Deferring discussion of the (admittedly involved) practical aspects of simulating PDMPs, we can begin to highlight some of their appealing properties, focusing on the interplay with non-smooth targets. Firstly, most PDMPs of practical interest move, by design, at a constant speed  $s_*$ : no matter how strange an energy landscape one seeks to explore, there is no risk of spatial explosivity (though other pathologies concerning e.g. jumps require a more subtle discussion), and one even has the almost-sure bound  $\|X_t - X_0\| \leq s_* \cdot t$ . Secondly, while discretisation of steep drifts in ODEs and SDEs is often ‘direction-agnostic’ in its effects – a strongly-stabilising steep drift towards the center of the space can be as disastrous as a steep drift which causes the process to explode into the tails – the feedback mechanism in PDMPs is of a different flavour. More specifically, when the process is moving in directions for which the potential  $U$  is steeply decreasing, the process simply carries on in this favourable direction; by contrast, when the potential is steeply *increasing*, one rapidly experiences a ‘bounce’ event, and rectifies the situation. Non-smoothness and boundary effects are similarly muted in their impact on the movement of the process; if you start to move in a bad direction, then whatever the provenance of this badness, the response of the process is to propose a more fruitful direction, as quickly as possible.

In order to indicate the potential usefulness of these processes in sampling, we return to the setting of Example 3.4 and the box-constrained Gaussian defined in (9). In Figure 20, we present the numerical results of simulating Zig-Zag process on this target, forcing the process to change direction whenever it hits the boundary of the box. For example, if the  $i^{\text{th}}$  coordinate of the process moves in direction  $v_i = +1$  and hits the boundary at  $x_i = +1$ , the velocity’s  $i^{\text{th}}$  coordinate will immediately switch to  $v_i = -1$ . We present the traceplot and the histogram of the first coordinate samples, having ran the process until  $N = 10^6$  changes of direction occurred, either by hitting the boundary, or due to a random direction switch. This makes the comparison fair to the MALA algorithm, as number of evaluations of the gradient of the log-density (which is one of the most computationally expensive parts of the algorithms) were roughly the same for both algorithms. Comparing Figure 20 against Figure 9, it is evident from the traceplot that the process explores

the state space much faster while the histogram indicates that the process has created much more accurate samples from the target.



(a) Traceplot of the first coordinate.



(b) Histogram of the first coordinate

Figure 20: Zig-Zag Sampler on the  $10^4$ -dimensional box-constrained Gaussian target (9), ran for  $N = 10^6$  direction switches.

Due to their crucial historical role in the genesis of PDMC methods, we pause briefly to explicitly discuss the interaction between PDMPs and boundary effects, in the form of hard constraints. Some of the most impactful early developments in Piecewise-Deterministic Monte Carlo were made in the Statistical Physics literature under the name of ‘Event-Chain Monte Carlo’ (ECMC), where they were used for (among other things) studying the dynamics of large ensembles of hard spheres; see [13, 67, 84, 102, 106] for some indicative references from this community. To be more precise, write  $\mathbf{T}^d$  for the torus in dimension  $d$  with associated metric  $\mathbf{d}_{\mathbf{T}}$ , pick some large ‘number of particles’  $N \gg 1$ , and consider some positive radius  $R > 0$ . One then studies the uniform distribution over collections of  $N$  particles  $(x_1, x_2, \dots, x_N) \in (\mathbf{T}^d)^N$ , subject to the constraint that for all  $1 \leq i < j \leq N$ , there holds the ‘exclusion principle’  $\mathbf{d}_{\mathbf{T}}(x_i, x_j) \geq 2 \cdot R$ . That is, viewing each  $x_i \in \mathbf{T}^d$  as the centroid of a sphere of radius  $R$ , one imposes that any two distinct spheres cannot overlap. Such systems are of physical interest as a foundational model for phase transitions between the liquid and solid states of matter.

In this context, the ‘potential’ for the system is essentially non-existent, and so the challenge of sampling is governed entirely by the non-convex hard-sphere constraint. The core ECMC strategy thus proceeds by selecting a single particle, assigning it a velocity, propagating it in the direction of that velocity until it collides with another stationary particle, and then ‘bounces’ that particle into motion, transferring its velocity away. Interacting with the boundary no longer raises questions of how to tune an appropriate step-size to safely remain inside the feasible set, and becomes entirely a question of computing when the next collision occurs. Even without expanding upon the improved mixing behaviour conferred by the use of these piecewise-deterministic dynamics, this is already of substantial interest in view of the more pedestrian concern of how Markov process-based algorithms might interact more simply with hard constraints. In any case, the qualitative story here remains valid for other challenges which have been discussed here, including steep gradients, non-smooth targets, other boundary effects (including ‘stickiness’), and more. For some illustrative examples in statistical contexts, see e.g. [63, 66, 38].

Nevertheless, one must make some comments on numerical aspects of PDMPs. While discretisation *per se* might be avoidable, there are various challenges associated with the simulation of PDMPs, chief among them being the simulation of the ‘bounce’ events. For rather structured models, there are good strategies available for exact simulation by either harnessing the ‘shape’ of the event rate (e.g. in the case of convexity assumptions; see Example 1 in [30], or [130]) or by effective use of bounds on the event rate within thinning procedures [90]. Various other numerical strategies have been proposed for making this task feasible for

wider classes of model; these are arguably most successful at making ‘impossible’ models ‘possible’, rather than at making ‘possible’ models faster; see e.g. [54, 114, 41, 3] for examples to this effect.

Moving beyond the difficulty of simulating these events, there is also the compounding concern of event frequency. While theoretical results reveal that PDMPs can converge to their invariant measure quite quickly when viewed as continuous-time processes (see e.g. [48, 45, 96, 97]), the practical complexity of PDMPs is really measured by the ‘clock’ of how many events have occurred (or in some cases, how many events have been proposed), which changes the picture somewhere. Some heuristics based on product-form model problems suggest that for sampling problems on  $\mathbf{R}^d$ , one might expect to witness  $\Theta(d^{1/2})$  events per unit time at stationarity, even in the best (reasonable) case. Along similar lines, for potentials which vary irregularly, one might expect to witness events at higher frequencies; see [23] for a study of this phenomenon on ill-conditioned Gaussian target distributions. In any case, there is no free lunch, and the principle of ‘conservation of difficulty’ is borne out to some extent: even for rough targets, PDMPs can admit stable dynamics, but one will eventually have to work hard to simulate those stable dynamics.

### 3.3 Open Problems

We here highlight some problems which are recurrent in this sub-literature, and whose resolution would be of some major practical and theoretical interest.

#### **Robustness to Roughness and Hyperparameter Adaptation**

A number of works [8, 95] make the empirical and intuitive comment that when using MCMC algorithms based on more ‘robust’ proposals in this way, adaptation of algorithmic hyperparameters (e.g. step-size, covariance matrix, etc.) is more stable and efficient. It is very appealing to attach some theory to this observation, which in practice, turns out to provide a very strong argument in favour of such methods.

#### **Robustness to Roughness without Overconservatism**

Some of the lightweight approaches to robustification effectively amount to replacing a superlinear drift with a bounded drift. We know that processes of bounded drift will typically not be any better than exponentially ergodic (e.g. it becomes difficult to satisfy a Logarithmic Sobolev Inequality). Can we constrain our drift more gently so that we retain numerical stability, but without hindering our ability to return to the bulk of the state space exponentially quickly? Proximal MCMC offers a partial solution here, but makes certain implicit assumptions about the feasibility of solving certain optimisation problems efficiently. A more generic and lightweight solution would be of substantial interest. We note some recent answers in this direction [83, 100] which have been obtained by enriching the usual ‘taming’ framework.

#### **Robustness to Roughness by Local Adaptation**

Our solutions here have effectively focused on resolving instabilities with rather global changes to the dynamics. It would be of substantial interest to devise robustification strategies with a more localised character, i.e. adaptively detecting challenges in the local geometry, and modifying the dynamics only when needed. This characterisation is rather too vague to be useful at present, but it stands to reason that such a strategy could be made both precise and practical.

## Robustness to Roughness by Discretisation

Several proposed strategies have the flavour of discretising the spatial domain, and hence reduce the overall sampling problem to a sequence of discrete-space Markov chains; this is the case for the Barker-Langevin approach of [95] and some of its relatives (e.g. [29, 46]). This simple strategy can be quite successful, even when done in an ad-hoc way. What are principles for performing such a discretisation ‘optimally’?

## Model Problems for Roughness

Finally, our classification of ‘rough’ problems is rather binary in nature. It would be useful to develop a richer spectrum of model problems which reflect the variety of sampling challenges posed in modern applications, and to develop a more refined language for delineating between their distinct pathologies. This should then provide a clearer path towards resolving said challenges.

# 4 Heavy-Tailedness

## 4.1 Formulation of the Pathology and Examples

In this section we will focus on a class of targets that assign a significant amount of mass at the tails of the state space. We will be working with the following definition.

**Definition 4.1.** *We will call a density  $\pi \in C^1(\mathbf{R}^d)$  heavy-tailed if*

$$\lim_{r \rightarrow \infty} \sup \{ \|\nabla \log \pi(x)\| : \|x\| \geq r \} = 0.$$

Note that this is more restrictive than various other standard definitions of heavy-tailedness which focus only on the probability mass which is contained in the tails of the distribution, e.g. that  $\int \pi(dx) \exp(s \cdot \|x\|) = \infty$  for all  $s > 0$ . We choose to work with this formulation because of the central practical role of this gradient in many popular ‘local’ MCMC algorithms. If one sought to further highlight this distinction, then one might rename our definition as ‘flat-tailed’.

Morally, one can think of this class as the densities with tails heavier than any exponential distribution. Observe, for example, that a symmetrised one-dimensional exponential distribution on  $\mathbf{R}$  is of the form  $\pi(x) = \frac{a}{2} \cdot \exp(-a|x|)$  for some  $a > 0$ , whereby it holds that  $|\nabla \log \pi(x)| = a$ , i.e. the gradient of the potential acts with a force of constant magnitude.

Distributions with heavy tails arise naturally in settings when one tries to model extreme or rare events (for example in finance [52], see also [92] for more general applications). They also often arise in practical Bayesian statistics, for example in models that aim to account for the behaviour of outliers in a more robust manner (e.g. [85]). The resulting posteriors typically exhibit heavy tails and one needs to create samples from that posterior in order to perform inference. The same phenomenon can occur when one uses a shrinkage prior (e.g. the horseshoe prior [140]); this can be natural in settings wherein one is trying to estimate many parameters, of which only a small subset are expected to be significantly non-zero.

From a sampling point of view, MCMC techniques tend to struggle with such targets, as the dynamics of the chain fail to systematically explore the tails of the distribution, wherein much of the probability mass resides. Furthermore, in contrast to the targets considered in the previous section, for targets with heavy tails, even ideal continuous processes (such as the Langevin dynamics) tend to be quite slow at exploring. For example, let us recall the example of Overdamped Langevin diffusion (2), given by

$$dX_t = \nabla \log \pi(X_t) dt + \sqrt{2} dW_t.$$

Since  $\pi$  is heavy-tailed, when  $X_t$  is exploring the (appreciable) probability mass contained in the tails, it holds that  $\nabla \log \pi(X_t) \approx 0$ , and so the dynamics of  $X_t$  are largely governed by the Brownian motion term, with little impact from the shape of  $\pi$  itself. This results into backtracking dynamics and prohibits the process from returning to areas of higher density sufficiently fast, leading to inefficient space exploration.

Similar random-walk behaviour on heavy-tailed targets is observed for other continuous-time processes discussed earlier in this work, such as the Bouncy Particle Sampler or the Zig-Zag process (see e.g. the discussion in [134, 133]). The problem there is further emphasised by the fact that by construction, these PDMPs move around the space with constant speed, which precludes the possibility of transitioning between the tails and the centre of mass at a sufficiently fast rate to e.g. ensure exponential ergodicity.

As a result of this inefficient behaviour, MCMC proposals that arise as numerical discretisations of the continuous-time processes explored earlier cannot hope to have reliable performance on heavy-tailed targets, as the underlying process which they seek to approximate fails to exhibit such a performance. We will showcase this with the following two examples.

**Example 4.1** (Cauchy distribution). *Our first example in the heavy-tailed category is the Cauchy distribution, i.e. when the target distribution is of the form*

$$\pi(x) \propto (1 + x^2)^{-1}. \quad (15)$$

*Observe that under  $\pi$ , no moment of order  $\geq 1$  exists, reflecting that  $\pi$  places a great deal of mass around large values of  $|x|$ . Observe also that  $|\nabla \log \pi(x)| \leq \min\{1, 2 \cdot |x|^{-1}\}$ , i.e. for large  $|x|$ , the landscape of  $\pi$  is quite flat.*

To showcase the sampling problem in the setting of Example 4.1, in Figure 21 we present the results of a numerical simulation from the one-dimensional Cauchy distribution using the MALA algorithm. The first two plots show the traceplot. It can be seen that there are long and infrequent excursions at the tails, which typically indicate unstable algorithmic performance. We also present the autocorrelation plot, which shows how rapidly correlation between samples decays as a function of the time-distance between samples. Ideally, the correlation should decay fast, which would indicate a fast generation of independent samples, without the need to run the algorithm for many iterations. In this example, autocorrelations appear to decay quite fast, but we will later see variants of the algorithm, designed to target heavy tailed targets that will improve the autocorrelation decay for the same target. Finally, the last plot demonstrates the mean squared relative error in estimating the probability of the event  $\{X \geq 5\}$ , with  $X$  following a Cauchy distribution. The plot reports the average squared relative errors over 100 independent realisations of the chain, as a function of algorithmic iterations. It is clear that the squared relative errors decay slowly as the algorithm evolves in time, barely reaching a value less than 0.8 in the first 10000 iterations. This is indicative of the fact that the algorithm is not well-suited to explore the tails of the distribution. We emphasise here that the step-size of the algorithm (in this case  $h = 10$ ) was carefully chosen to maximise the algorithmic performance in terms of Effective Sample Size (see e.g. [32]), and the algorithmic behaviour can be considerably worse with under-optimised step-sizes.

We now present a somewhat more realistic statistical model which exhibits heavy tails both in the parameters and in the observations.

**Example 4.2** (Cauchy Regression with Horseshoe Prior). *In the context of the linear model, one approach to robustly capture the behaviour of outliers (e.g. [85]) is to model the observational errors as following the Cauchy distribution. In this model, given parameters  $\beta = (\beta_1, \beta_2, \dots, \beta_p)^\top$ , and design matrix  $\mathbf{X} \in \mathbf{R}^{n \times p}$ , the response vector  $\mathbf{Y} \in \mathbf{R}^n$  is given by*

$$\mathbf{Y} = \mathbf{X}\beta + \epsilon, \quad (16)$$

*with the errors  $\epsilon = (\epsilon_1, \dots, \epsilon_n)^\top$  following i.i.d. Cauchy distributions as in (15). In order to allow additional flexibility to the model, a widely used prior for the parameters  $\beta$  is the horseshoe prior, which is defined hierarchically by setting*

$$\{\beta_j\} \mid \{\lambda_j\} \stackrel{\text{ind}}{\sim} \mathcal{N}(0, \lambda_j), \quad \lambda_j \stackrel{\text{iid}}{\sim} \text{Cauchy}^+. \quad (17)$$

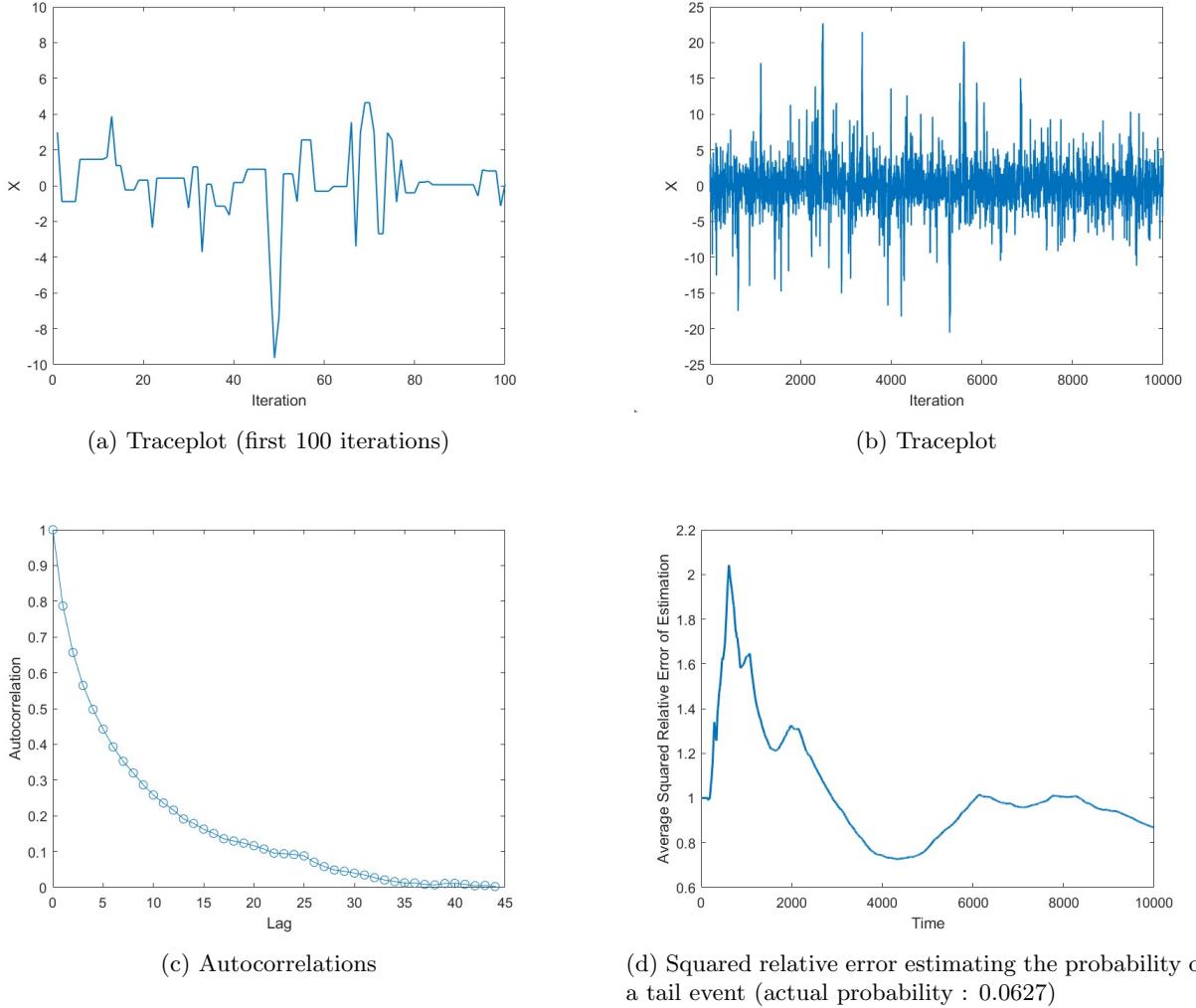


Figure 21: MALA algorithm on a one-dimensional Cauchy target.

where  $\text{Cauchy}^+$  is the standard Cauchy random variable, conditioned on being positive. This has the effect of allowing some of the regression coefficients to take values close to zero, thus potentially reducing the effective dimensionality of the parameter space, with other coefficients free to take much larger values (due to the relatively heavy marginal tails). In this case, the joint posterior over  $(\beta, \lambda)$  takes the form

$$\pi(\beta, \lambda | \mathbf{X}, \mathbf{Y}) \propto \left( \prod_{j=1}^p \frac{1}{1 + \lambda_j^2} \right) \cdot \left( \prod_{j=1}^p \frac{1}{\sqrt{2\lambda_j}} \exp\left(-\frac{1}{2\lambda_j}\beta_j^2\right) \right) \cdot \left( \prod_{i=1}^n \frac{1}{1 + \left(Y_i - \sum_{j=1}^p \beta_j X_{i,j}\right)^2} \right), \quad (18)$$

which presents heavy tails for the parameters  $\lambda$  and for the marginal of  $\beta$  (see e.g. [35]).

The slow behaviour of MALA in Example 4.2 is showcased in Figure 22, where we summarise the findings of a numerical simulation with  $p = 5$  parameters and  $n = 20$  observations. The true values of  $\beta_j$  were drawn from a horseshoe distribution (17), the covariates  $X_{i,j}$ 's were drawn from the uniform distribution in  $\{0, 1, 2, 3\}$ , and conditionally upon these values, response data  $Y_i$  was drawn from the model (16). We used the MALA algorithm to draw samples from the marginal posterior (18). Figure 22a shows the traceplot of the first coordinate  $\beta_1$ , while Figure 22b shows the autocorrelation plot of  $\beta_1$ . Evidently, the autocorrelations decay as a very slow rate, indicating large dependence between the generated samples, and poor mixing.

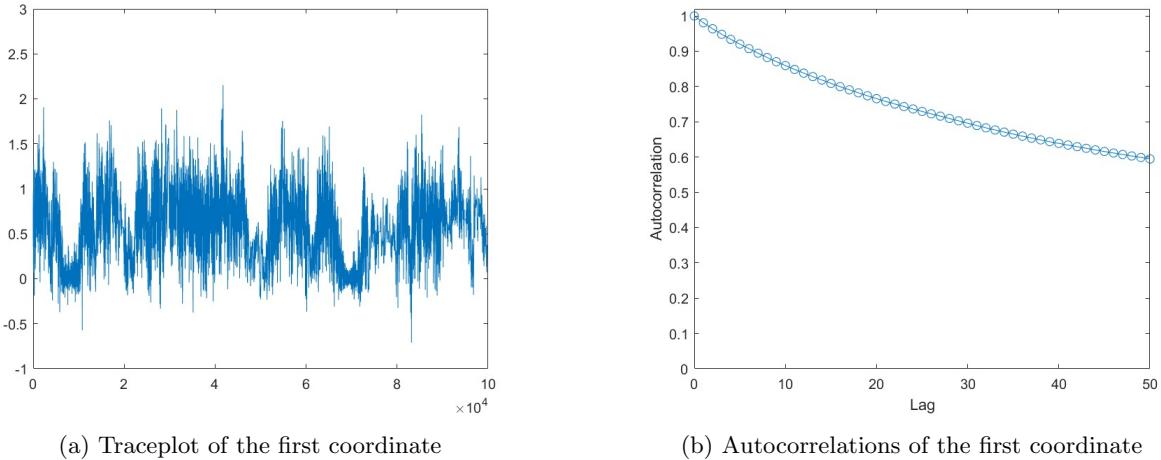


Figure 22: MALA algorithm on a 5-dimensional Cauchy-Regression target.

In the next section, we will review contemporary approaches that aim to produce more robust and efficient algorithms for this setting. We note that by contrast with the case of rough targets, the solutions are necessarily somewhat less generic. This can be understood as a symptom of the general phenomenon that in order to alleviate the effects of heavy tails, one often requires an understanding of the degree of heavy-tailedness in the problem. With this being said, the overall strategies are reasonably general in spirit, but do require a greater degree of instance-specific tuning than is typical of the light-tailed world.

## 4.2 Proposed Solutions

We will discuss two approaches that have been proposed in order to improve the performance of MCMC algorithms on heavy tails. These involve transforming either the state space or the time scale in which the process runs. In the first approach, one transforms the state space via a function  $f$  which, if carefully chosen, induces a push-forward density  $\tilde{\pi}(x)$  with significantly lighter tails than  $\pi$ . The idea is to then run a more efficient MCMC algorithm to target the light-tailed  $\tilde{\pi}$ , and then push the generated samples back through the function  $f^{-1}$  in order to recover approximate samples from  $\pi$ . The second approach also involves changing the target to a density  $\tilde{\pi}$  with lighter tails, but instead of filtering the samples through a space transformation, one changes the time-scale in which the process moves, effectively discarding some of the samples. This biases the  $\tilde{\pi}$ -generated samples in an appropriate fashion so that the remaining samples are distributed according to  $\pi$ .

### 4.2.1 Space Transformations

In the context of heavy-tailed sampling, the core idea behind space-transformation strategies is to identify an appropriate invertible function  $f : \mathbf{R}^d \rightarrow X$ ,  $X \subset \mathbf{R}^d$  so that the push-forward measure  $\tilde{\pi}$  on  $X$ , defined by  $\tilde{\pi}(A) = \pi(f^{-1}(A))$  has lighter tails than  $\pi$ . One then takes advantage of MCMC algorithms that are more efficient when targeting  $\tilde{\pi}$ , in order to construct samples  $X_1, \dots, X_n$  approximately distributed according to  $\tilde{\pi}$ . Pulling these samples back to the original space by application of  $f^{-1}$  then yields samples  $f^{-1}(X_1), \dots, f^{-1}(X_n)$  that are approximately distributed according to  $\pi$ .

**Example 4.3** (Cauchy target). *To showcase how this approach can be useful, let us consider the one-dimensional Cauchy distribution from Example 4.1, along with the space transformation*

$$f(x) = \text{sign}(x) \cdot \log(1 + |x|). \quad (19)$$

Consider the push-forward measure  $\tilde{\pi}(A) = \pi(f^{-1}(A))$ , where  $\pi$  is the Cauchy distribution. For the density of  $\tilde{\pi}$ , the change-of-variables formula yields that

$$\tilde{\pi}(y) = \pi(f^{-1}(y)) \cdot |\mathrm{D}f^{-1}(y)| = \frac{\pi(f^{-1}(y))}{|\mathrm{D}f \circ f^{-1}(y)|}.$$

We calculate that  $f^{-1}(y) = \text{sign}(x) \cdot (\exp|y|-1)$ , and  $\mathrm{D}f(x) = (1+|x|)^{-1}$ , so that one obtains

$$\pi(f^{-1}(y)) \propto \frac{1}{1 + (f^{-1}(y))^2}, \quad \mathrm{D}f \circ f^{-1}(y) = \exp|y|,$$

and therefore

$$\tilde{\pi}(y) \propto \frac{\exp|y|}{1 + (\exp|y|-1)^2} = (2 \cosh y - 1)^{-1} = \exp(-|y| + O(1))$$

with lighter tails, similar to those of Laplace distribution. One can then use an MCMC algorithm to target the distribution  $\tilde{\pi}$ , in order to efficiently generate samples  $Y_1, \dots, Y_n$  that are approximately distributed according to  $\tilde{\pi}$ , and then interpret the samples  $X_i := f^{-1}(Y_i) = \text{sign}(Y_i) \cdot (\exp|Y_i|-1)$ ,  $i = 1, \dots, n$  as approximate samples from  $\pi$ .

There have been various suggestions on which space transformation to use, and this may depend on the MCMC algorithm one uses to generate samples from the transformed target  $\tilde{\pi}$ . For target distributions whose tails decay at a polynomial rate, one can see that transformations of logarithmic growth (such as the one considered in (19)) can be effective choices, as they push the target forward onto a distribution which is asymptotically approximately log-concave, without impacting the smoothness properties of the implied potential too aggressively.

Multivariate examples of this form of transformations were explored in [82], where the authors suggested the use of isotropic transforms, i.e. functions of the form

$$f(x) = h(\|x\|) \cdot \frac{x}{\|x\|}, \quad x \in \mathbf{R}^d \tag{20}$$

for some appropriate increasing  $h \in C^1(\mathbf{R}_+)$ .

In order to visualise the impact of such a space transformation, on the left plot of Figure 23 we simulated 25 points uniformly at random in the box  $[-10, 10]^2$ . We then considered the function  $f$  as in (20), with  $h(r) = \log(1+r)$ , essentially generalising (19) to two dimensions. We apply the function  $f$  to the uniformly simulated points, and the resulting points are shown in the right plot of Figure 23. It is evident that the transformed points are much more concentrated towards the origin, indicating that their distribution has lighter tails.

Parameterising  $f$  following Equation 20 offers some simplicity, in the sense that it reduces the task of searching for an effective high-dimensional transport map into a problem which is essentially univariate in nature. In concordance with the discussion above, the authors of [82] suggest that for target distributions with polynomially-decaying tails, one ought to take  $h$  which grows logarithmically for large values of  $\|x\|$ . They then show that when using the Random Walk Metropolis algorithm (see e.g. [127]) on the transformed target, one can obtain much-improved rates of convergence, compared to applying the same method to the original target. In particular, the authors prove that the transformed algorithm is exponentially ergodic in scenarios where the base method is not.

Similarly, [69] introduced the Transformed Unadjusted Langevin Algorithm. The space transformation is also isotropic, i.e. of the form of (20), and the function  $h$  is taken to have polylogarithmic growth for large arguments. The authors then make use of ULA (3) to explore the transformed target, and demonstrate theoretically that by a careful design of the transform map, one can gain improved rates of convergence.

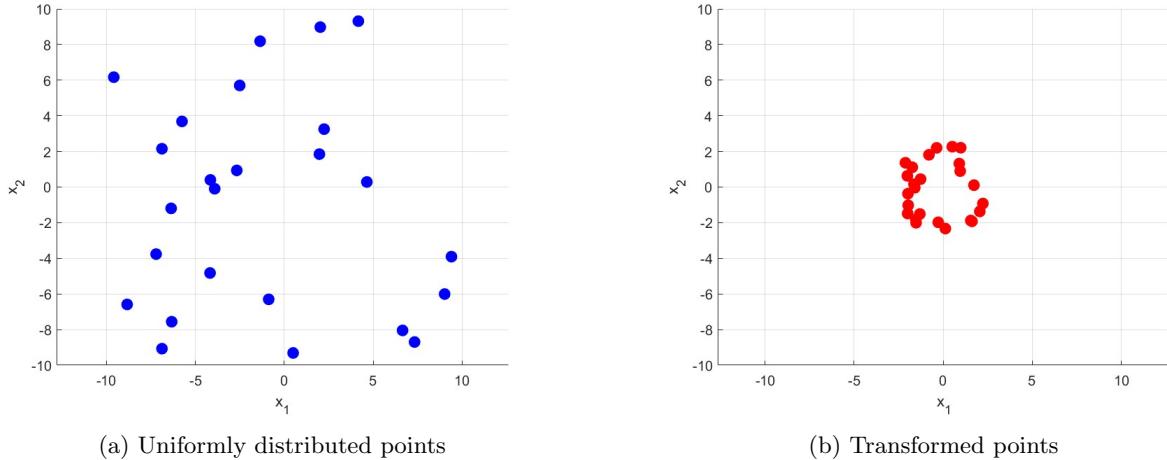


Figure 23: Uniformly distributed points on the square  $[-10, 10]^2$  and the transformed points after applying the transformation  $f(x) = \log(1 + \|x\|) \cdot x / \|x\|$ .

Similarly, by introducing a Metropolis-adjustment, one can run an instance of MALA on the transformed target, ensuring that the samples are asymptotically consistent, before transferring the samples back to the original space via the inverse transformation. We use this approach in our next simulation, which uses a space transformation with  $f$  as in (19) to target the Cauchy distribution from Example 4.3. As explained in that Example, the resulting push-forward measure  $\tilde{\pi}$  will have lighter tails, therefore, it will be easier for MALA to explore the area where most of the target mass concentrates. We present our results in Figure 24. This can be compared against Figure 21, where we ran MALA directly on the heavy tailed target. The first two plots show the traceplot, and in contrast to Figure 21, the algorithm seems to have much more frequent and short tail excursions, indicating that the process explores the state space sufficiently fast. This is further emphasised by the autocorrelation plot, which shows a much faster decay of autocorrelations, compared to the decay observed in Figure 21. Finally, the last plot demonstrates the mean squared relative error in estimating the probability of the tail event  $\{X \geq 5\}$ , with  $X$  following a Cauchy distribution. We use the same approach as the one discussed when running a MALA algorithm directly on the Cauchy target, reporting the average squared relative errors over 100 independent realisations of the chain. Comparing with Figure 21, it is evident the errors decay much faster, having reached a value less than 0.1 before 1000 iterations, and staying consistently small thereafter. As a note, we comment that the step-size of the algorithm did not influence the algorithmic performance as extensively as in the previous case, where MALA was targeting the Cauchy target directly. Here we have picked the step-size to be  $h = 0.5$ .

We also consider the Transformed MALA on the Cauchy regression model, with Horseshoe prior, as in Example 4.2. As in Section 4.1, we take  $p = 5$  and  $n = 20$  observations. We again use a spherically symmetric transformation  $f$  as in (20), with  $h(r) = \log(1 + r)$ , as in Figure 23. We ran MALA on the transformed target and obtain our primal samples by applying the inverse transformation. The step-size was chosen so that the average acceptance rate was between 50 – 60%. As in Example 4.2, we focus on the first coordinate of the process  $\beta_1$  and present the traceplot and the autocorrelation plot of the algorithm in Figure 25. In this case, the traceplot seems a bit more stable compared to Figure 22 where MALA was directly targeting the posterior. On the other hand, while there is a small improvement on the rate of decay of the autocorrelations (which reach the value 0.45 after 50-lag, compared to the value of 0.60 the MALA had reached), the autocorrelation decay is still extremely slow. This indicates that while space transformation can improve the performance of an algorithm, the multi-dimensionality of the problem can interact with the transformation to create additional problems. One problem that is of particular interest in this setting is that different directions in the target density can exhibit different rates of decay, which one expects to be somewhat generic behaviour in complex, high-dimensional contexts. The goal of the transformation  $f$  is then to create a transformed target with rates of decay which are as uniform as possible, and a spherically-

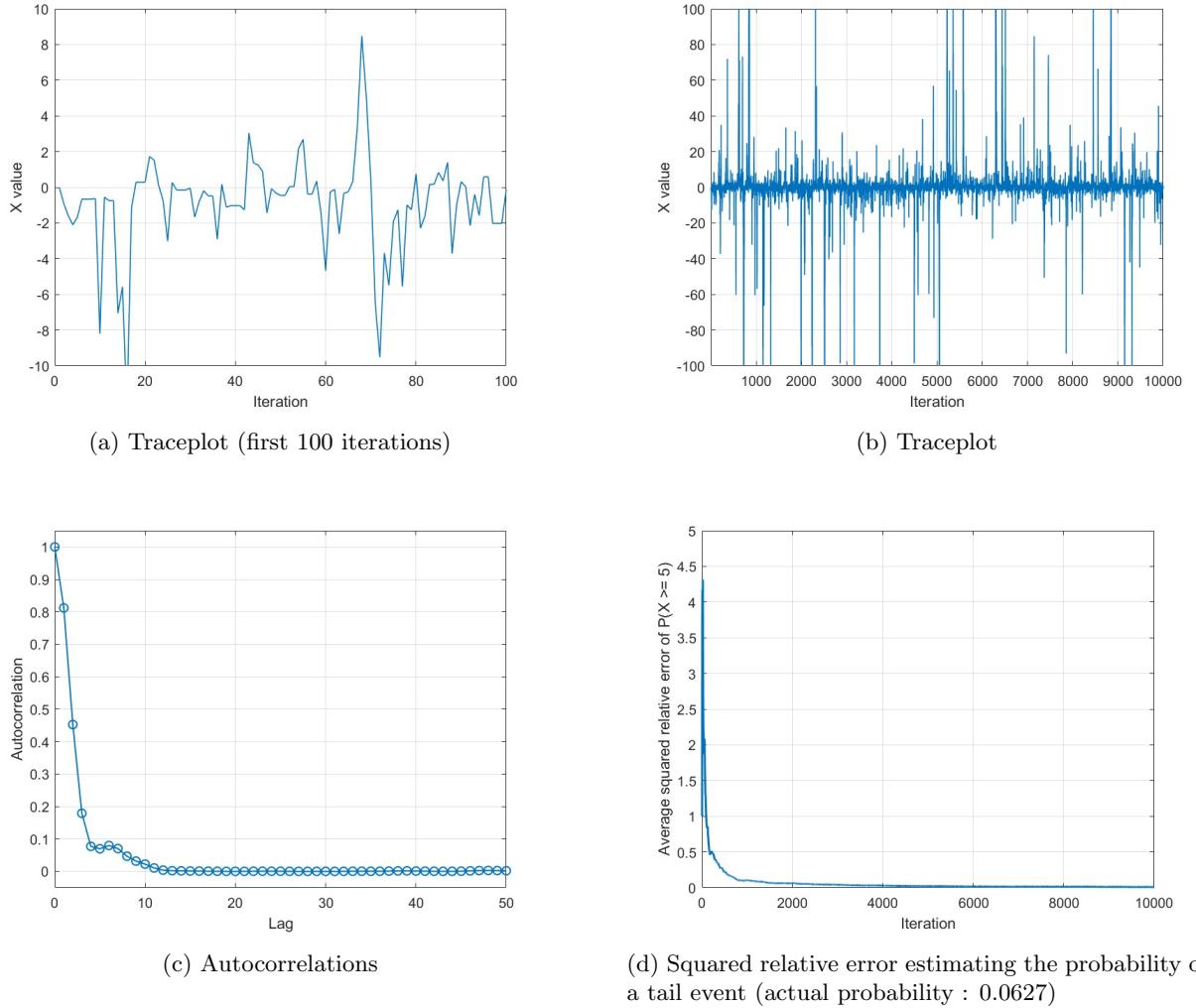


Figure 24: Transformed MALA algorithm on a one-dimensional Cauchy target.

symmetric transform is limited in its ability to achieve this. In this regard, some extra care is needed when designing transformations for high-dimensional problems.

A strategy with a slightly different flavour is to instead construct space transformations which push  $\mathbf{R}^d$  into a compact set, and then run the MCMC there. One benefit of this strategy is that if the compactifying map leads to a reasonably well-behaved pushed-forward target, then there is reason to expect that the convergence profile of the chain will depend much more gently on the initialisation, on the basis that all points within the support of the target are within a certain *bounded* distance of one another.

In this direction, [139] recently suggested a space transformation that maps  $\mathbf{R}^d$  to the unit sphere in one dimension higher, i.e.  $X = \mathbf{S}^d \subset \mathbf{R}^{d+1}$ . This is done using the stereographic projection (see e.g. [87]) depicted in Figure 26. A relatively standard MCMC algorithm is then used to sample from the transformed target on the sphere, and the authors show that under reasonable conditions (which accommodate various degrees of heavy-tailedness), one recovers exponentially-fast convergence to the target in total variation distance. Perhaps even more impressively, this rate of convergence in total variation distance can be bounded *independently of the starting position* of the algorithm, i.e. the process is *uniformly ergodic*. On the other hand, due to the dramatic fashion in which the stereographic projection warps the state space, some care is again required in centering and scaling the target distribution. For further work which studies the challenges

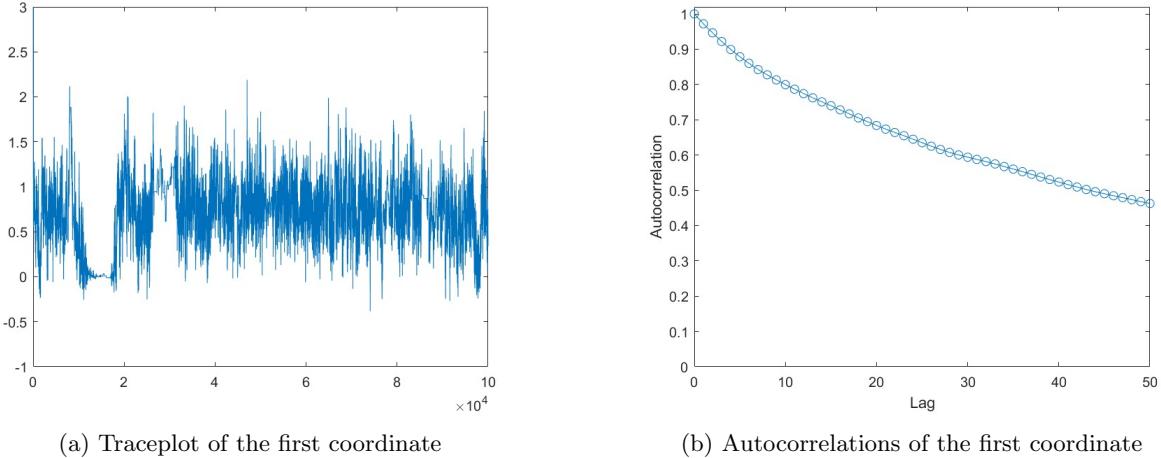


Figure 25: Transformed MALA algorithm on a 5-dimensional Cauchy-Regression target.

associated with this adaptation, see also [11].

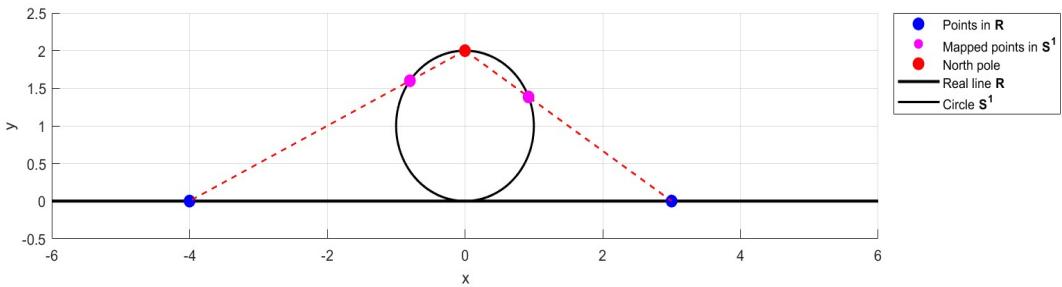


Figure 26: Stereographic projection from  $\mathbf{R}$  to the circle  $\mathbf{S}^1$ . The x-axis is the real line. The figure shows how each point in  $\mathbf{R}$  (blue points) is mapped to the circle (magenta points). The magenta points are at the intersection between the circle and the line connecting each real (blue) point to the north pole (red point).

Finally, we note that other more elaborate space transformations have been suggested in the MCMC literature and applied with some success (see e.g. [78, 113, 32, 74]), but often without focusing on heavy-tailed target applications. In view of our focus here, we therefore omit further discussion of these approaches.

#### 4.2.2 Time Transformations

Rather than transforming the state space, a second approach to deal with the heavy-tailedness of the target is to apply a time-transformation (also known as a *time change*) to the process, whereby the evolution of time itself within the algorithm is altered. Equivalently, one introduces a state-dependent ‘speed function’ which dictates how quickly the process moves given its current position. Choosing an appropriate speed function can thus allow the process to move around the space faster and explore the tails more efficiently.

For example, the long but infrequent excursions of the process at the tails (showcased in Section 4.1), can become much shorter but persistent when the movement is regulated through a speed function, leading to more consistent estimation.

A general framework for how to define a time-changed process can be found in Chapter 6 of [53], while a unifying framework for employing these processes for MCMC sampling was developed in [15]. A process  $(X_t)_{t \geq 0}$  on  $\mathbf{R}^d$  is defined to be a *time-change* of a process  $(Y_t)_{t \geq 0}$  if for all  $t \geq 0$

$$X_t = Y_{r(t)}, \text{ where } r(t) := \int_0^t s(X_u) du \quad (21)$$

for some function  $s : \mathbf{R}^d \rightarrow (0, +\infty)$ . The process  $Y$  is called the *base process* and the function  $s$  is the *speed function*. Note that on a heuristic level, it holds that

$$\frac{dX_t}{dt} = \frac{dY_{r(t)}}{dr(t)} \cdot \frac{dr(t)}{dt} = \dot{Y}_{r(t)} \cdot s(X_t).$$

Thus, the process  $X$  traverses the path of the process  $Y$  but does so ‘ $s$  times faster’. Figure 27 illustrates this phenomenon. The background path is the one of a ‘canonical’ Zig-Zag process (which here plays the role of the base process  $Y$ ), moving with constant speed. The red points denote the position of the time-changed process  $X$  at times  $0.5, 1, 1.5, 2, 2.5, 3, \dots$ , after applying a time-change with speed  $s(x) = (1 + |x|^2)^{\frac{1}{2}}$ . While these points are equi-distant in time, it is evident that they are not equi-distant in space, since the process  $X$  moves with varying speed, depending on the current position. The time-changed process is targeting a bivariate, isotropic Cauchy distribution, i.e.

$$\pi(x) \propto (1 + x_1^2 + x_2^2)^{-\frac{3}{2}}.$$

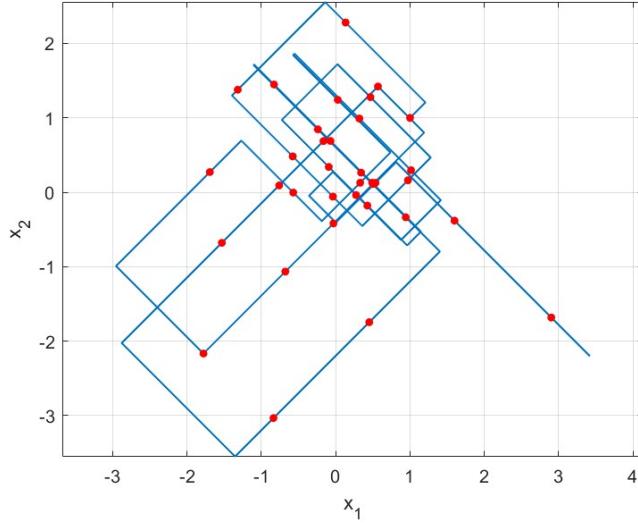


Figure 27: Scatter plot of a time-changed Zig-Zag process. The red points denote the position of the process every  $h = 0.5$  time units. Target: 2-dimensional Cauchy,  $\pi(x) \propto (1 + x_1^2 + x_2^2)^{-\frac{3}{2}}$ . Speed function:  $s(x) = (1 + x_1^2 + x_2^2)^{\frac{1}{2}}$ .

Conceptually, using a time-changed process  $X$  has a lot in common with the space-transformations discussed in Section 4.2.1. A useful perspective to consider when devising time-changed processes for the purpose of sampling is the following: assuming that  $s$  is  $\pi$ -integrable, continuous and bounded below by a positive constant, it holds that

$$X \text{ is } \pi\text{-invariant} \iff Y \text{ is } \tilde{\pi}\text{-invariant}, \quad (22)$$

where the measure  $\tilde{\pi}$  is defined as

$$\tilde{\pi}(dx) = \frac{1}{\pi(s)} s(x) \pi(dx) \quad (23)$$

with  $\pi(s) = \int_{\mathbf{R}^d} s(y) \pi(dy)$  the induced normalising constant.

Bearing this in mind, if one wants to use the time-changed process  $X$  to create samples from  $\pi$ , one ought to proceed as follows: construct a base process  $Y$  that targets the modified measure  $\tilde{\pi}$ , and then apply an appropriate time-change to adjust the speed of the process, ensuring that one spends more time in relevant areas of the space, ultimately biasing the samples towards  $\pi$ . In this regard, at a high level, one can say that both space- and time-transformation address the problem of the target having heavy tails by targeting a different measure  $\tilde{\pi}$  and then ‘correcting’ the samples.

An important distinction between the two approaches, however, lies in the type of tails induced on the transformed distribution  $\tilde{\pi}$ . Space transformations typically lead to a lighter-tailed  $\tilde{\pi}$ , under the expectation that the process targeting  $\tilde{\pi}$  will explore the target more efficiently, thereby yielding higher-quality samples after applying the inverse transform. By contrast, when applying time transformations, the tails of  $\tilde{\pi}$  as in (23) become even heavier, provided that  $s(x) \xrightarrow{|x| \rightarrow \infty} \infty$ . This allows the base process  $Y$  (targeting  $\tilde{\pi}$ ) to exhibit long excursions in the tails. These regions are subsequently visited more rapidly by  $X$ , leading to improved tail exploration. Indeed, [15] shows that under assumptions on the speed function  $s$ , applying a time-change can lead to exponential and even uniform ergodicity, even for heavy tailed targets, where the base process  $Y$  would converge much slower.

On the other hand, using a time-transformed process can be challenging from a practical point of view and can significantly increase the computational cost of the algorithm. This increased difficulty arises from the fact that one must be able to analytically integrate the function  $s$  along the path of the process (see, for example, the definition of the quantity  $r$  in (21)). From this perspective, base processes that move in simple ways can be more promising candidates for use in a time-transformed setting, at least from a computational standpoint. Two prominent examples are the Zig-Zag process and the Bouncy Particle Sampler, introduced in Section 3.2.6, since between jumping events these processes move in straight lines. As shown in [134], where a time-transformed Zig-Zag process was introduced under the name “Speed Up Zig-Zag,” between jumping events, determining the position of the time-changed process at a given time amounts to solving a one-dimensional ordinary differential equation (ODE), which the authors argue can be done exactly for speed functions of the form

$$s(x) = (1 + |x|^2)^{\frac{1+k}{2}}, \quad k \in \mathbf{N}.$$

Of course in a practical setting, when the target is high dimensional, one might want to use other speed functions. For example, when the tails of the target have different rates of decay across different directions, a well-designed speed function ought to reflect this.

To showcase the performance of time-transformed processes, we revisit our examples from Sections 4.1 and 4.2.1. We first consider the Cauchy target from Example 4.1. We run a time-transformed Zig-Zag process, using the speed function  $s(x) = (1 + x^2)^{\frac{1}{2}}$ . For a somewhat equal comparison with Figure 21 we ran the algorithm for the same number of target evaluations as in Section 4.1. We present our results in Figure 28, presenting four plots of the same type as when we ran MALA directly on the target (Figure 21). It is evident from the traceplot that the process explores the space well, having frequent and short tail excursions, while the autocorrelation plot decays sufficiently fast. Estimation of the tail event  $\{X \geq 5\}$  also becomes very accurate after a few algorithmic iterations.

We also consider the time-transformed Zig-Zag on Example 4.2, the Cauchy regression model with a Horse-shoe prior, with  $p = 5$  and  $n = 20$  observations. As in the one-dimensional Cauchy target, we ran the algorithm for the same number of target evaluations as in Section 4.1. We note, however, that the time-changed process required substantially more wall-clock computation time to run (approximately two orders of magnitude more), as the routine must repeatedly optimise complicated functions.

Nevertheless, focusing on the first coordinate  $\beta_1$ , we present the trace plot and the autocorrelation plot of

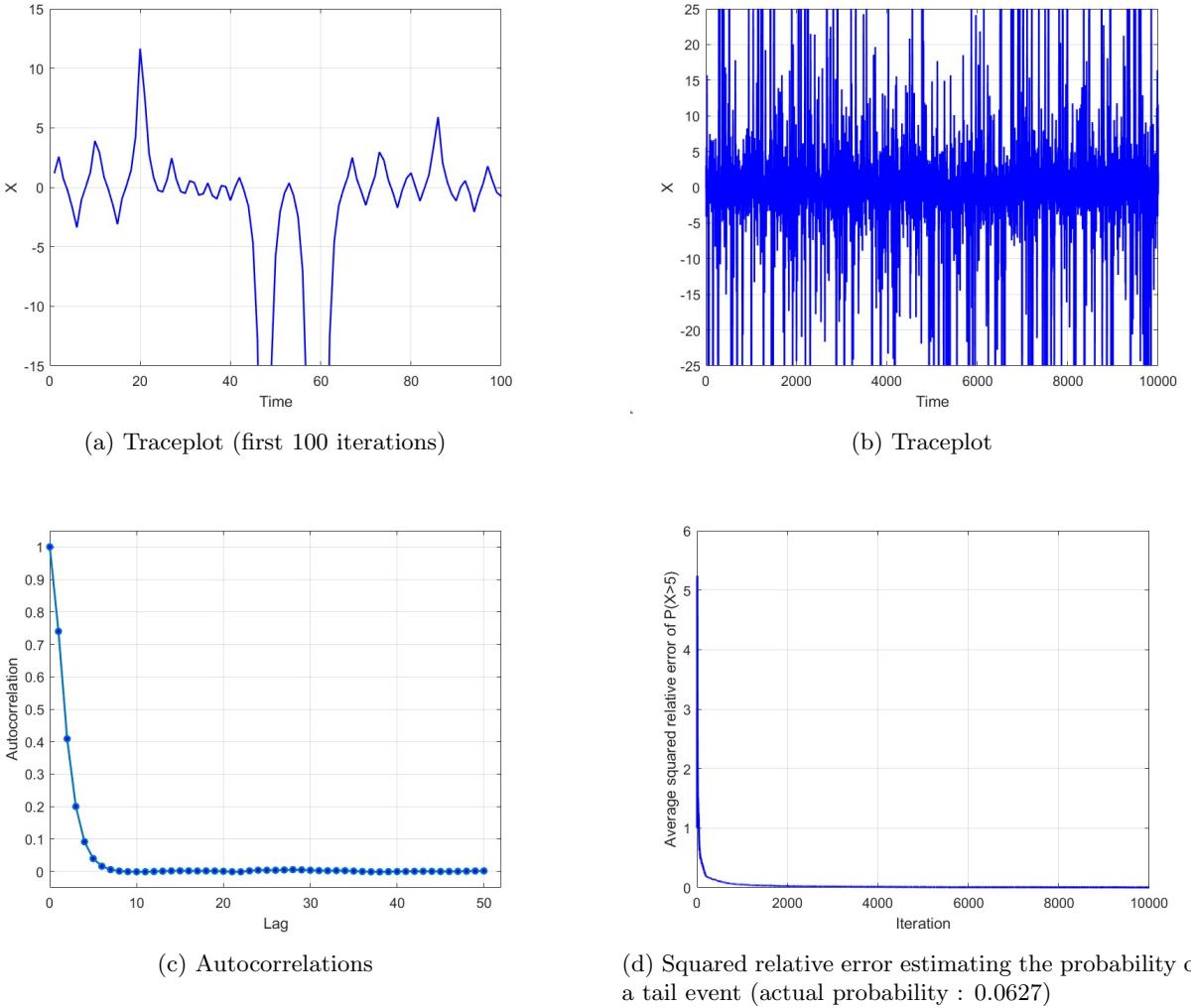


Figure 28: Time-transformed Zig-Zag algorithm on a one-dimensional Cauchy target. Speed function:  $s(x) = (1 + x^2)^{\frac{1}{2}}$ .

the algorithm in Figure 29. Compared with Figures 22 and 25, the trace plot appears substantially more stable, and the autocorrelation exhibits a markedly faster decay.

Time-transforming a process for sampling was considered in [123] where the authors studied time-changes of diffusion processes, using a speed function of the form  $s(x) = \pi(x)^{-a}$ ,  $a \in (0, 1)$ . For more recent uses of time-changing diffusion for sampling, see also [88, 89]. There are also connections with the idea of self-normalised importance sampling (see e.g. [121]), in the sense that rescaling time according to a speed function can have a similar effect to adjusting the weight of particular samples from the base process. In that direction, [2] use as base process a discrete time Markov chain with stationary distribution different than the target and adjust it by letting it spend a random number of iterations at each state, with the expected occupation time playing a similar role to the (reciprocal) speed function; see also [103].

We note also that there exist numerical integrators for SDEs which involve ‘adaptive time-stepping’ routines, where the time-discretisation parameter  $h$  is effectively taken to depend on the location of the process; see e.g. [56] for a nice recent exposition. While these are not based on time-changing per se – indeed, the ultimate goal of such methods is to make a strong approximation of the original SDE – there is at the very least some conceptual common ground in their construction. Conversely, we observe that such methods

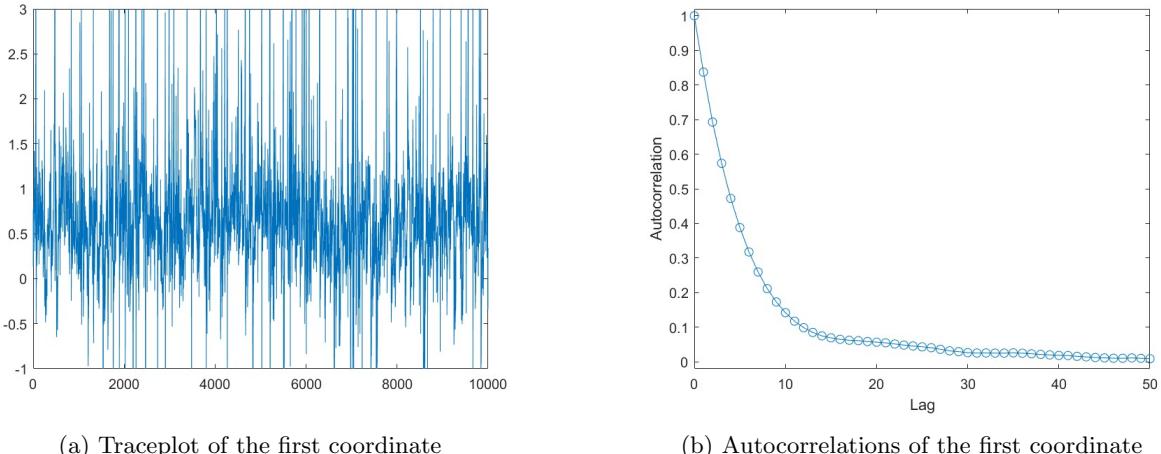


Figure 29: Time-transformed Zig-Zag algorithm on a 5-dimensional Cauchy-Regression target. Speed function:  $s(x) = (1 + x^2)^{\frac{1}{2}}$ .

appear to have garnered most interest in the context of rough targets, with the application to flat targets an apparently more recent development.

### 4.3 Open Problems

We again pause to discuss some problems which are not entirely addressed in the present literature, but would be of some interest to resolve.

#### Heavy Tails in High Dimension

High dimensions are understood to make generic sampling tasks more challenging, and this is particularly true in the case of heavy-tailed problems. One heuristic in this direction is that in high dimension, there are many more ways in which to be heavy-tailed (or not). In particular, heterogeneity of tail-heaviness across directions is a real challenge; see the literature on variational inference for some relevant discussion [81, 91, 70]. There are also qualitative differences between spherically-symmetric tails and product-form tails which are well-understood in some contexts, but are rather under-discussed in the sampling literature at present. The language of functional inequalities and isoperimetry gives some guidance in this direction (see e.g. [25, 64]), but does not really tackle the practical question of detecting ‘what the problem is’ in a specific instance. Along similar lines, it appears important to define a richer class of model problems which are reflective of practical challenges, but also tractable enough that their analysis might yield some actionable insight.

#### Heavy Tails with Unknown Heaviness

Many of these problems are reasonably solvable when we know precisely how heavy the tails are. In a practical setting, this is hard to do without solving the sampling problem well. It would be useful to find adaptive strategies for introducing the tails more gradually, and resolving problems as they appear. Techniques for tail index estimation from the extremes literature (e.g. the classical Hill estimator [73] and its extensions) could potentially be useful tools in either case, having similarly been applied successfully in the context of density estimation [86].

## **Reliable Transformations for Negating Heavy Tails**

Effective transformations of space or time are difficult to come up with in the absence of strong a priori information. Existing works have proposed the use of radial transformations, and it would also be rather natural to propose coordinate-wise transformations. These transformations are plausibly learnable in practice, but also flexible enough to accommodate an interesting range of target behaviours. It would be of great value to identify other families of structured target distributions for which ‘good’ transformations naturally suggest themselves.

## **Negating Heavy Tails without Inducing Roughness**

While confining the target more can resolve the issues with tail heaviness, a careless approach to confinement can easily exacerbate issues with roughness. It seems to be widely-known at an informal level that ‘one must be careful’ in such scenarios, but it appears difficult to automate this intuition. Even for the relatively safe and non-intrusive strategy of applying a transformation which coincides with the identity mapping on a large ball of radius  $R \gg 0$ , and then compresses mass much more tightly outside of this ball, there do not appear to be default solutions in the literature for setting this critical radius  $R$ . Practical resolutions to these challenges could be of great utility for navigating heavy-tailed problems.

## **5 Conclusion**

While classical MCMC still ‘works’, there is an increasing realisation that the sampling problems arising in contemporary applications often fall outside of the scope of the ‘usual’ assumptions which we have in mind when developing methods and theory. There has been substantial progress on deriving algorithms which achieve a degree of robustness with respect to some of these pathologies, while remaining performant in ‘benign’ settings. We have reviewed a number of existing solutions and identified some areas for potential contributions and growth.

Our focus has been on sampling problems for which ‘size’ is somehow the primary issue, whether through the richness of information about the geometry of the target distribution for rough targets, or the sparsity of information provided by flat targets. This has allowed for a relatively clean treatment, whereby careful inflation or attenuation of this information can lead to more stable algorithms. Nevertheless, despite this relative simplicity, various challenges remain, largely relating to the difficulty of estimating these features ‘on-the-fly’ as part of the sampling process, rather than by a priori mathematical computations. We see this as an interesting avenue for future contributions.

In any case, it bears mentioning that there are many other sampling problems of practical interest for which ‘size’ is a relatively minor concern, and ‘shape’ is entirely crucial. One instance of this is the task of multi-modal sampling, for which the literature on multi-canonical or ‘tempering / annealing’ strategies is vast. Another slightly more vague instance would be the phenomenon of target distributions with ‘heterogeneous local geometry’, for which different parts of the state space require drastically different ‘shapes’ of dynamics to explore efficiently; rigorous progress on this topic is still rather nascent. We certainly do not suggest that the developments detailed in this survey are designed with such (arguably more difficult) challenges in mind, though one hopes that they will at least be compatible with them.

The past decade or so has been immensely fruitful from the perspective of establishing theoretical guarantees for MCMC in the ‘well-conditioned, log-concave’ setting, with the quantitative analysis of various algorithms coming along in leaps and bounds during this period. In search of similar foundational guarantees for methods which apply to a broader class of sampling problems, we are optimistic that this burgeoning literature on ‘robust’ MCMC will provide ample inspiration, leading to improved practical solutions with a firm theoretical basis. We are excited to see how things progress from here.

## References

- [1] David F Anderson and Thomas G Kurtz. *Stochastic analysis of biochemical systems*, volume 674. Springer, 2015.
- [2] Charly Andral, Randal Douc, Hugo Marival, and Christian P. Robert. The importance Markov chain. *Stochastic Processes and their Applications*, 171:104316, 2024.
- [3] Charly Andral and Kengo Kamatani. Automated techniques for efficient sampling of Piecewise-Deterministic Markov Processes. *arXiv preprint arXiv:2408.03682*, 2024.
- [4] Christophe Andrieu, Anthony Lee, and Sam Livingstone. A general perspective on the Metropolis-Hastings kernel. *arXiv preprint arXiv:2012.14881*, 2020.
- [5] Christophe Andrieu, Anthony Lee, Sam Power, and Andi Q Wang. Explicit convergence bounds for Metropolis Markov chains: Isoperimetry, spectral gaps and profiles. *The Annals of Applied Probability*, 34(4):4022–4071, 2024.
- [6] Christophe Andrieu, Anthony Lee, Sam Power, and Andi Q Wang. Weak Poincaré inequalities for Markov chains: theory and applications. *The Annals of Applied Probability*, to appear.
- [7] Filippo Ascolani, Hugo Lavenant, and Giacomo Zanella. Entropy contraction of the Gibbs sampler under log-concavity. *arXiv preprint arXiv:2410.00858*, 2024.
- [8] Yves F Atchadé. An adaptive version for the Metropolis-adjusted Langevin algorithm with a truncated drift. *Methodology and Computing in Applied Probability*, 8(2):235–254, 2006.
- [9] Anthony Alfred Barker. Monte Carlo calculations of the radial distribution functions for a proton-electron plasma. *Australian Journal of Physics*, 18(2):119–134, 1965.
- [10] Claude JP Bélisle, H Edwin Romeijn, and Robert L Smith. Hit-and-Run algorithms for generating multivariate distributions. *Mathematics of Operations Research*, 18(2):255–266, 1993.
- [11] Cameron Bell, Krzysztof Łatuszyński, and Gareth O. Roberts. Adaptive Stereographic MCMC, 2025.
- [12] Matej Benko, Iwona Chlebicka, Jørgen Endal, and Błażej Miasojedow. Langevin Monte Carlo beyond Lipschitz Gradient Continuity. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 15541–15549, 2025.
- [13] Etienne P Bernard, Werner Krauth, and David B Wilson. Event-chain Monte Carlo algorithms for hard-sphere systems. *Physical Review E—Statistical, Nonlinear, and Soft Matter Physics*, 80(5):056704, 2009.
- [14] Andrea Bertazzi, Joris Bierkens, and Paul Dobson. Approximations of Piecewise-Deterministic Markov Processes and their convergence properties. *Stochastic Processes and their Applications*, 154:91–153, 2022.
- [15] Andrea Bertazzi and Giorgos Vasdekis. Sampling with time-changed Markov processes, 2025.
- [16] Julian Besag. Comments on “Representations of knowledge in complex systems” by U Grenander and MI Miller. *J. Roy. Statist. Soc. Ser. B*, 56(591-592):4, 1994.
- [17] Alexandros Beskos, Omiros Papaspiliopoulos, and Gareth O Roberts. Retrospective exact simulation of diffusion sample paths with applications. *Bernoulli*, 12(6):1077–1098, 2006.
- [18] Alexandros Beskos, Omiros Papaspiliopoulos, Gareth O Roberts, and Paul Fearnhead. Exact and computationally efficient likelihood-based estimation for discretely observed diffusion processes (with discussion). *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 68(3):333–382, 2006.

- [19] Alexandros Beskos and Gareth O. Roberts. Exact simulation of diffusions. *The Annals of Applied Probability*, 15(4):2422 – 2444, 2005.
- [20] Michael Betancourt. A conceptual introduction to Hamiltonian Monte Carlo, 2018.
- [21] Joris Bierkens and Andrew Duncan. Limit theorems for the Zig-Zag Process. *Advances in Applied Probability*, 49(3):791–825, 2017.
- [22] Joris Bierkens, Paul Fearnhead, and Gareth Roberts. The Zig-Zag Process and super-efficient sampling for Bayesian analysis of big data. *The Annals of Statistics*, 47(3):1288 – 1320, 2019.
- [23] Joris Bierkens, Kengo Kamatani, and Gareth O Roberts. Scaling of piecewise deterministic Monte Carlo for anisotropic targets. *Bernoulli*, 31(3):2323–2350, 2025.
- [24] Jose Blanchet and Fan Zhang. Exact simulation for multivariate Itô diffusions. *Advances in Applied Probability*, 52(4):1003–1034, 2020.
- [25] Sergey Bobkov. Large deviations and isoperimetry over convex probability measures with heavy tails. *Electronic Journal of Probability*, 12:1072–1100, 2007.
- [26] Nawaf Bou-Rabee and Martin Hairer. Nonasymptotic mixing of the MALA algorithm. *IMA Journal of Numerical Analysis*, 33(1):80–110, 2013.
- [27] Nawaf Bou-Rabee and Jesús María Sanz-Serna. Randomized Hamiltonian Monte Carlo. *The Annals of Applied Probability*, 27(4):2159–2194, 2017.
- [28] Nawaf Bou-Rabee and Eric Vanden-Eijnden. Pathwise accuracy and ergodicity of Metropolized integrators for SDEs. *Communications on Pure and Applied Mathematics: A Journal Issued by the Courant Institute of Mathematical Sciences*, 63(5):655–696, 2010.
- [29] Nawaf Bou-Rabee and Eric Vanden-Eijnden. *Continuous-Time Random Walks for the Numerical Solution of Stochastic Differential Equations*, volume 256, no. 1228 of *Memoirs of the American Mathematical Society*. American Mathematical Society, Providence, RI, 2018.
- [30] Alexandre Bouchard-Côté, Sebastian J Vollmer, and Arnaud Doucet. The Bouncy Particle Sampler: A nonreversible, rejection-free Markov chain Monte Carlo method. *Journal of the American Statistical Association*, 113(522):855–867, 2018.
- [31] Stephen P Boyd and Lieven Vandenberghe. *Convex Optimization*. Cambridge university press, 2004.
- [32] Steve Brooks, Andrew Gelman, Galin Jones, and Xiao-Li Meng, editors. *Handbook of Markov Chain Monte Carlo*. Chapman & Hall/CRC, New York, 2011.
- [33] Nicolas Brosse, Alain Durmus, Éric Moulines, and Sotirios Sabanis. The tamed unadjusted Langevin algorithm. *Stochastic Processes and their Applications*, 129(10):3638–3663, 2019.
- [34] Sébastien Bubeck. Convex optimization: Algorithms and complexity. *Foundations and Trends® in Machine Learning*, 8(3-4):231–357, 2015.
- [35] Carlos M Carvalho, Nicholas G Polson, and James G Scott. Handling sparsity via the horseshoe. In *Artificial intelligence and statistics*, pages 73–80. PMLR, 2009.
- [36] Lotfi Chaari, Jean-Yves Tourneret, Caroline Chaux, and Hadj Batatia. A Hamiltonian Monte Carlo method for non-smooth energy sampling. *IEEE Transactions on Signal Processing*, 64(21):5585–5594, 2016.
- [37] Xiang Cheng, Niladri S Chatterji, Peter L Bartlett, and Michael I Jordan. Underdamped Langevin MCMC: A non-asymptotic analysis. In *Conference on learning theory*, pages 300–323. PMLR, 2018.
- [38] Augustin Chevallier, Sam Power, Andi Q Wang, and Paul Fearnhead. PDMP Monte Carlo methods for piecewise-smooth densities. *Advances in Applied Probability*, 56(4):1153–1194, 2024.

- [39] Sinho Chewi, Murat A Erdogdu, Mufan Li, Ruoqi Shen, and Matthew S Zhang. Analysis of Langevin Monte Carlo: From Poincaré to Log-Sobolev. *Foundations of Computational Mathematics*, pages 1–51, 2024.
- [40] Sinho Chewi, Chen Lu, Kwangjun Ahn, Xiang Cheng, Thibaut Le Gouic, and Philippe Rigollet. Optimal dimension dependence of the Metropolis-adjusted Langevin algorithm. In *Conference on Learning Theory*, pages 1260–1300. PMLR, 2021.
- [41] Alice Corbella, Simon EF Spencer, and Gareth O Roberts. Automatic Zig-Zag sampling in practice. *Statistics and Computing*, 32(6):107, 2022.
- [42] Francesca R Crucinio, Alain Durmus, Pablo Jiménez, and Gareth O Roberts. Optimal scaling results for Moreau-Yosida Metropolis-adjusted Langevin algorithms. *Bernoulli*, 31(3):1889–1907, 2025.
- [43] Arnak S Dalalyan. Theoretical guarantees for approximate sampling from smooth and log-concave densities. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 79(3):651–676, 2017.
- [44] Mark HA Davis. Piecewise-Deterministic Markov processes: A general class of non-diffusion stochastic models. *Journal of the Royal Statistical Society: Series B (Methodological)*, 46(3):353–376, 1984.
- [45] George Deligiannidis, Daniel Paulin, Alexandre Bouchard-Côté, and Arnaud Doucet. Randomized Hamiltonian Monte Carlo as scaling limit of the Bouncy Particle Sampler and dimension-free convergence rates. *The Annals of Applied Probability*, 31(6):2612–2662, 2021.
- [46] Samuel Duffield, Maxwell Aifer, Denis Melanson, Zach Belateche, and Patrick J Coles. Lattice Random Walk Discretisations of Stochastic Differential Equations. *arXiv preprint arXiv:2508.20883*, 2025.
- [47] Alain Durmus and Andreas Eberle. Asymptotic bias of inexact Markov chain Monte Carlo methods in high dimension. *The Annals of Applied Probability*, 34(4):3435–3468, 2024.
- [48] Alain Durmus, Arnaud Guillin, and Pierre Monmarché. Geometric ergodicity of the Bouncy Particle Sampler. *The Annals of Applied Probability*, 30(5):2069–2098, 2020.
- [49] Alain Durmus and Éric Moulines. Nonsymptotic convergence analysis for the unadjusted Langevin algorithm. *The Annals of Applied Probability*, 27(3):1551 – 1587, 2017.
- [50] Alain Durmus and Éric Moulines. High-dimensional Bayesian inference via the unadjusted Langevin algorithm. *Bernoulli*, 25(4A):2854 – 2882, 2019.
- [51] Weinan E, Tiejun Li, and Eric Vanden-Eijnden. *Applied Stochastic Analysis*, volume 199. American Mathematical Soc., 2021.
- [52] Paul Embrechts, Claudia Klüppelberg, and Thomas Mikosch. *Modelling Extremal Events: for Insurance and Finance*, volume 33 of *Stochastic Modelling and Applied Probability*. Springer-Verlag, Berlin, 1st edition, 1997.
- [53] Stewart N Ethier and Thomas G Kurtz. *Markov processes: characterization and convergence*, volume 282. John Wiley & Sons, 2009.
- [54] Paul Fearnhead, Sebastiano Grazzi, Chris Nemeth, and Gareth O Roberts. Stochastic Gradient Piecewise Deterministic Monte Carlo Samplers. *arXiv preprint arXiv:2406.19051*, 2024.
- [55] Paul Fearnhead, Omilos Papaspiliopoulos, and Gareth O Roberts. Particle filters for partially-observed diffusions. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 70(4):755–777, 2008.
- [56] James Foster and Andraž Jelinčič. On the convergence of adaptive approximations for stochastic differential equations. *arXiv preprint arXiv:2311.14201*, 2023.

- [57] Alan E Gelfand. Gibbs sampling. *Journal of the American statistical Association*, 95(452):1300–1304, 2000.
- [58] Saul Brian Gelfand and Sanjoy K Mitter. Weak convergence of Markov chain sampling methods and annealing algorithms to diffusions. *Journal of Optimization Theory and Applications*, 68(3):483–498, 1991.
- [59] Andrew Gelman, Walter R Gilks, and Gareth O Roberts. Weak convergence and optimal scaling of random walk Metropolis algorithms. *The annals of applied probability*, 7(1):110–120, 1997.
- [60] Michael B Giles. Multilevel Monte Carlo methods. *Acta numerica*, 24:259–328, 2015.
- [61] Daniel T Gillespie. Exact stochastic simulation of coupled chemical reactions. *The journal of physical chemistry*, 81(25):2340–2361, 1977.
- [62] Daniel T Gillespie. Approximate accelerated stochastic simulation of chemically reacting systems. *The Journal of chemical physics*, 115(4):1716–1733, 2001.
- [63] Jacob Vorstrup Goldman, Torben Sell, and Sumeetpal Sidhu Singh. Gradient-based Markov chain Monte Carlo for Bayesian inference with non-differentiable priors. *Journal of the American Statistical Association*, 117(540):2182–2193, 2022.
- [64] Nathael Gozlan, Cyril Roberto, and Paul-Marie Samson. Isoperimetry for product of heavy tails distributions. *Progress in analysis and its applications*, pages 470–478, 2010.
- [65] Ernst Hairer, Christian Lubich, and Gerhard Wanner. Geometric numerical integration illustrated by the Störmer–Verlet method. *Acta numerica*, 12:399–450, 2003.
- [66] Luke Hardcastle, Samuel Livingstone, and Gianluca Baio. Diffusion Piecewise-Exponential models for survival extrapolation using Piecewise-Deterministic Monte Carlo. *arXiv preprint arXiv:2505.05932*, 2025.
- [67] Julian Harland, Manon Michel, Tobias A Kampmann, and Jan Kierfeld. Event-chain Monte Carlo algorithms for three-and many-particle interactions. *Europhysics Letters*, 117(3):30001, 2017.
- [68] W. K. Hastings. Monte Carlo Sampling Methods Using Markov Chains and Their Applications. *Biometrika*, 57(1):97–109, 1970.
- [69] Ye He, Krishnakumar Balasubramanian, and Murat A Erdogdu. An analysis of transformed unadjusted Langevin algorithm for heavy-tailed sampling. *IEEE Transactions on Information Theory*, 70(1):571–593, 2023.
- [70] Tennessee Hickling and Dennis Prangle. Flexible tails for normalizing flows. In Aarti Singh, Maryam Fazel, Daniel Hsu, Simon Lacoste-Julien, Felix Berkenkamp, Tegan Maharaj, Kiri Wagstaff, and Jerry Zhu, editors, *Proceedings of the 42nd International Conference on Machine Learning*, volume 267 of *Proceedings of Machine Learning Research*, pages 23155–23178. PMLR, 13–19 Jul 2025.
- [71] Desmond Higham and Peter Kloeden. *An introduction to the numerical simulation of stochastic differential equations*. SIAM, 2021.
- [72] Desmond J Higham. An algorithmic introduction to numerical simulation of stochastic differential equations. *SIAM review*, 43(3):525–546, 2001.
- [73] Bruce M Hill. A simple general approach to inference about the tail of a distribution. *The annals of statistics*, pages 1163–1174, 1975.
- [74] Max Hird and Samuel Livingstone. Quantifying the Effectiveness of Linear Preconditioning in Markov Chain Monte Carlo. *Journal of Machine Learning Research*, 26(119):1–51, 2025.

- [75] Max Hird, Samuel Livingstone, and Giacomo Zanella. A fresh Take on ‘Barker Dynamics’ for MCMC. In Alexander Keller, editor, *Monte Carlo and Quasi-Monte Carlo Methods*, pages 169–184, Cham, 2022. Springer International Publishing.
- [76] Jean-Baptiste Hiriart-Urruty and Claude Lemaréchal. *Convex Analysis and Minimization Algorithms II: Advanced Theory and Bundle Methods*, volume 306 of *Grundlehren der mathematischen Wissenschaften*. Springer-Verlag, Berlin, Heidelberg, 1993.
- [77] Liam Hodgkinson, Robert Salomone, and Fred Roosta. Implicit Langevin algorithms for sampling from log-concave densities. *Journal of Machine Learning Research*, 22(136):1–30, 2021.
- [78] Matthew Hoffman, Pavel Sountsov, Joshua V. Dillon, Ian Langmore, Dustin Tran, and Srinivas Va-sudevan. NeuTra-lizing Bad Geometry in Hamiltonian Monte Carlo using Neural Transport, 2019.
- [79] Martin Hairer, Arnulf Jentzen, and Peter E Kloeden. Strong and weak divergence in finite time of Euler’s method for stochastic differential equations with non-globally Lipschitz continuous coefficients. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 467(2130):1563–1576, 2011.
- [80] Martin Hairer, Arnulf Jentzen, and Peter E Kloeden. Strong convergence of an explicit numerical method for SDEs with nonglobally Lipschitz continuous coefficients. *Annals of Applied Probability*, 2012.
- [81] Priyank Jaini, Ivan Kobyzev, Yaoliang Yu, and Marcus Brubaker. Tails of Lipschitz triangular flows. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 4673–4681. PMLR, 13–18 Jul 2020.
- [82] Leif T Johnson and Charles J Geyer. Variable transformation to obtain geometric ergodicity in the random-walk Metropolis algorithm. *The Annals of Statistics*, pages 3050–3076, 2012.
- [83] Tim Johnston and Sotirios Sabanis. A strongly monotonic polygonal euler scheme. *Journal of Complexity*, 80:101801, 2024.
- [84] Werner Krauth. Event-chain Monte Carlo: Foundations, applications, and prospects. *Frontiers in Physics*, 9:663457, 2021.
- [85] Kenneth L. Lange, Roderick J. A. Little, and Jeremy M. G. Taylor. Robust Statistical Modeling Using the t Distribution. *Journal of the American Statistical Association*, 84(408):881–896, 1989.
- [86] Mike Laszkiewicz, Johannes Lederer, and Asja Fischer. Marginal Tail-Adaptive Normalizing Flows. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 12020–12048. PMLR, 17–23 Jul 2022.
- [87] John M. Lee. *Introduction to Smooth Manifolds*, volume 218 of *Graduate Texts in Mathematics*. Springer, New York, NY, 2nd edition, 2013.
- [88] Benedict Leimkuhler, René Lohmann, and Peter Whalley. A Langevin sampling algorithm inspired by the Adam optimizer, 2025.
- [89] Alix Leroy, Benedict Leimkuhler, Jonas Latz, and Desmond J Higham. Adaptive stepsize algorithms for Langevin dynamics. *SIAM Journal on Scientific Computing*, 46(6):A3574–A3598, 2024.
- [90] Peter A.W. Lewis and Gerald S. Shedler. Simulation methods for poisson processes in nonstationary systems. In *Proceedings of the 10th Conference on Winter Simulation - Volume 1*, WSC ’78, page 155–163. IEEE Press, 1978.

- [91] Feynman Liang, Michael Mahoney, and Liam Hodgkinson. Fat-Tailed Variational Inference with Anisotropic Tail Adaptive Flows. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 13257–13270. PMLR, 17–23 Jul 2022.
- [92] Eckhard Limpert, Werner A. Stahel, and Markus Abbt. Log-normal Distributions across the sciences: Keys and Clues. *BioScience*, 51(5):341–352, 05 2001.
- [93] Samuel Livingstone, Michael F Faulkner, and Gareth O Roberts. Kinetic energy choice in Hamiltonian/Hybrid Monte Carlo. *Biometrika*, 106(2):303–319, 2019.
- [94] Samuel Livingstone, Nikolas Nüsken, Giorgos Vasdekis, and Rui-Yang Zhang. Skew-symmetric schemes for stochastic differential equations with non-Lipschitz drift: an unadjusted Barker algorithm, 2024.
- [95] Samuel Livingstone and Giacomo Zanella. The Barker proposal: Combining robustness and efficiency in gradient-based MCMC. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 84(2):496–523, 2022.
- [96] Jianfeng Lu and Lihan Wang. Complexity of zigzag sampling algorithm for strongly log-concave distributions. *Statistics and Computing*, 32(3):48, 2022.
- [97] Jianfeng Lu and Lihan Wang. On explicit  $L^2$ -convergence rate estimate for piecewise deterministic Markov processes in MCMC algorithms. *The Annals of Applied Probability*, 32(2):1333–1361, 2022.
- [98] Xiaoyu Lu, Valerio Perrone, Leonard Hasenclever, Yee Whye Teh, and Sebastian Vollmer. Relativistic Monte Carlo. In *Artificial Intelligence and Statistics*, pages 1236–1245. PMLR, 2017.
- [99] Iosif Lytras and Panayotis Mertikopoulos. Tamed Langevin sampling under weaker conditions. *arXiv preprint arXiv:2405.17693*, 2024.
- [100] Iosif Lytras and Sotirios Sabanis. Taming under isoperimetry. *Stochastic Processes and their Applications*, page 104684, 2025.
- [101] Iosif Lytras, Sotirios Sabanis, and Ying Zhang. kTULA: A Langevin sampling algorithm with improved KL bounds under super-linear log-gradients. *arXiv preprint arXiv:2506.04878*, 2025.
- [102] AC Maggs and Werner Krauth. Large-scale dynamics of event-chain Monte Carlo. *Physical Review E*, 105(1):015309, 2022.
- [103] Sonia Malefaki and George Iliopoulos. On convergence of properly weighted samples to the target distribution. *Journal of Statistical Planning and Inference*, 138(4):1210–1225, 2008.
- [104] Lorenzo Mauri and Giacomo Zanella. Robust approximate sampling via stochastic gradient Barker dynamics. In *International Conference on Artificial Intelligence and Statistics*, pages 2107–2115. PMLR, 2024.
- [105] Nicholas Metropolis, Arianna W Rosenbluth, Marshall N Rosenbluth, Augusta H Teller, and Edward Teller. Equation of state calculations by fast computing machines. *The journal of chemical physics*, 21(6):1087–1092, 1953.
- [106] Manon Michel, Sebastian C Kapfer, and Werner Krauth. Generalized event-chain Monte Carlo: Constructing rejection-free global-balance algorithms from infinitesimal steps. *The Journal of chemical physics*, 140(5), 2014.
- [107] Grigori N Milstein and Michael V Tretyakov. *Stochastic numerics for mathematical physics*, volume 39. Springer, 2004.
- [108] Grigori N Milstein and Michael V Tretyakov. Numerical integration of stochastic differential equations with nonglobally Lipschitz coefficients. *SIAM journal on numerical analysis*, 43(3):1139–1154, 2005.

- [109] Viacheslav Natarovskii, Daniel Rudolf, and Björn Sprungk. Quantitative spectral gap estimate and Wasserstein contraction of simple slice sampling. *The Annals of Applied Probability*, 31(2):pp. 806–825, 2021.
- [110] Radford M Neal. Slice sampling. *The annals of statistics*, 31(3):705–767, 2003.
- [111] Radford M Neal. MCMC using Hamiltonian dynamics. *Handbook of Markov Chain Monte Carlo*, 2(11):2, 2011.
- [112] Akihiko Nishimura, Zhenyu Zhang, and Marc A Suchard. Zigzag path connects two Monte Carlo samplers: Hamiltonian counterpart to a Piecewise-Deterministic Markov process. *Journal of the American Statistical Association*, 120(550):1077–1089, 2025.
- [113] Kjartan Kloster Osmundsen, Tore Selland Kleppe, and Roman Liesenfeld. Importance Sampling-Based Transport Map Hamiltonian Monte Carlo for Bayesian Hierarchical Models. *Journal of Computational and Graphical Statistics*, 30(4):906–919, 2021.
- [114] Filippo Pagani, Augustin Chevallier, Sam Power, Thomas House, and Simon Cotter. NuZZ: Numerical Zig-Zag for general models. *Statistics and Computing*, 34(1):61, 2024.
- [115] Omiros Papaspiliopoulos. Monte Carlo probabilistic inference for diffusion processes: A methodological framework. In David Barber, A. Taylan Cemgil, and Silvia Chiappa, editors, *Bayesian Time Series Models*, pages 82–103. Cambridge University Press, Cambridge, 2011.
- [116] Marcelo Pereyra. Proximal Markov chain Monte Carlo algorithms. *Statistics and Computing*, 26:745–760, 2016.
- [117] Marcelo Pereyra, Luis Vargas Mieles, and Konstantinos C Zygalakis. Accelerating proximal Markov chain Monte Carlo by using an explicit stabilized method. *SIAM Journal on Imaging Sciences*, 13(2):905–935, 2020.
- [118] Natesh S. Pillai, Andrew M. Stuart, and Alexandre H. Thiéry. Optimal scaling and diffusion limits for the Langevin algorithm in high dimensions. *The Annals of Applied Probability*, 22(6):2320–2356, 2012.
- [119] Eckhard Platen. An introduction to numerical methods for stochastic differential equations. *Acta numerica*, 8:197–246, 1999.
- [120] Sam Power, Daniel Rudolf, Björn Sprungk, and Andi Q Wang. Weak Poincaré inequality comparisons for ideal and hybrid slice sampling. *arXiv preprint arXiv:2402.13678*, 2024.
- [121] Christian P. Robert and George Casella. *Monte Carlo Statistical Methods*. Springer-Verlag, Berlin, Heidelberg, 2005.
- [122] 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, 01 2002.
- [123] Gareth O. Roberts and Osnat Stramer. Langevin Diffusions and Metropolis-Hastings Algorithms. *Methodology and Computing in Applied Probability*, 2002.
- [124] Gareth O. Roberts and Richard L. Tweedie. Exponential convergence of Langevin distributions and their discrete approximations. *Bernoulli*, 2(4):341 – 363, 1996.
- [125] Peter J Rossky, Jimmie D Doll, and Harold L Friedman. Brownian dynamics as smart Monte Carlo simulation. *The Journal of Chemical Physics*, 69(10):4628–4633, 1978.
- [126] Daniel Rudolf. Hit-and-Run for numerical integration. In *Monte Carlo and Quasi-Monte Carlo Methods 2012*, pages 597–612. Springer, 2013.
- [127] Chris Sherlock, Paul Fearnhead, and Gareth O. Roberts. The Random Walk Metropolis: Linking Theory and Practice Through a Case Study. *Statistical Science*, 25(2):172 – 190, 2010.

- [128] Apratim Shukla, Dootika Vats, and Eric C Chi. MCMC Importance Sampling via Moreau-Yosida Envelopes. *arXiv preprint arXiv:2501.02228*, 2025.
- [129] Apratim Shukla, Dootika Vats, and Eric C Chi. Proximal Hamiltonian Monte Carlo. *arXiv preprint arXiv:2510.22252*, 2025.
- [130] Matthew Sutton and Paul Fearnhead. Concave-convex PDMP-based sampling. *Journal of Computational and Graphical Statistics*, 32(4):1425–1435, 2023.
- [131] Lukasz Szpruch. V-stable tamed Euler schemes. *arXiv preprint arXiv:1310.0785*, 2013.
- [132] Luke Tierney. A note on Metropolis-Hastings kernels for general state spaces. *Annals of applied probability*, pages 1–9, 1998.
- [133] Giorgos Vasdekis and Gareth O Roberts. A note on the polynomial ergodicity of the one-dimensional Zig-Zag process. *Journal of Applied Probability*, 59(3):895–903, 2022.
- [134] Giorgos Vasdekis and Gareth O Roberts. Speed up Zig-Zag. *The Annals of Applied Probability*, 33(6A):4693 – 4746, 2023.
- [135] Jure Vogrinc and Wilfrid S. Kendall. Counterexamples for optimal scaling of Metropolis–Hastings chains with rough target densities. *The Annals of Applied Probability*, 31(2):972 – 1019, 2021.
- [136] Jure Vogrinc, Samuel Livingstone, and Giacomo Zanella. Optimal design of the Barker proposal and other locally balanced Metropolis–Hastings algorithms. *Biometrika*, 110(3):579–595, 2023.
- [137] Andre Wibisono. Proximal Langevin algorithm: Rapid convergence under isoperimetry. *arXiv preprint arXiv:1911.01469*, 2019.
- [138] Keru Wu, Scott Schmidler, and Yuansi Chen. Minimax mixing time of the Metropolis-adjusted Langevin algorithm for log-concave sampling. *Journal of Machine Learning Research*, 23(270):1–63, 2022.
- [139] Jun Yang, Krzysztof Latuszyński, and Gareth O. Roberts. Stereographic Markov chain Monte Carlo. *The Annals of Statistics*, 52(6):2692 – 2713, 2024.
- [140] Weixin Yao, Yan Wei, and Chun Yu. Robust mixture regression using the t-distribution. *Computational Statistics and Data Analysis*, 71:116–127, 2014.
- [141] Giacomo Zanella. Informed proposals for local MCMC in discrete spaces. *Journal of the American Statistical Association*, 115(530):852–865, 2020.