

# Zig-Zag Sampling

A MCMC Game-Changer

Hirofumi Shiba 

Institute of Statistical Mathematics

the University of Tokyo

9/10/2024

# Today's Menu

- 1 The Zig-Zag Sampler: What Is It?
- 2 The Algorithm: How to Use It?
- 3 Proof of Concept: How Good Is It?

# I The Zig-Zag Sampler: V

A continuous-time variant of MCMC algorithm



Trajectory for Zig-Zag Sampler. Please attribute Hirofumi Shiba

Hirofumi Shiba

# 1.1 Keywords: PDMP (1/2)

PDMP (Piecewise Deterministic<sup>1</sup> Markov Process)

- 1. Mostly **deterministic** with the exception of jumps which happens at random times
  - 2. **Continuous-time**, instead of discrete-time
- Plays a **complementary role** to SDEs / Diffusions

Property	PDMP
Exactly simulatable?	✓
Subject to discretization errors?	✗
Driving noise	Poisson

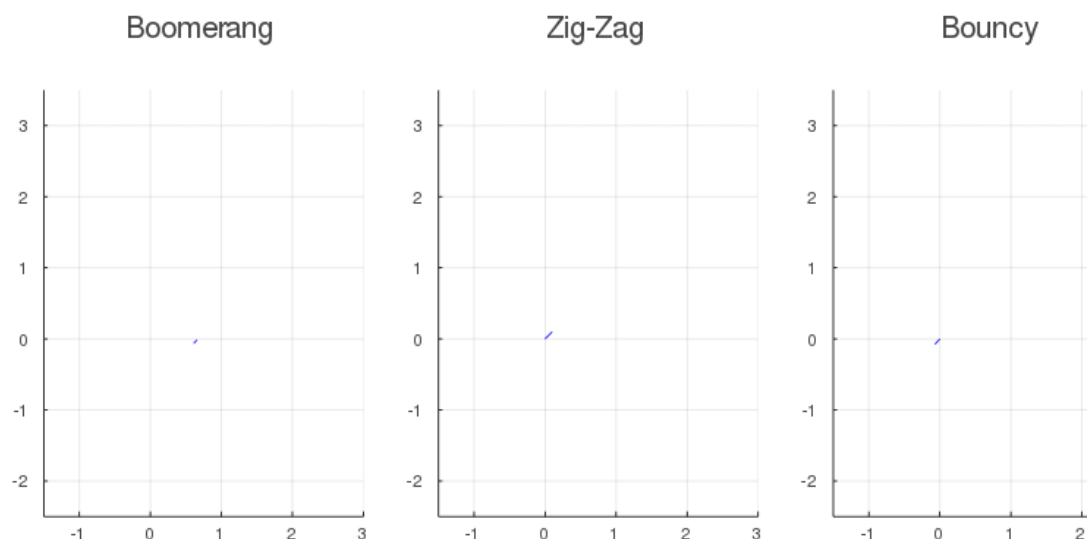
## History of PDMP Applications

1. First applications: control theory, operations research,
2. Second applications: Monte Carlo simulation in materials science (de With, 2012)
3. Third applications: Bayesian statistics (Bouchard-Côté et al., 2015)

1. Mostly **deterministic** with the exception of random jumps
2. **Continuous-time**, instead of discrete-time processes

## I.2 Keywords: PDMP (2/2)

- We will concentrate on Zig-Zag sampler (Fearnhead, et al., 2019)
- Other PDMPs: Bouncy sampler (Bouchard), Boomerang sampler (Bierkens et al., 2020)



The most famous three PDMPs. Animated by (Grazzi, 2020)  
Hirofumi Shiba

## I.3 Menu

### What We've Learned

The new algorithm 'Zig-Zag Sampler' is based on continuous

### What We'll Learn in the Rest of this Section I

We will review 3 instances of the standard (discrete-time) MCMC: **MH**, and **MALA**.

1. Review: **MH** (Metropolis-Hastings) algorithm
2. Review: **Lifted MH**, A method bridging **MH** and Zig-Zag
3. Comparison: **MH** vs. **Lifted MH** vs. Zig-Zag
4. Review: **MALA** (Metropolis Adjusted Langevin Algorithm)
5. Comparison: Zig-Zag vs. **MALA**

## I.4 Review: Metropolis-Hastings (I/2)

(Metropolis et al., 1953)-(Hastings, 1970)

Input: Target distribution  $p$ , (symmetric) proposal distribution

1. Draw a  $X_t \sim q(-|X_{t-1})$

2. Compute

$$\alpha(X_{t-1}, X_t) = \frac{p(X_t)}{p(X_{t-1})}$$

3. Draw a uniform random number  $U \sim U([0, 1])$ .

4. If  $\alpha(X_{t-1}, X_t) \leq U$ , then  $X_t \leftarrow X_{t-1}$ . Do nothing otherwise

5. Return to Step 1.

MH algorithm works even without  $p$ 's normalizing constant. Ho



## I.5 Review: Metropolis-Hastings (2/2)

Alternative View: MH is a generic procedure to turn a into a **Markov chain converging to  $p$** .

### The Choice of Proposal $q$

- Random Walk Metropolis (Metropolis et al., 1953): Uniform

$$q(y|x) = q(y - x) \in \left\{ \frac{dU([0, 1])}{d\lambda}(y - x), \frac{dN}{d\lambda}(y - x) \right\}$$

- Hybrid / Hamiltonian Monte Carlo (Duane et al., 1987): H

$$q(y|x) = \delta_{x+\epsilon\rho}, \quad \epsilon > 0, \quad \rho : \text{momentum defined}$$

- Metropolis-adjusted Langevin algorithm (MALA) (Besag, 1994)

$$q(-|X_t) := \text{the transition probability of } X_t \text{ where } dX_t =$$



## I.6 Problem: **Reversibility**

**Reversibility** (a.k.a detailed balance):

$$p(x)q(x|y) = p(y)q(y|x)$$

In words:

Probability[Going  $x \rightarrow y$ ] = Probability

→ Harder to explore the entire space

→ Slow mixing of **MH**

From the beginning of 21th century, many efforts have been made

## 1.7 Lifting (1/3)

**Lifting**: A method to make MH's dynamics **irreversible**

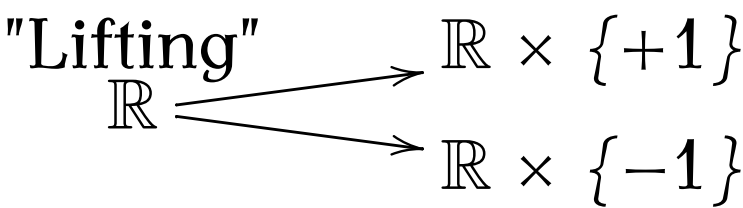
How?: By adding an auxiliary variable  $\sigma \in \{\pm 1\}$ , called

### **Lifted MH** (Turitsyn et al., 2011)

Input: Target  $p$ , **two** proposals  $q^{(+1)}, q^{(-1)}$ , and **momentum**  $\sigma$

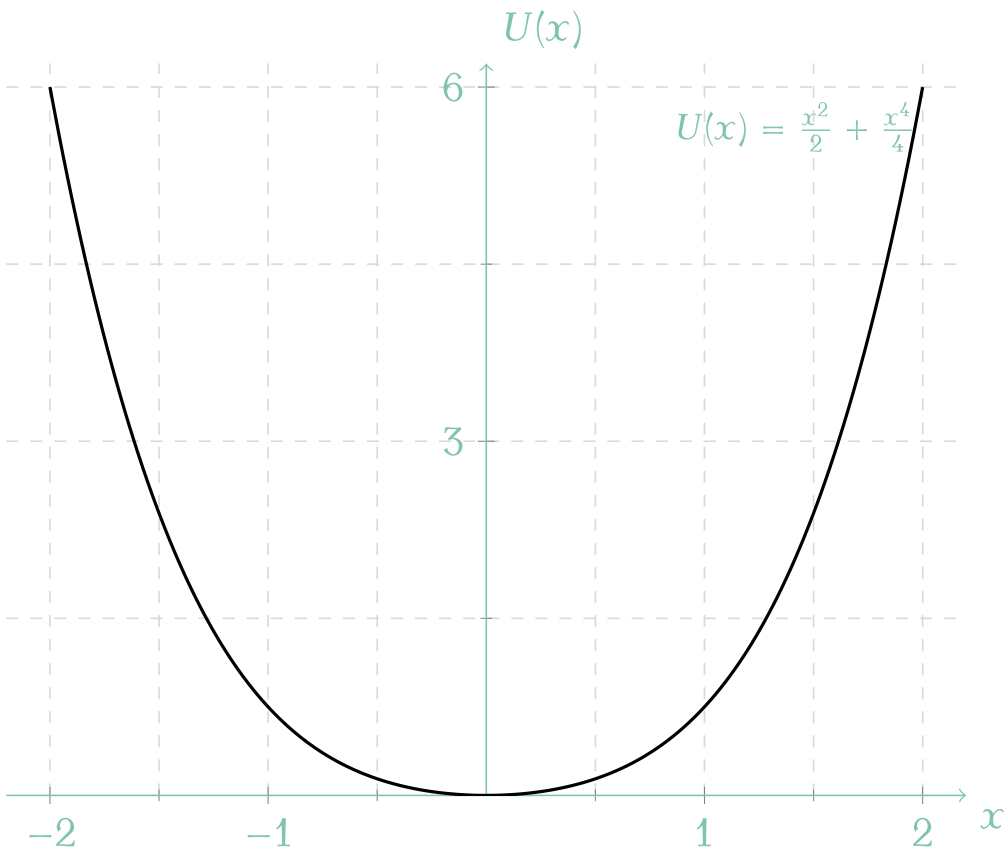
1. Draw  $X_t$  from  $q^{(\sigma)}$
2. Do a MH step
3. If accepted, go back to Step 1.
4. If rejected, **flip the momentum** and go back to Step 1.

I.8 Lifting (2/3)



$q^{(+1)}$ : Only

$q^{(-1)}$ : Only



→ Once g  
continues

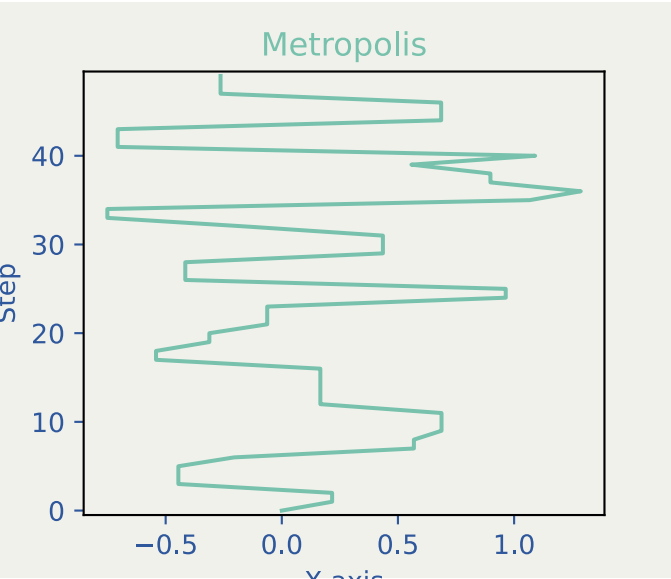
→ This is i

Probabil

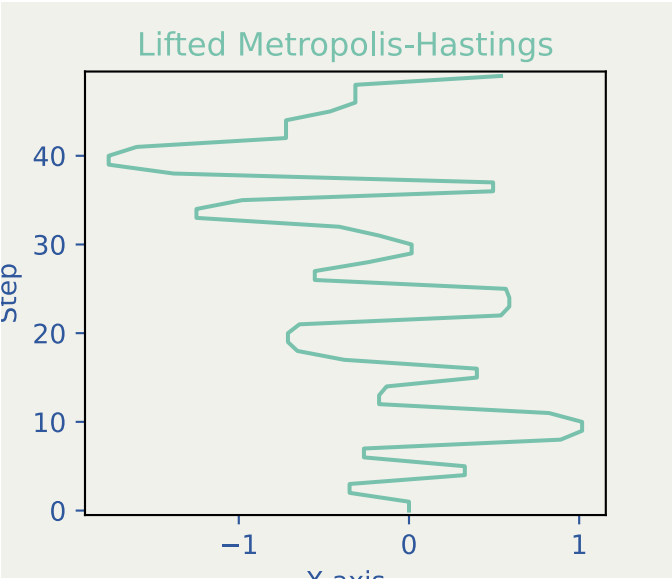
$\neq \mathbb{P}$

## 1.9 Lifting (3/3)

Reversible dynamic of MH has ‘irreversified’



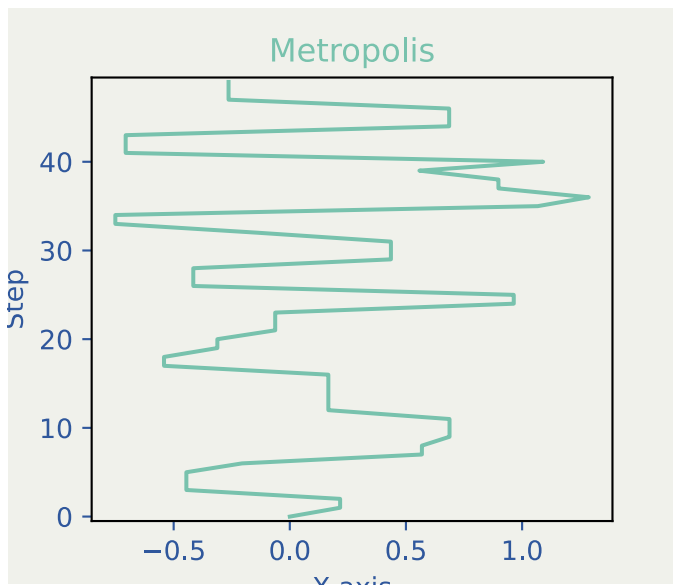
MH



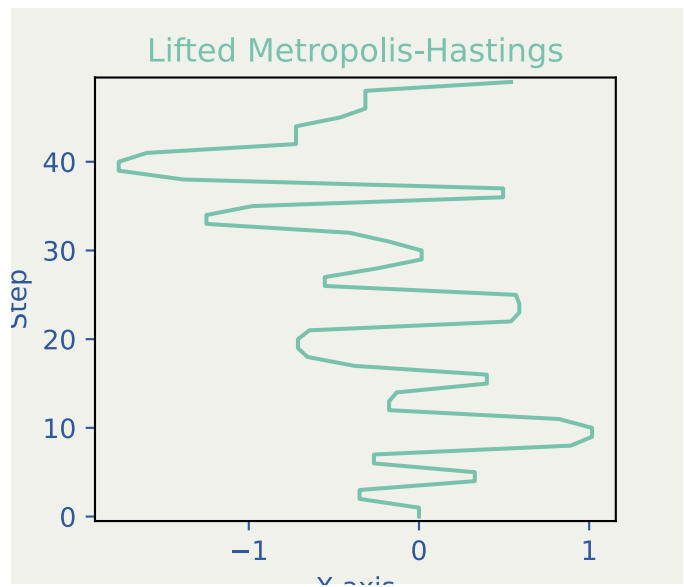
Lifted MH

\*Irreversibility actually improves the efficiency of MCMC, as w

## I.10 Comparison: MH vs. LMH vs. Zig-



MH



Lifted MH

Zig-Zag corresponds to the **limiting case** of size of proposal  $q$  goes to zero, as we'll learn

→ Zig-Zag has a maximum **irreversibility**.

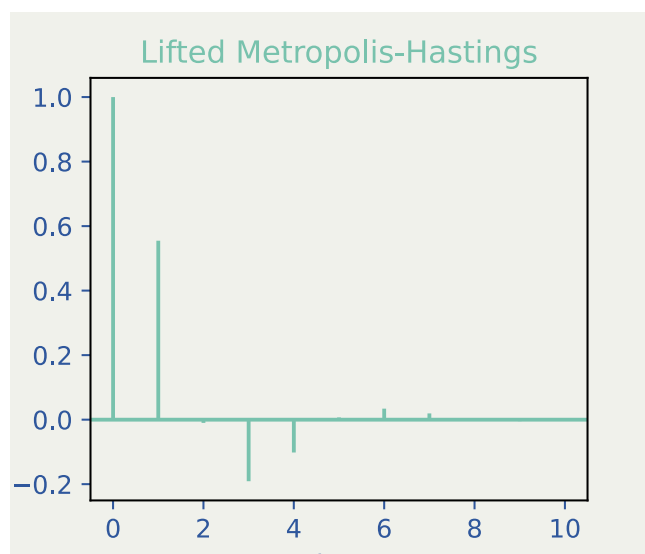
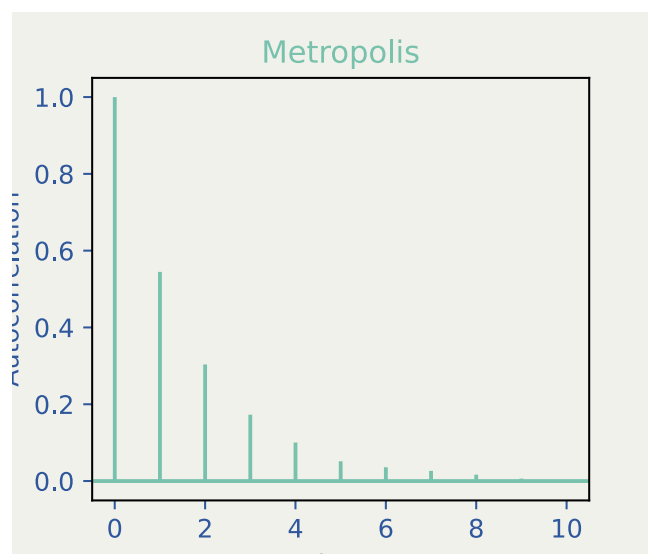
# 1.1 Comparison: MH vs. LMH vs. ZigZag

**Irreversibility** actually improves the efficiency

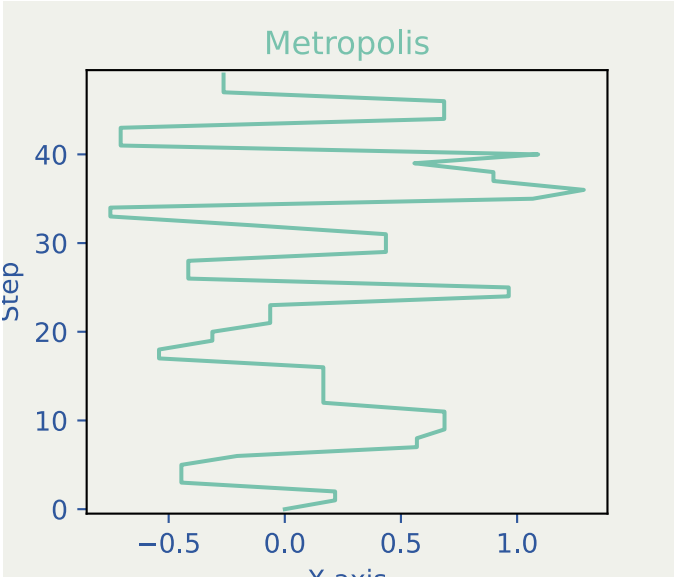
Faster decay of **autocorrelation**  $\rho_t \approx C e^{-\lambda t}$

1. faster mixing of MCMC

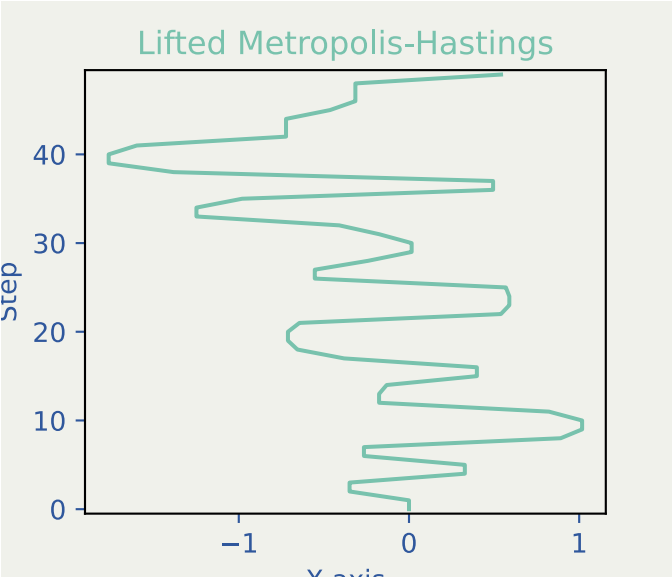
2. lower variance of Monte Carlo estimates







MH



Lifted MH

## 1.12 Review: MALA

**Langevin diffusion:** A diffusion process defined by the following stochastic differential equation:

$$dX_t = \nabla \log p(X_t) dt + \sqrt{2\beta^{-1}} dB_t.$$

**Langevin diffusion** itself converges to the target distribution  $p$ :

$$\|p_t - p\|_{L^1} \rightarrow 0, \quad t \rightarrow \infty.$$

Two MCMC algorithms derived from **Langevin diffusion**:

ULA (Unadjusted Langevin Algorithm)

Use the discretization of  $(X_t)$ . **Discretization errors** are introduced.

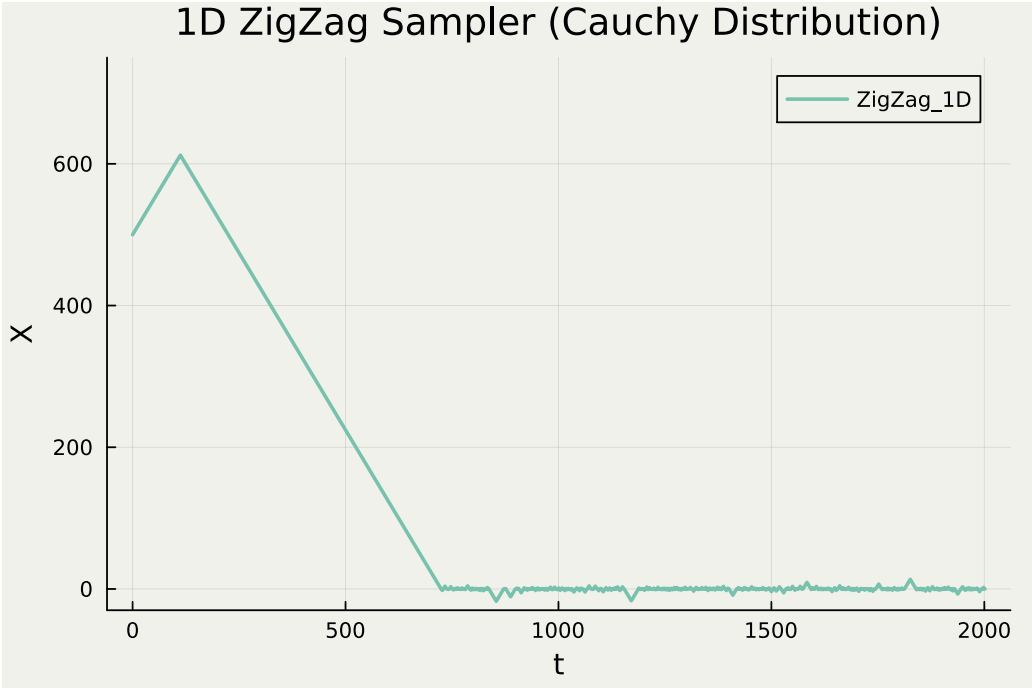
MALA (Metropolis Adjusted Langevin Algorithm)

Use ULA as a proposal in MH, erasing the errors by the Metropolis adjustment.

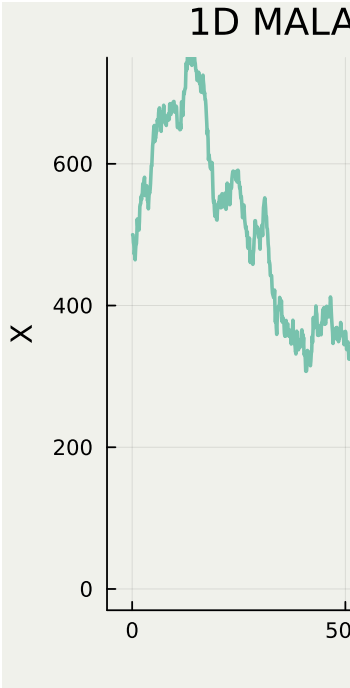
1. under fairly general conditions on  $p$ .

# 1.13 Comparison: Zig-Zag vs. MALA (

How fast do they go back to high-probability



Zig-Zag

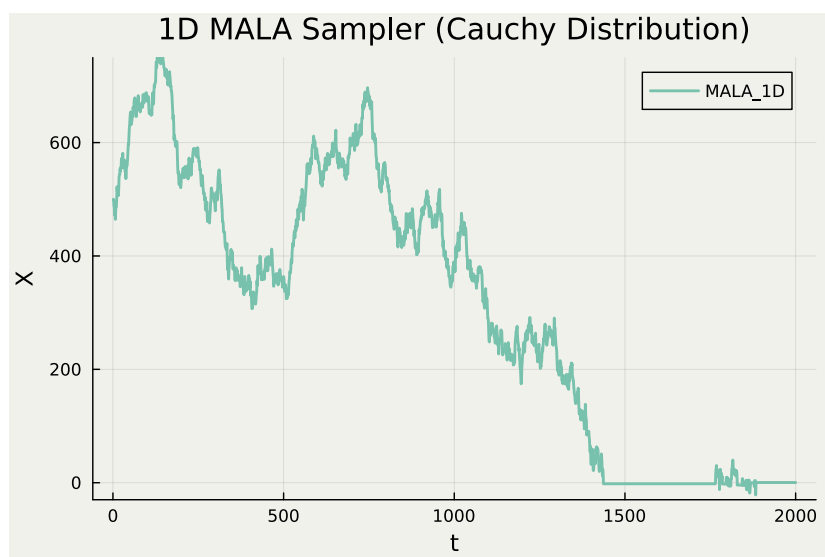


MALA

Irreversibility of Zig-Zag accelerates its conv

I. The target here is the standard Cauchy distribution  $C(0, 1)$ , distribution. Its heavy tails hinder the convergence of MCMC

# 1.14 Comparison: Zig-Zag vs. MALA (2)



☐ **Caution: Fake C**

The left plot looks  
**actually is not.**

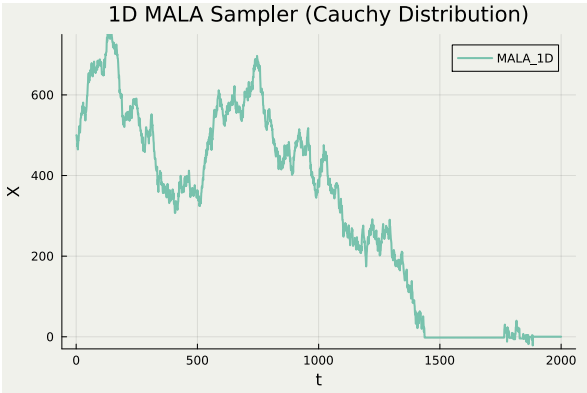
MALA trajectory

MH, including MALA, is actually a discrete-t

The plot is obtained by connecting the poin

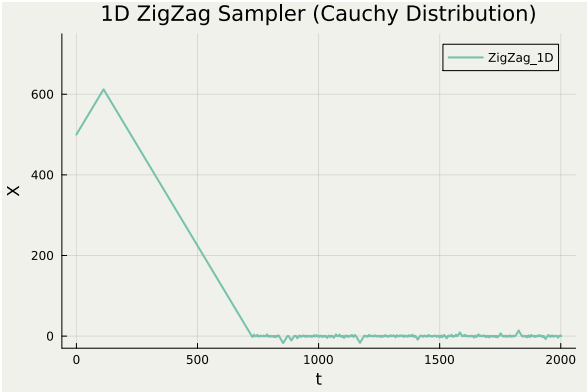
# 1.15 Comparison: Zig-Zag vs. MALA (3)

Monte Carlo estimation is also done differently



MALA outputs  $(X_n)_{n \in \mathbb{N}}$

$$\frac{1}{N} \sum_{n=1}^N f(X_n) \xrightarrow{N \rightarrow \infty} \int f(x) \mu(dx)$$

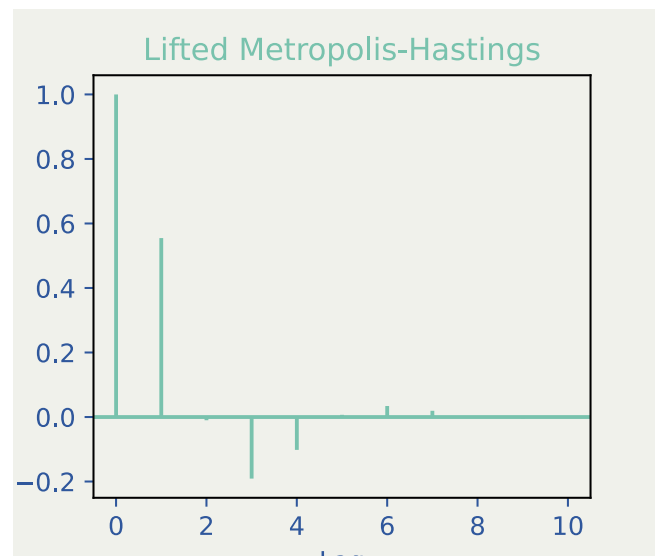
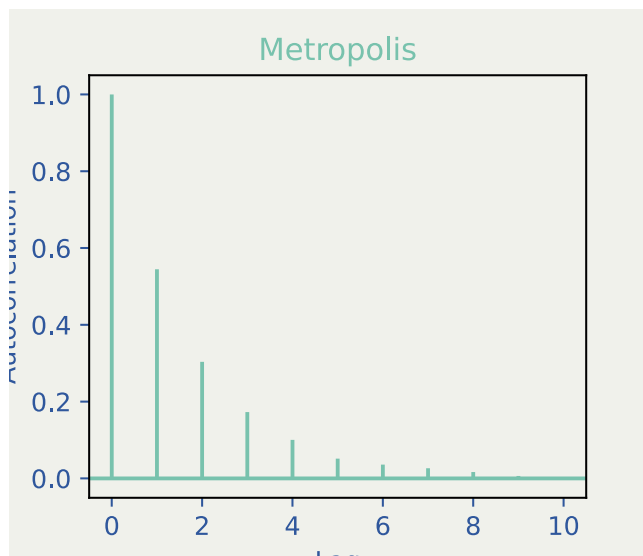


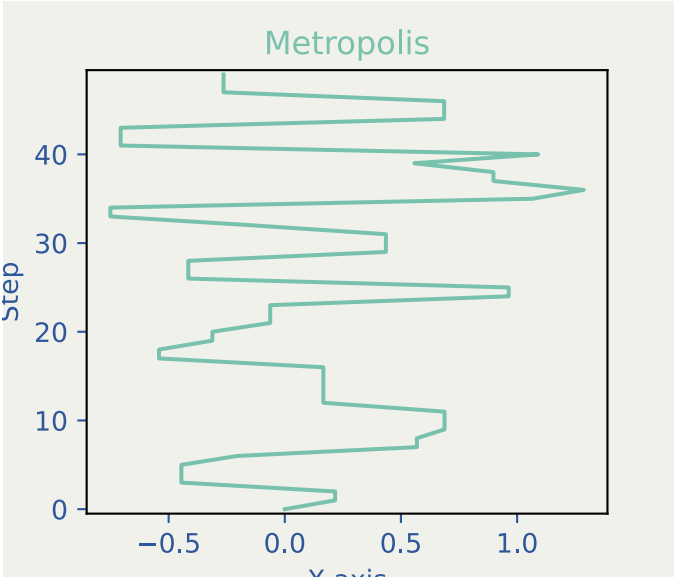
Zig-Zag outputs  $(X_t)_{t \in [0, \infty)}$

$$\int_0^T f(X_t) dt \xrightarrow{T \rightarrow \infty} \int f(x) \mu(dx)$$

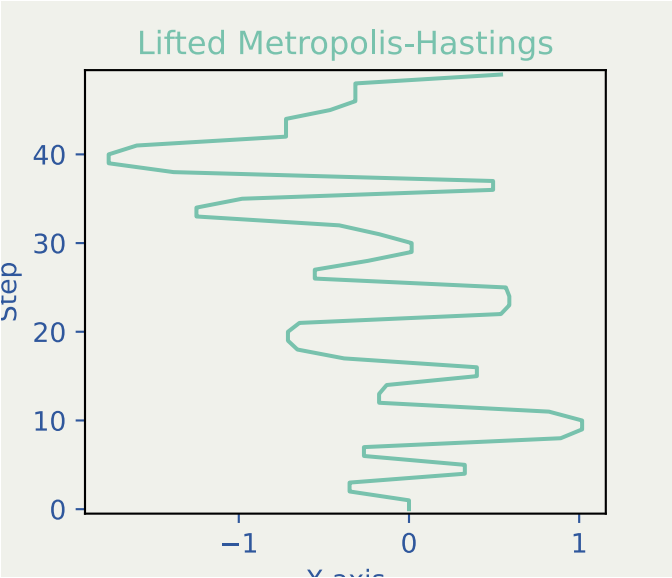
## 1.16 Recap of Section 1

- Zig-Zag Sampler's trajectory is a PDMP.
- PDMP, by design, has maximum **irreversibility**.
- **Irreversibility** leads to faster convergence comparisons against **MH**, **Lifted MH**, and





MH



Lifted MH

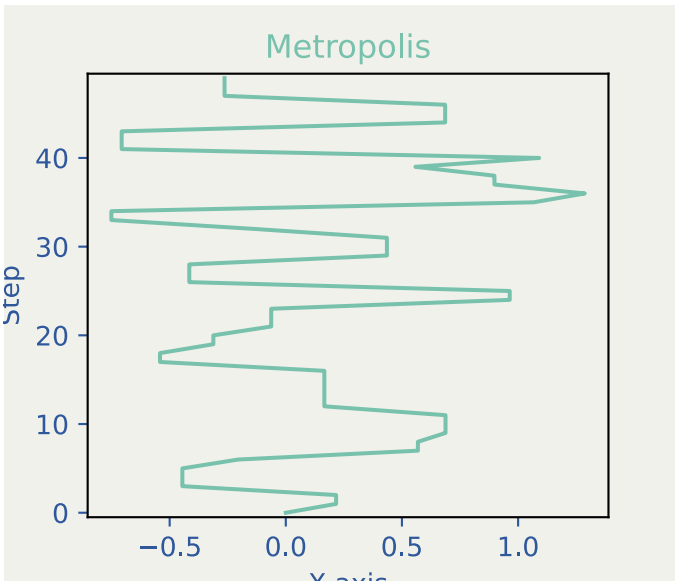


# 2 The Algorithm: How to

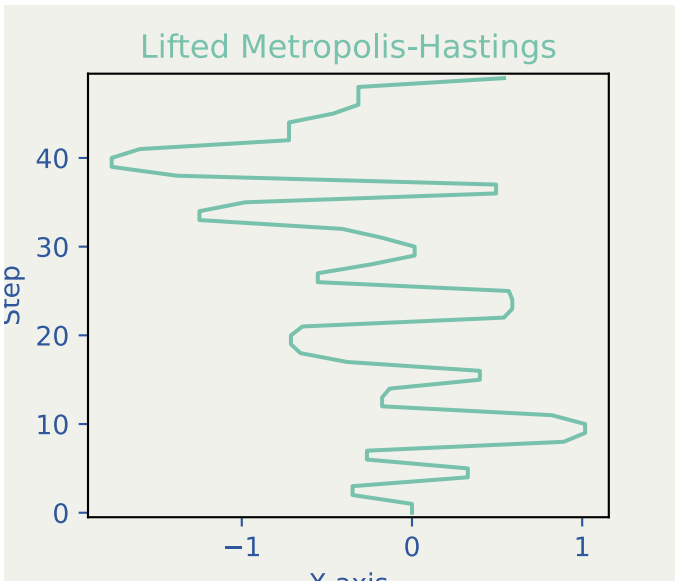
Fast and exact simulation of continuous traj

## 2.1 Review: MH vs. LMH vs. Zig-Zag (L)

As we've learned before, Zig-Zag corresponds to the **case of lifted MH** as the step size of proposal



MH

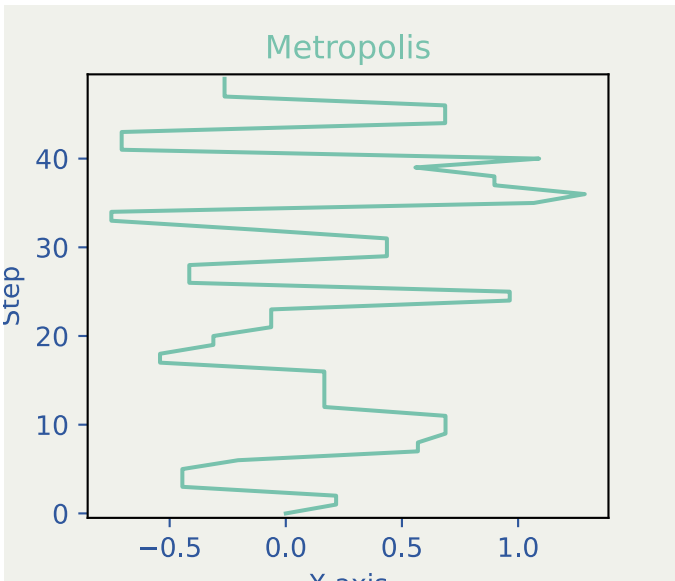


Lifted MH

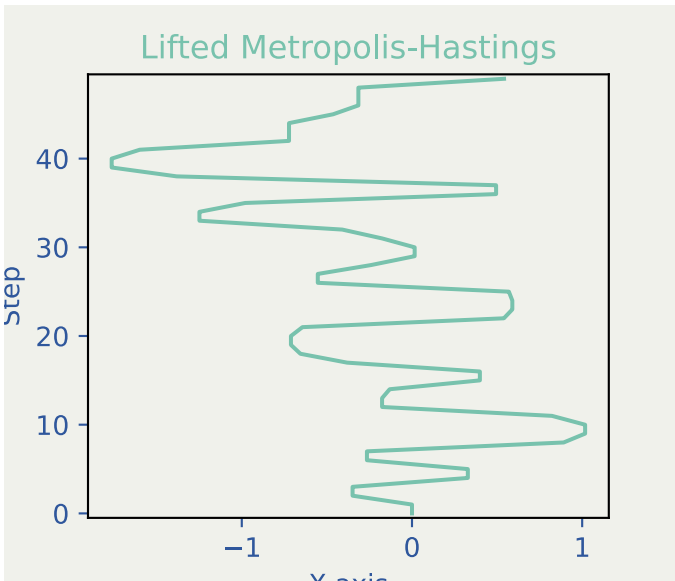
## 2.2 Review: MH vs. LMH vs. Zig-Zag (2

‘Limiting case of lifted MH’ means that we o

we should flip the momentum  $\sigma \in \{$



MH



Lifted MH

## 2.3 Algorithm (1/2)

‘Limiting case of lifted MH’ means that we only move when we see a new point. we should flip the momentum  $\sigma \in \{\pm 1\}$

(Id ' **Zig Zag sampler** Bierkens, Fearnhead, et al., 2019)

**Input:** Gradient  $\nabla \log p$  of log target density  $p$

For  $n \in \{1, 2, \dots, N\}$ :

1. Simulate an first arrival time  $T_n$  of a **Poisson point process**
2. Linearly interpolate until time  $T_n$ :

$$X_t = X_{T_{n-1}} + \sigma(t - T_{n-1}), \quad t \in [T_{n-1}, T_n]$$

3. Go back to Step 1 with the momentum  $\sigma \in \{\pm 1\}$  flipped

1. Multidimensional extension is straightforward, but we won't

## 2.4 Algorithm (2/2)

### (Fundamental Property of Zig-Zag Sampler (Id) 2019)

Let  $U(x) := -\log p(x)$ . Simulating a **Poisson point process** w

$$\lambda(x, \sigma) := \left( \sigma U'(x) \right)_+ + \gamma(x)$$

ensures the Zig-Zag sampler converges to the target  $p$ , where  $\gamma$  is a non-negative function.

Its ergodicity is ensured as long as there exists  $c > 0$  such that<sup>1</sup>

$$p(x) \leq C|x|^{-c}.$$

<sup>1</sup> I. With some regularity conditions on  $U$ . (See Hirofumi Shiba Bierkens, Robert

## 2.5 Core of the Algorithm

Given a rate function

$$\lambda(x, \sigma) := \left( \sigma U'(x) \right)_+ + \gamma$$

how to simulate a corresponding **Poisson** process

### **What We'll Learn in the Rest of this** Section 2

1. What is **Poisson Point Process**?
2. How to Simulate It?
3. Core Technique: **Poisson Thinning**

**Take Away: Zig-Zag sampling reduces to **Poisson****

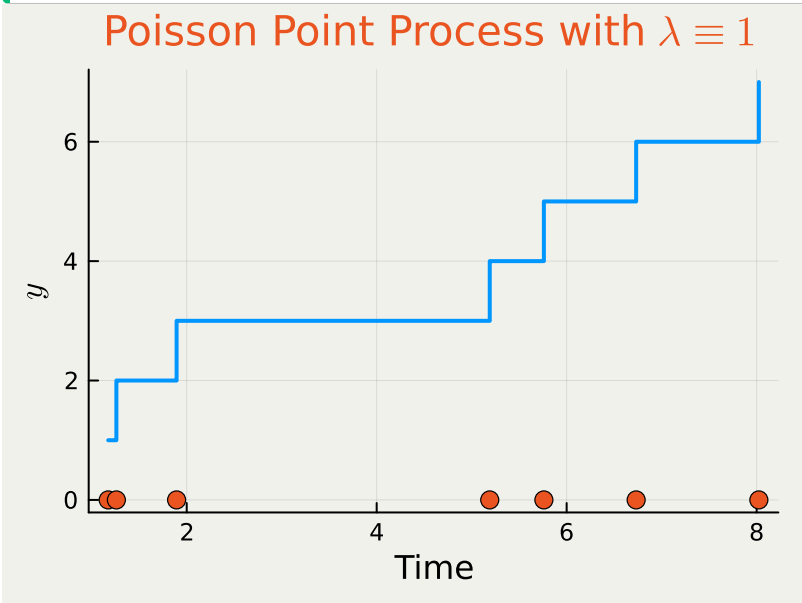
## 2.6 Simulating Poisson Point Process (

What is a **Poisson Point Process** with rate  $\lambda$ ?

The number of points in  $[0, t]$  follows a Poisson distribution with

$$N([0, t]) \sim \text{Pois}(M(t)), \quad M(t) := \int_0^t \lambda \, ds$$

We want to know when the first point  $T_1$  falls on  $[0, \infty)$ .



When  $\lambda(x, \sigma) \equiv$

- blue line: Poisson
- red dots: Poisson

satisfying  $N_t = N$   
Hirofumi Shiba

## 2.7 Simulating Poisson Point Process (2)

### Proposition (Simulation of Poisson Point Process)

The first arrival time  $T_1$  of a Poisson Point Process with rate

$$T_1 \stackrel{d}{=} M^{-1}(E), \quad E \sim \text{Exp}(1), \quad M(t) := \int_0^t \lambda(x, \sigma) dx$$

where  $\text{Exp}(1)$  denotes the exponential distribution with parameter 1.

Since  $\lambda(x, \sigma) := \left( \sigma U'(x) \right)_+ + \gamma(x)$ ,  $M$  can be complicated.

→ Inverting  $M$  can be impossible.

→ We need more general techniques: Poisson



## 2.8 Poisson Thinning (1/2)

(Lewis and Shedler, 1979)

To obtain the first arrival time  $T_1$  of a **Poisson Point Process** v

1. Find a bound  $M$  that satisfies

$$m(t) := \int_0^t \lambda(x_s, \sigma_s) ds \leq M(t).$$

2. Simulate a point  $T$  from the **Poisson Point Process** with int

3. Accept  $T$  with probability  $\frac{m(T)}{M(T)}$ .

- $m(t)$ : Defined via  $\lambda(x, \sigma) := \left( \sigma U'(x) \right)_+ + \gamma(x)$ .

- $M(t)$ : Simple upper bound  $m \leq M$ , such that  $M^{-1}$  i

## 2.9 Poisson Thinning (2/2)

In order to simulate a **Poisson Point Process**

$$\lambda(x, \sigma) := \left( \sigma U'(x) \right)_+ + \gamma$$

we find a **invertible upper bound**  $M$  that sat

$$\int_0^t \lambda(x_s, \sigma_s) ds = m(t) \leq M$$

for all possible Zig-Zag trajectories  $\{(x_s, \sigma_s)\}$

## 2.10 Recap of Section 2

1. Continuous-time MCMC, based on PDMP different algorithm and strategy.
2. To simulate PDMP is to simulate **Poisson**
3. The core technology to simulate **Poisson**  
**Poisson Thinning**.
4. **Poisson Thinning** is about finding an **upper**  
tractable inverse  $M^{-1}$ ; Typically a polynomial
5. The **upper bound**  $M$  has to be given on a

# 3 Proof of Concept: How It?

Quick demonstration of the state-of-the-art  
toy example.

### 3.1 Review: The 3 Steps of Zig-Zag Sampling

Given a target  $p$ ,

1. Calculate the negative log-likelihood  $U(x)$
2. Fix a refresh rate  $\gamma(x)$  and compute the refresh

$$\lambda(x, \sigma) := \left( \sigma U'(x) \right)_+ + \gamma(x)$$

3. Find an invertible upper bound  $M$  that satisfies

$$\int_0^t \lambda(x_s, \sigma_s) ds =: m(t) \leq M$$

## 3.2 Model: 1d Gaussian Mean Reconstruction

### Setting

- Data:  $y_1, \dots, y_n \in \mathbb{R}$  acquired by

$$y_i \stackrel{\text{iid}}{\sim} \mathcal{N}(x_0, \sigma^2), \quad i \in [n],$$

with  $\sigma > 0$  known,  $x_0 \in \mathbb{R}$  unknown.

- Prior:  $\mathcal{N}(0, \rho^2)$  with known  $\rho > 0$ .
- Goal: Sampling from the posterior

$$p(x) \propto \left( \prod_{i=1}^n \phi(x|y_i, \sigma^2) \right) \phi(x|0, \rho^2),$$

where  $\phi(x|y, \sigma^2)$  is the  $\mathcal{N}(y, \sigma^2)$  density.

The negative

$$U(x) = -\log p(x) = -\log \left( \prod_{i=1}^n \phi(x|y_i, \sigma^2) \right) \phi(x|0, \rho^2)$$

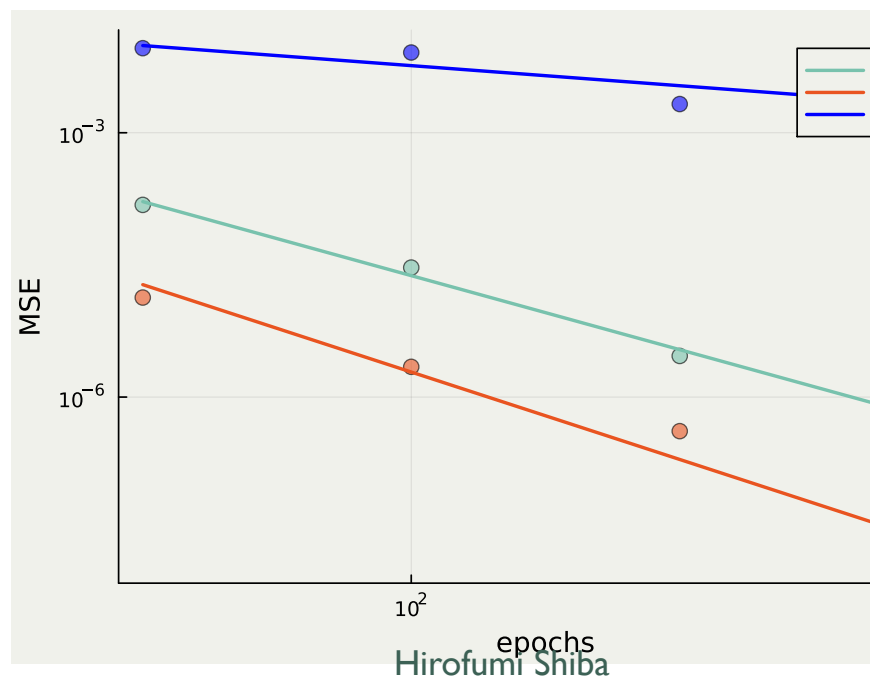
$$U'(x) = \frac{x}{\rho^2}$$

$$U''(x) = \frac{1}{\rho^2}$$

## 3.3 Menu

In the rest of this Section 3, we'll learn:

1. Even a simple Zig-Zag Sampler with  $\gamma \equiv 0$
2. Incorporating sub-sampling, Zig-Zag with further improves the efficiency.



### 3.4 Simple Zig-Zag Sampler with $\gamma \equiv 0$

Fixing  $\gamma \equiv 0$ , we obtain the upper bound  $M$

$$\begin{aligned} m(t) &= \int_0^t \lambda(x_s, \sigma_s) ds = \int_0^t \left( \sigma U'(x_s) \right) ds \\ &\leq \left( \frac{\sigma x}{\rho^2} + \frac{\sigma}{\sigma^2} \sum_{i=1}^n (x - y_i) \right) + t \left( \frac{\sigma}{\sigma^2} \right) \\ &=: (a + bt)_+ = M(t), \end{aligned}$$

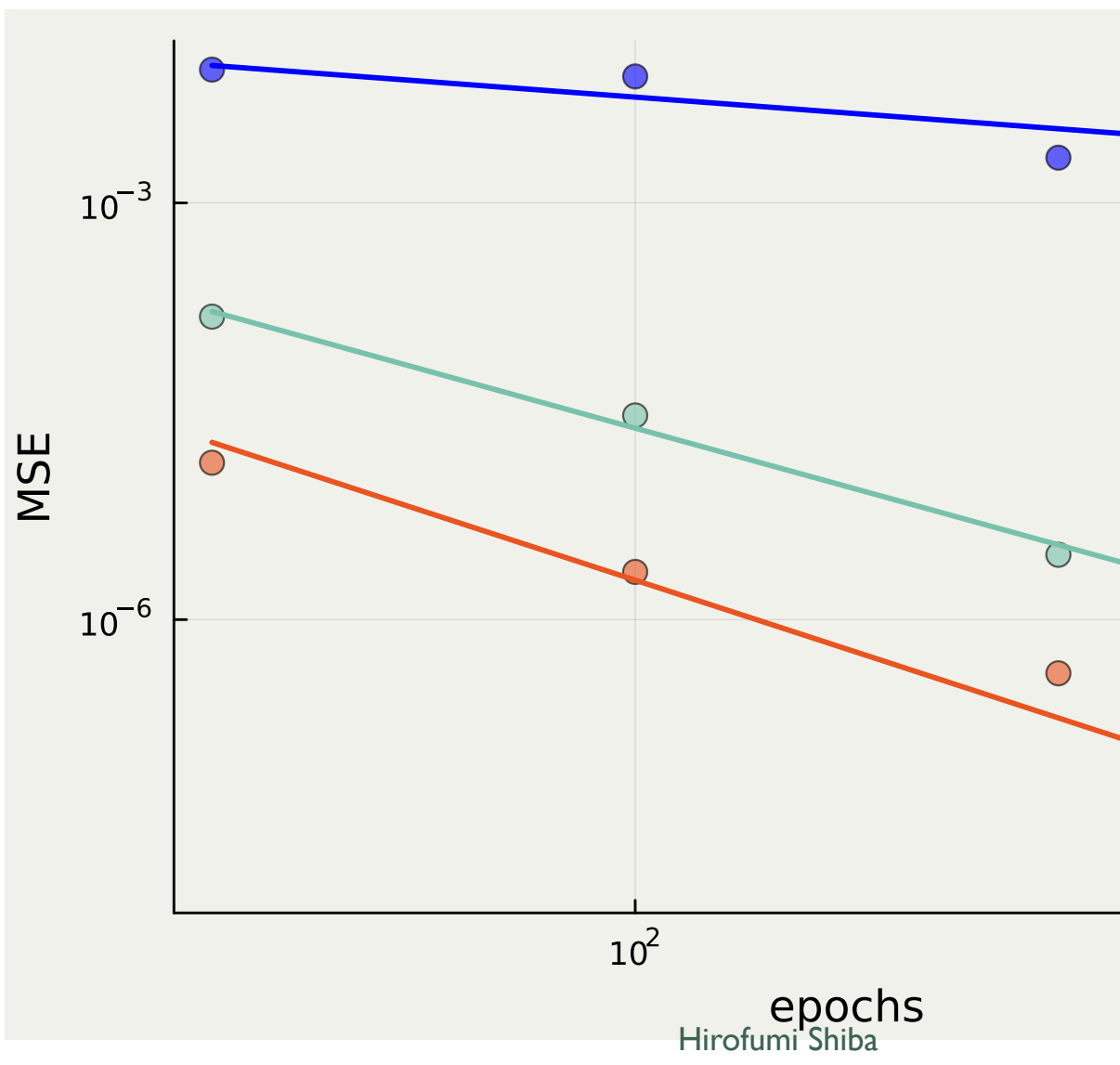
where



$$a = \frac{\sigma x}{\rho^2} + \frac{\sigma}{\sigma^2} \sum_{i=1}^n (x - y_i), \quad b =$$

# 3.5 Result: 1d Gaussian Mean Reconstruction

We generated 100 samples from  $N(x_0, \sigma^2)$  v



## 3.6 MSE per Epoch: The Vertical Axis

MSE (Mean Squared Error) of  $\{X_i\}_{i=1}^n$  is def

$$\frac{1}{n} \sum_{i=1}^n (X_i - x_0)^2.$$

Epoch: Unit computational cost.

**The following is considered as one epoch:**

- One evaluation of a likelihood ratio

$$\frac{p(X_{n+1})}{p(X_n)}.$$

- One evaluation of a **Poisson Point Process**.

## 3.7 Good News!

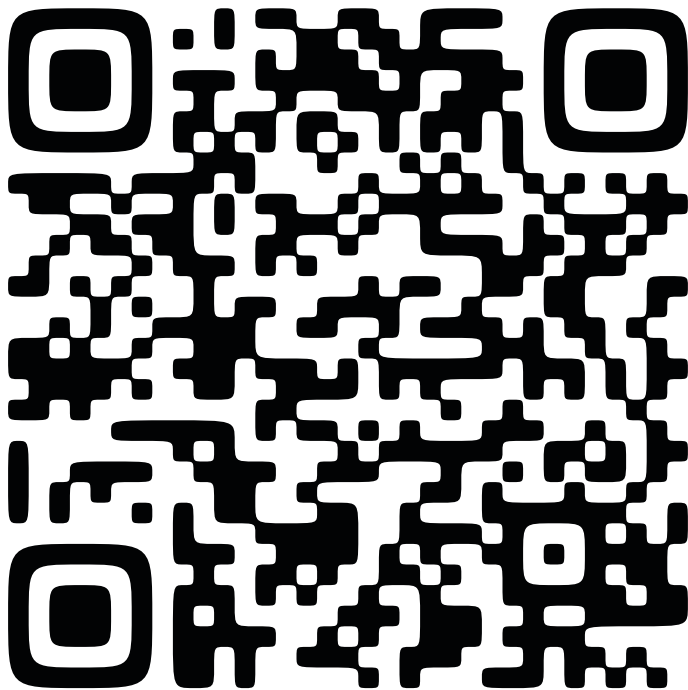
Case-by-case construction of an upper bound is often complicated / demanding.

Therefore, we are trying to automate the work.

### **Automatic Zig-Zag**

1. Automatic Zig-Zag (Corbella et al., 2022)
2. Concave-Convex PDMP (Sutton and Fearnhead, 2023)
3. NuZZ (numerical Zig-Zag) (Pagani et al., 2024)

# References



Slides and codes are available here

- Besag, J. E. (1994). *Co*  
*of Knowledge in*  
*Grenander and M*  
*Statistical Society.*  
56(4), 591–592.
- Bierkens, J., Fearnhead  
*The Zig-Zag Pro*  
*Sampling for Baye*  
*Annals of Statistics*
- Bierkens, J., Grazi, S.,  
G. O. (2020). *The*  
*Proceedings of the*  
*on Machine Learn*
- Bierkens, J., Roberts, C  
*Ergodicity of the*  
*Applied Probability*
- Bouchard-Côté, A., Vo  
Hirofumi Shiba

(2018). The bound on the error of nonreversible rejection for the monte carlo method. *Statistical Association*

Corbella, A., Spencer, S. (2022). Automating the analysis of data. *Statistics and Computing*

Dai, H., Pollock, M., and Shiba, H. (2018). Carlo Fusion. *Journal of the Royal Statistical Society* 56(1), 174–191.

Davis, M. H. A. (1984). On the convergence of markov processes to equilibrium. *diffusion stochastic processes* *Statistical Society* 46(3), 353–388.

Davis, M. H. A. (1993). *Stochastic optimization*, Vol. 1.

Duane, S., Kennedy, A. D., Roweth, D. (1987). *Letters B*, 195(2),

Fearnhead, P., Grazi, G. O. (2024). Stochastic models for general models. *Stochastic models for general models*.  
Grazi, S. (2020). Piecewise deterministic models. *Piecewise deterministic models*.  
Hastings, W. K. (1970). Monte Carlo methods using Markov chains. *Monte Carlo methods using Markov chains*.  
Lewis, P. A. W., and Sherris, C. H. (1968). Simulation of non-homogeneous Poisson processes by thinning. *Simulation of non-homogeneous Poisson processes by thinning*.  
Metropolis, N., Rosenbluth, A. W., Rosenbluth, M. N., Teller, A. H., and Teller, J. C. (1953). Monte Carlo simulation of state calculations. *Monte Carlo simulation of state calculations*.  
Pagani, F., Chevallier, A., and Cotter, S. (2024). Stochastic models for general models. *Stochastic models for general models*.  
Hirofumi Shiba

61.

- Peters, E. A. J. F., and d  
Rejection-free m  
general potentials  
Scott, S. L., Blocker, A  
H. A., George, E.  
(2016). Bayes and  
monte carlo algo  
*Management Scien*  
*Management, 11,*  
Srivastava, S., Cevher,  
(2015). WASP: S  
of subset posterio  
N. Vishwanathan,  
*eighteenth interna*  
*intelligence and sta*  
San Diego, Califo  
Sutton, M., and Fearnh  
convex PDMP-ba



1425–1435.

Turitsyn, K. S., Chertkov, M., and Welling, M.,  
Irreversible Monte Carlo Sampling via  
Efficient Sampling of Gradients,  
*Phenomena*, 240(1):1–14, 2019.

Welling, M., and Teh, Y. W.,  
Stochastic Gradient Descent  
via stochastic gradients,  
*Proceedings of the  
International Conference  
on international conference on  
machine learning*, pages 681–688. NIPS  
Omnipress, 2011.

# Appendix: Scalability by Subsampling

Construction of **ZZ-CV** (Zig-Zag with Cont

### 3.8 Review: 1d Gaussian Mean Reconstruction

$U'$  has an alternative form:

$$U'(x) = \frac{x}{\rho^2} + \frac{1}{\sigma^2} \sum_{i=1}^n (x - y_i) =: \frac{1}{\rho^2} x + \frac{n}{\sigma^2} (x - \bar{y})$$

where

$$U'_i(x) = \frac{x}{\rho^2} + \frac{n}{\sigma^2} (x - y_i)$$

→ We only need one sample  $y_i$  to evaluate

## 3.9 Randomized Rate Function

Instead of

$$\lambda_{\text{ZZ}}(x, \sigma) = \left( \sigma U'(x) \right)_+$$

we use

$$\lambda_{\text{ZZ-CV}}(x, \sigma) = \left( \sigma U'_I(x) \right)_+, \quad I \sim$$

Then, the latter is an unbiased estimator of the former

$$\mathbb{E}_{I \sim \text{U}([n])} \left[ \lambda_{\text{ZZ-CV}}(x, \sigma) \right] = \lambda_{\text{ZZ}}(x, \sigma)$$

### 3.10 Last Step: Poisson Thinning

Find an invertible upper bound  $M$  that satisfies

$$\int_0^t \lambda_{\text{ZZ-CV}}(x_s, \sigma_s) ds =: m_I(t) \leq M(t),$$

It is harder to bound  $\lambda_{\text{ZZ-CV}}$ , since it is now (random function).

### 3.1 | Upper Bound $M$ with Control Va

#### Preprocessing (once and for all)

1. Find

$$x_* := \operatorname{argmin}_{x \in \mathbb{R}} U(x)$$

2. Compute

$$U'(x_*) = \frac{x_*}{\rho^2} + \frac{1}{\sigma^2} \sum_{i=1}^n (x_* - y_i).$$

Then, with  
parameter

$$m_i(t) \leq$$

where

$$a = (\sigma U'(x_*))_+ + \|U'\|_{\text{Lip}} \|x - x_*\|_p,$$

And  $m_i$  is redefined as

$$m_i(t) = U'(x_*) + U'_i(x) - U$$

### 3.12 Subsampling with Control Variate

Zig-Zag sampler with the random rate function

$$\lambda_{\text{ZZ-CV}}(x, \sigma) = \left( \sigma U'_I(x) \right)_+, \quad I$$

and the upper bound

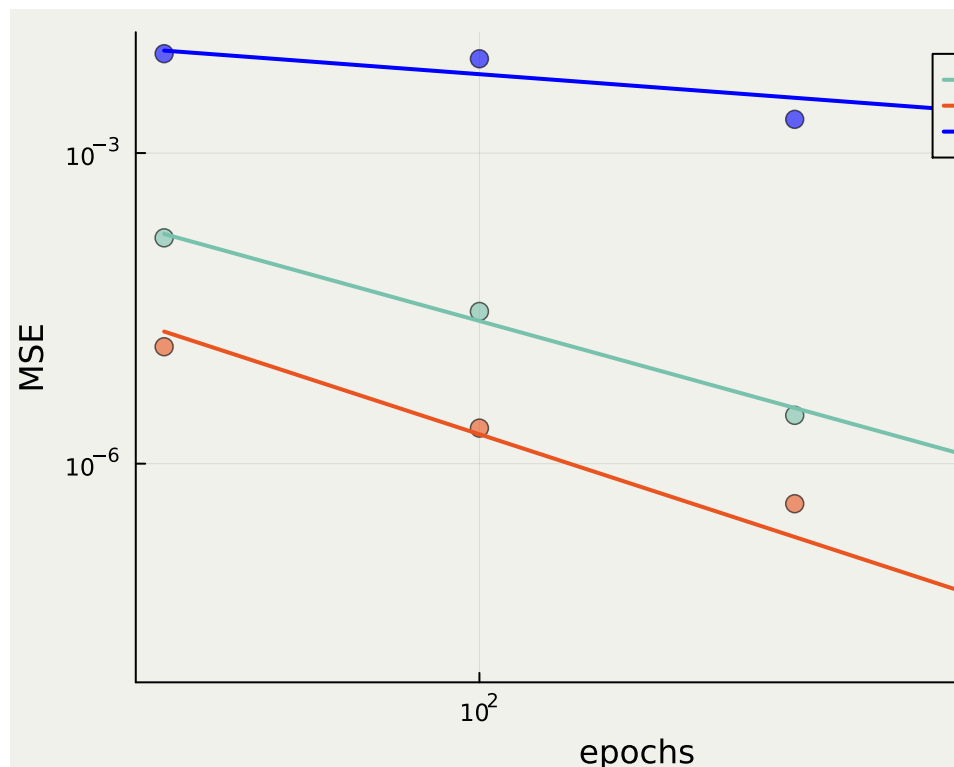
$$M(t) = a + bt$$

is called **Zig-Zag with Control Variates** (Bierman et al., 2019).



### 3.13 Zig-Zag with Control Variates

1. has  $O(1)$  efficiency as the sample size  $n$  grows
2. is exact (no bias).



I. As long as the preprocessing step is properly done.

### 3.14 Scalability (1/3)

There are currently two main approaches to scalability for large data.

#### 1. Devide-and-conquer

Devide the data into smaller **chunks** and process each **chunk**.

#### 2. Subsampling

Use a subsampling estimate of the likelihood. Methods that require the entire data.

### 3.15 Scalability (2/3) by Devide-and-co

Devide the data into smaller chunks and run chunk.

Unbiased?	Method
×	WASP
×	Consensus Monte Carlo
✓	Monte Carlo Fusion

# 3.16 Scalability (3/3) by Subsampling

Use a subsampling estimate of the likelihood  
require the entire data.

Unbiased?	Method
×	Stochastic Gradient MCMC
✓	Zig-Zag with Subsampling
×	Stochastic Gradient PDMP

0 reactions



0 comments

Write

Preview

Sign in to comment