

## **Piecewise Deterministic Monte Carlo** CIRM, Advances in Computational Statistical Physics

Joris Bierkens 21 September 2018

#### Acknowledgements

#### Collaborators



Kengo Kamatani

Gareth Roberts

Pierre-André Zitt



# Outline

#### 1 The Zig-Zag process and Bouncy Particle Sampler

**2** Simulation and Subsampling

**3** Convergence and efficiency

#### **Next Section**

#### 1 The Zig-Zag process and Bouncy Particle Sampler

2 Simulation and Subsampling

**3** Convergence and efficiency

(B., Roberts, Ann. Appl. Prob. 2017, https://arxiv.org/abs/1509.00302) The Zig-Zag process is a continuous time Markov process with states  $(X(t), V(t)) \in \mathbb{R} \times \{-1, +1\}.$ 

(B., Roberts, Ann. Appl. Prob. 2017, https://arxiv.org/abs/1509.00302) The Zig-Zag process is a continuous time Markov process with states  $(X(t), V(t)) \in \mathbb{R} \times \{-1, +1\}.$ 

X(t) moves in the direction V(t), so  $X(t) = X(0) + \int_0^t V(s) ds$ .

(B., Roberts, Ann. Appl. Prob. 2017, https://arxiv.org/abs/1509.00302) The Zig-Zag process is a continuous time Markov process with states  $(X(t), V(t)) \in \mathbb{R} \times \{-1, +1\}.$ 

X(t) moves in the direction V(t), so  $X(t) = X(0) + \int_0^t V(s) ds$ .

V(t) switches sign with switching intensity  $\lambda(X(t), V(t))$ , i.e. the first switching time T has distribution

$$\mathbb{P}(T \ge t) = \exp\left(-\int_0^t \lambda(X(s), V(s)) \, ds\right).$$

(B., Roberts, Ann. Appl. Prob. 2017, https://arxiv.org/abs/1509.00302) The Zig-Zag process is a continuous time Markov process with states  $(X(t), V(t)) \in \mathbb{R} \times \{-1, +1\}.$ 

X(t) moves in the direction V(t), so  $X(t) = X(0) + \int_0^t V(s) ds$ .

V(t) switches sign with switching intensity  $\lambda(X(t), V(t))$ , i.e. the first switching time T has distribution

$$\mathbb{P}(T \ge t) = \exp\left(-\int_0^t \lambda(X(s), V(s)) \, ds\right).$$



Joris Bierkens (TU Delft)

(B., Roberts, Ann. Appl. Prob. 2017, https://arxiv.org/abs/1509.00302)

• Potential  $U(x) = -\log \pi(x)$ 

(B., Roberts, Ann. Appl. Prob. 2017, https://arxiv.org/abs/1509.00302)

- Potential  $U(x) = -\log \pi(x)$
- $\pi$  is stationary if and only if  $\lambda(x, +1) \lambda(x, -1) = U'(x)$  for all x.

(B., Roberts, Ann. Appl. Prob. 2017, https://arxiv.org/abs/1509.00302)

- Potential  $U(x) = -\log \pi(x)$
- $\pi$  is stationary if and only if  $\lambda(x, +1) \lambda(x, -1) = U'(x)$  for all x.
- Equivalently,

$$\lambda(x,v) = \underbrace{\max(0,vU'(x))}_{+} + \underbrace{\gamma(x)}_{-}, \quad \gamma(x) \ge 0.$$

canonical switching rate

excess switching rate

(B., Roberts, Ann. Appl. Prob. 2017, https://arxiv.org/abs/1509.00302)

- Potential  $U(x) = -\log \pi(x)$
- $\pi$  is stationary if and only if  $\lambda(x, +1) \lambda(x, -1) = U'(x)$  for all x.
- Equivalently,

$$\lambda(x,v) = \max_{\text{canonical switching rate}} (0,vU'(x)) + \underbrace{\gamma(x)}_{\text{excess switching rate}}, \quad \gamma(x) \ge 0.$$

#### **Example:** Gaussian distribution $\mathcal{N}(0, \sigma^2)$

- Density  $\pi(x) \propto \exp(-x^2/(2\sigma^2))$
- Potential  $U(x) = x^2/(2\sigma^2)$
- Derivative  $U'(x) = x/\sigma^2$
- Switching rates  $\lambda(x, v) = (vx/\sigma^2)_+ + \gamma(x)$

(B., Fearnhead, Roberts, https://arxiv.org/abs/1607.03188)

- Target  $\pi(x) = \exp(-U(x))$  on  $\mathbb{R}^d$ .
- Set of directions  $v \in \{-1, +1\}^d$ .
- Switching rates  $\lambda_i(x, v) = (v_i \partial_i U(x))_+ + \gamma_i(x, v)$ , for i = 1, ..., d.

(B., Fearnhead, Roberts, https://arxiv.org/abs/1607.03188)

- Target  $\pi(x) = \exp(-U(x))$  on  $\mathbb{R}^d$ .
- Set of directions  $v \in \{-1, +1\}^d$ .
- Switching rates  $\lambda_i(x, v) = (v_i \partial_i U(x))_+ + \gamma_i(x, v)$ , for  $i = 1, \dots, d$ .
- The excess switching rate γ<sub>i</sub>(x, v) should not depend on the *i*-th component of v.













#### The Bouncy Particle Sampler

(Bouchard-Côté, Vollmer, Doucet, JASA, 2017, https://arxiv.org/abs/1510.02451)

The *Bouncy Particle Sampler (BPS)* is a second canonical example of a PDMP which can be used for sampling. State space is  $\mathbb{R}^d \times \mathbb{R}^d$  with stationary distribution  $\pi(x) dx \otimes \mathcal{N}(0, \sigma^2 I_n)$ . The BPS bounces off at random contours of  $\pi$ , through specular reflection. Additionally, the momentum gets refreshed after at rate  $\lambda_{ref}$ .



$$\mathcal{L}f(x,v) = \langle v, \nabla_x f(x) \rangle + \sum_{i=1}^d \left[ \left( v_i \partial_{x_i} U(x) \right)^+ + \gamma_i(x,v) \right] \left( f(x,F_iv) - f(x,v) \right),$$

where  $F_i v$  flips the *i*-th component of v

$$\mathcal{L}f(x,v) = \langle v, \nabla_x f(x) \rangle + \sum_{i=1}^d \left[ \left( v_i \partial_{x_i} U(x) \right)^+ + \gamma_i(x,v) \right] \left( f(x,F_iv) - f(x,v) \right),$$

where  $F_i v$  flips the *i*-th component of vBouncy Particle Sampler:

$$egin{aligned} \mathcal{L}f(x,v) &= \langle v, 
abla_x f(x) 
angle + \langle v, 
abla_x U(x) 
angle^+ (f(x,R(x)v) - f(x,v)) \ &+ \lambda_{\mathrm{ref}} \int_V [f(x,u) - f(x,v)] 
u(dv), \end{aligned}$$

$$\mathcal{L}f(x,v) = \langle v, \nabla_x f(x) \rangle + \sum_{i=1}^d \left[ \left( v_i \partial_{x_i} U(x) \right)^+ + \gamma_i(x,v) \right] \left( f(x,F_iv) - f(x,v) \right),$$

where  $F_i v$  flips the *i*-th component of vBouncy Particle Sampler:

$$egin{aligned} \mathcal{L}f(x,v) &= \langle v, 
abla_x f(x) 
angle + \langle v, 
abla_x U(x) 
angle^+ (f(x,R(x)v) - f(x,v)) \ &+ \lambda_{ ext{ref}} \int_V [f(x,u) - f(x,v)] 
u(dv), \end{aligned}$$

where

$$R(x)v = v - 2 \frac{\langle \nabla_x U(x), v \rangle}{\|\nabla U(x)\|^2} \nabla U(x),$$

$$\mathcal{L}f(x,v) = \langle v, \nabla_x f(x) \rangle + \sum_{i=1}^d \left[ \left( v_i \partial_{x_i} U(x) \right)^+ + \gamma_i(x,v) \right] \left( f(x,F_iv) - f(x,v) \right),$$

where  $F_i v$  flips the *i*-th component of vBouncy Particle Sampler:

$$egin{aligned} \mathcal{L}f(x,v) &= \langle v, 
abla_x f(x) 
angle + \langle v, 
abla_x U(x) 
angle^+ (f(x,R(x)v) - f(x,v)) \ &+ \lambda_{ ext{ref}} \int_V [f(x,u) - f(x,v)] 
u(dv), \end{aligned}$$

where

$$R(x)v = v - 2 \frac{\langle \nabla_x U(x), v \rangle}{\|\nabla U(x)\|^2} \nabla U(x),$$

and e.g

- $V = \mathbb{R}^d$ ,  $u(dv) \sim \mathcal{N}(0, I_d)$ , or
- $V = S^{d-1}, \nu(dv) \sim U(S^{d-1}).$

#### **Next Section**

#### **1** The Zig-Zag process and Bouncy Particle Sampler

**2** Simulation and Subsampling

3 Convergence and efficiency



(B., Fearnhead, Roberts, Ann. Stat. 2018, https://arxiv.org/abs/1607.03188)



Joris Bierkens (TU Delft)





#### Subsampling

(B., Fearnhead, Roberts, https://arxiv.org/abs/1607.03188)



Joris Bierkens (TU Delft)

#### Subsampling





#### Subsampling



Joris Bierkens (TU Delft)

#### **Next Section**

**1** The Zig-Zag process and Bouncy Particle Sampler

2 Simulation and Subsampling

**3** Convergence and efficiency

# Ergodicity of the multi-dimensional Zig-Zag process

(B., Roberts, Zitt, 2017, https://arxiv.org/abs/1712.09875) Ergodicity means that we should be able to construct *admissable paths* between any  $(x^1, v^1)$  and  $(x^2, v^2)$  in  $\mathbb{R}^d \times \{-1, +1\}^d$ .

Can we do this for the Zig-Zag process?

# Ergodicity of the multi-dimensional Zig-Zag process

(B., Roberts, Zitt, 2017, https://arxiv.org/abs/1712.09875) Ergodicity means that we should be able to construct admissable paths between any  $(x^1, v^1)$  and  $(x^2, v^2)$  in  $\mathbb{R}^d \times \{-1, +1\}^d$ .

Can we do this for the Zig-Zag process?

**Remark**: The Bouncy Particle Sampler is *not ergodic*, without refreshments. For BPS with refreshment, see (Deligiannidis, Bouchard-Côté, Doucet, *Annals of Statistics*, 2018).



(Bouchard-Côté. Vollmer. Doucet. 2017)





















## Ergodicity of the ZZP

(B., Roberts, Zitt, 2017, https://arxiv.org/abs/1712.09875)

#### **Growth conditions**

(GC1)  $\lim_{|x|\to\infty} U(x) = \infty$ 

**Theorem.** Suppose  $U \in C^3(\mathbb{R}^d)$  has a non-degenerate local minimum.

 If (GC1) is satisfied, then the Zig-Zag process is a ψ-irreducible aperiodic T-process.

## Ergodicity of the ZZP

(B., Roberts, Zitt, 2017, https://arxiv.org/abs/1712.09875)

#### **Growth conditions**

(GC1)  $\lim_{|x|\to\infty} U(x) = \infty$ 

(GC2) For some constants c > d,  $c' \in \mathbb{R}$ ,  $U(x) \ge c \ln(|x|) - c'$  for all  $x \in \mathbb{R}^d$ .

**Theorem.** Suppose  $U \in C^3(\mathbb{R}^d)$  has a non-degenerate local minimum.

- If (GC1) is satisfied, then the Zig-Zag process is a ψ-irreducible aperiodic T-process.
- 2 If (GC2) is satisfied, then the law of (X(t), V(t)) converges to  $\pi$  in total variation for all initial conditions.

## Ergodicity of the ZZP

(B., Roberts, Zitt, 2017, https://arxiv.org/abs/1712.09875)

#### **Growth conditions**

(GC1)  $\lim_{|x|\to\infty} U(x) = \infty$ (GC2) For some constants c > d,  $c' \in \mathbb{R}$ ,  $U(x) \ge c \ln(|x|) - c'$  for all  $x \in \mathbb{R}^d$ . (GC3)

$$\lim_{|x|\to\infty}\frac{\max(1,\|\operatorname{Hess} U(x)\|)}{|\nabla U(x)|}=0\quad\text{and}\quad\lim_{|x|\to\infty}\frac{|\nabla U(x)|}{U(x)}=0.$$

**Theorem.** Suppose  $U \in C^3(\mathbb{R}^d)$  has a non-degenerate local minimum.

- If (GC1) is satisfied, then the Zig-Zag process is a ψ-irreducible aperiodic T-process.
- 2 If (GC2) is satisfied, then the law of (X(t), V(t)) converges to  $\pi$  in total variation for all initial conditions.
- If the excess switching rates (γ<sub>i</sub>)<sup>d</sup><sub>i=1</sub> are bounded and (GC3) is satisfied, then the ZZP is exponentially ergodic.

(B., Kamatani, Roberts, 2018, https://arxiv.org/abs/1807.11358)

Consider a factorized target distribution  $\pi(x) = \prod_{i=1}^{d} \pi_i(x_i)$  with  $\pi_i(y) = \exp(-U_i(y))$ .

(B., Kamatani, Roberts, 2018, https://arxiv.org/abs/1807.11358)

Consider a factorized target distribution  $\pi(x) = \prod_{i=1}^{d} \pi_i(x_i)$  with  $\pi_i(y) = \exp(-U_i(y))$ .

Switching rates:  $\lambda_i(x, v) = (v_i U'_i(x_i))_+$ .

(B., Kamatani, Roberts, 2018, https://arxiv.org/abs/1807.11358)

Consider a factorized target distribution  $\pi(x) = \prod_{i=1}^{d} \pi_i(x_i)$  with  $\pi_i(y) = \exp(-U_i(y))$ .

Switching rates:  $\lambda_i(x, v) = (v_i U'_i(x_i))_+$ .

Every component of the Zig-Zag process mixes at  $\mathcal{O}(1)$ .

(B., Kamatani, Roberts, 2018, https://arxiv.org/abs/1807.11358)

Consider a factorized target distribution  $\pi(x) = \prod_{i=1}^{d} \pi_i(x_i)$  with  $\pi_i(y) = \exp(-U_i(y))$ .

Switching rates:  $\lambda_i(x, v) = (v_i U'_i(x_i))_+$ .

Every component of the Zig-Zag process mixes at  $\mathcal{O}(1)$ .

Compare to RWM  $\mathcal{O}(d)$ , MALA  $\mathcal{O}(d^{1/3})$ , HMC  $\mathcal{O}(d^{1/4})$ .

(B., Kamatani, Roberts, 2018, https://arxiv.org/abs/1807.11358)

Setup:

- Multivariate standard normal marginal distribution for  $(X_t)$  on  $\mathbb{R}^d$ .
- BPS with speeds  $V_t \in d^{1/2}S^{d-1}$  (uniformly distributed)
- BPS with refreshment rate  $\rho d^{1/2}$ .
- Scaling limits for:
  - Angular momentum  $\langle X_t, V_t 
    angle$
  - Negative log target  $d^{1/2}(d^{-1}||X_t||^2-1)$
  - First coordinate X<sub>1,t</sub>

#### 

#### BPS limit: PDMP in $\mathbb{R}$ with generator $Lf(x) = f'(x) + (x)^+(f(-x) - f(x))$ .



BPS limit: PDMP in  $\mathbb{R}$  with generator  $Lf(x) = f'(x) + (x)^+(f(-x) - f(x))$ .

ZZ limit: non-Markovian Gaussian process with same covariance structure as the BPS limit.

Joris Bierkens (TU Delft)

#### Scaling with dimension – angular momentum

(B., Kamatani, Roberts, 2018, https://arxiv.org/abs/1807.11358)

BPS limit: PDMP in  $\mathbb{R}$  with generator  $Lf(x) = f'(x) + (x)^+(f(-x) - f(x))$ .

ZZ limit: non-Markovian Gaussian process with same covariance structure as the BPS limit.

Covariance kernel K(t), satisfying  $\int_0^\infty K(t) dt = 0$ .



## Scaling with dimension – log density

(B., Kamatani, Roberts, 2018, https://arxiv.org/abs/1807.11358)



BPS limit: Ornstein Uhlenbeck process  $dY_t = -\frac{\sigma^2(\rho)}{4}Y_t dt + \sigma(\rho) dW_t$ 

## Scaling with dimension – log density

(B., Kamatani, Roberts, 2018, https://arxiv.org/abs/1807.11358)



BPS limit: Ornstein Uhlenbeck process  $dY_t = -\frac{\sigma^2(\rho)}{4}Y_t dt + \sigma(\rho) dW_t$ 

ZZ limit: "time integral" of angular momentum limit process

Joris Bierkens (TU Delft)

## Scaling with dimension – log density

(B., Kamatani, Roberts, 2018, https://arxiv.org/abs/1807.11358)

BPS limit for log density process: OU process

$$dY_t = -\frac{\sigma^2(\rho)}{4}Y_t \, dt + \sigma(\rho) \, dW_t$$



Note that  $\sigma(0) = 0$  and  $\lim_{\rho \to \infty} \sigma(\rho) = 0$ , so a non-vanishing, finite diffusion coefficient is required! Maximized at  $\rho^* \approx 1.424$ 

#### Scaling with dimension – first coordinate

(B., Kamatani, Roberts, 2018, https://arxiv.org/abs/1807.11358)



BPS limit: Ornstein Uhlenbeck process  $dY_t = -\rho^{-1}Y_t + \sqrt{2\rho^{-1}} dW_t$ 

#### Scaling with dimension – first coordinate (B., Kamatani, Roberts, 2018, https://arxiv.org/abs/1807.11358) 10 100 1000 d BPS res dfDCSskalatorTimesh res dfDCSskalatorTimesh ΖZ 3 5 regul\$725skeletorTimes regul@ZZSskaletorTimes regul\$225skeletorTimes

BPS limit: Ornstein Uhlenbeck process  $dY_t = -\rho^{-1}Y_t + \sqrt{2\rho^{-1}} dW_t$ 

ZZ limit: 1-dimensional Zig-Zag process!

Joris Bierkens (TU Delft)

(B., Kamatani, Roberts, 2018, https://arxiv.org/abs/1807.11358)

#### • Mixing times:

Method	Angular momentum	Radius	1st Coordinate
ZZ	$\mathcal{O}(1)$	$\mathcal{O}(1)$	$\mathcal{O}(1)$
BPS	$\mathcal{O}(d^{-1/2})$	$\mathcal{O}(d^{1/2})$	$\mathcal{O}(d^{1/2})$
MALA	N/A	$\mathcal{O}(d^{1/3})$	$\mathcal{O}(d^{1/3})$

(B., Kamatani, Roberts, 2018, https://arxiv.org/abs/1807.11358)

Mixing times:

Method	Angular momentum	Radius	1st Coordinate
ZZ	$\mathcal{O}(1)$	$\mathcal{O}(1)$	$\mathcal{O}(1)$
BPS	$\mathcal{O}(d^{-1/2})$	$\mathcal{O}(d^{1/2})$	$\mathcal{O}(d^{1/2})$
MALA	N/A	$\mathcal{O}(d^{1/3})$	$\mathcal{O}(d^{1/3})$

• Computational effort:

Method	Switches/time	Comp. effort/switch	Effort/time
ZZ	$\mathcal{O}(d)$	$\mathcal{O}(1)$	$\mathcal{O}(d)$
BPS	$\mathcal{O}(d^{1/2})$	$\mathcal{O}(d)$	$\mathcal{O}(d^{3/2})$
MALA	1	$\mathcal{O}(d)$	$\mathcal{O}(d)$

(B., Kamatani, Roberts, 2018, https://arxiv.org/abs/1807.11358)

• Mixing times:

Method	Angular momentum	Radius	1st Coordinate
ZZ	$\mathcal{O}(1)$	$\mathcal{O}(1)$	$\mathcal{O}(1)$
BPS	$\mathcal{O}(d^{-1/2})$	$\mathcal{O}(d^{1/2})$	$\mathcal{O}(d^{1/2})$
MALA	N/A	$\mathcal{O}(d^{1/3})$	$\mathcal{O}(d^{1/3})$

Computational effort:

Method	Switches/time	Comp. effort/switch	Effort/time
ZZ	$\mathcal{O}(d)$	$\mathcal{O}(1)$	$\mathcal{O}(d)$
BPS	$\mathcal{O}(d^{1/2})$	$\mathcal{O}(d)$	$\mathcal{O}(d^{3/2})$
MALA	1	$\mathcal{O}(d)$	$\mathcal{O}(d)$

• Full computational effort before mixing:

Method	Angular momentum	Radius	1st Coordinate
ZZ	$\mathcal{O}(d)$	$\mathcal{O}(d)$	$\mathcal{O}(d)$
BPS	$\mathcal{O}(d)$	$\mathcal{O}(d^2)$	$\mathcal{O}(d^2)$
MALA	N/A	$\mathcal{O}(d^{4/3})$	$\mathcal{O}(d^{4/3})$

#### Piecewise Deterministic Monte Carlo

- We can use piecewise deterministic Markov processes for sampling
- Unbiased estimate for the log density gradient results in correct stationary distribution.
- Theoretical scaling results for Zig-Zag very promising, but relying upon
  - Tight computational bounds
  - Product distribution (independence) assumption
- Scaling results comparable to: Andrieu, Durmus, Nüsken, Roussel, *Hypocoercivity of Piecewise Deterministic Markov Process-Monte Carlo*, http://arxiv.org/abs/1808.08592
  - Efficiency of Zig-zag is estimated to depend on the amount of target anisotropy
  - Positive refreshment for Zigzag as a technical condition.
- R package (so far: Gaussian, Logistic Regression) RZigZag available on CRAN

#### References

http://jbierkens.nl/pdmps

Thank you!