

Towards not being afraid of the big bad data set

Gareth Roberts

(joint work with Paul Fearnhead, Adam Johansen & Murray Pollock)

Gareth.o.Roberts@warwick.ac.uk

July 21st, 2015



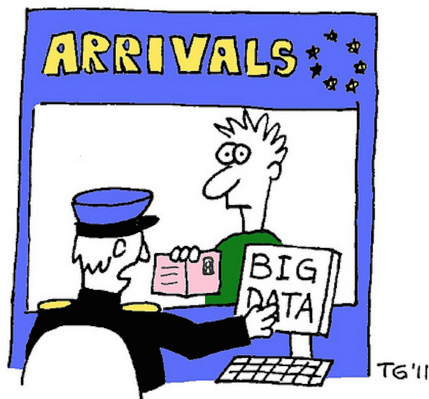
■ Motivation

- Diffusions for stationary distributions.
- Retrospective exact Monte Carlo for Diffusions
- Towards the ScaLE Algorithm

- Motivation
- Diffusions for stationary distributions.
- Retrospective exact Monte Carlo for Diffusions
- Towards the ScaLE Algorithm

- Motivation
- Diffusions for stationary distributions.
- Retrospective exact Monte Carlo for Diffusions
- Towards the ScaLE Algorithm

- Motivation
- Diffusions for stationary distributions.
- Retrospective exact Monte Carlo for Diffusions
- Towards the ScaLE Algorithm



"Your recent Amazon purchases, Tweet score and location history makes you 23.5% welcome here."



■ Big Data...

■ What is it?

- a) xxxBytes
- b) Storage
- c) Manipulation
- d) Interpretation
- e) Fit

■ Methodological / Computational Statistics (?)

- 'Robustness' / 'Scalability'
- Algorithmic Design
- 'Sufficiency'

■ Big Data...

■ What is it?

- a) xxxBytes
- b) Storage
- c) Manipulation
- d) Interpretation
- e) Fit

■ Methodological / Computational Statistics (?)

- 'Robustness' / 'Scalability'
- Algorithmic Design
- 'Sufficiency'

■ Big Data...

■ What is it?

- a) xxxBytes
- b) Storage
- c) Manipulation
- d) Interpretation
- e) Fit

■ Methodological / Computational Statistics (?)

- 'Robustness' / 'Scalability'
- Algorithmic Design
- 'Sufficiency'

■ Big Data...

■ What is it?

- a) xxxBytes
- b) Storage
- c) Manipulation
- d) Interpretation
- e) Fit

■ Methodological / Computational Statistics (?)

- 'Robustness' / 'Scalability'
- Algorithmic Design
- 'Sufficiency'

■ Big Data...

■ What is it?

- a) xxxBytes
- b) Storage
- c) Manipulation
- d) Interpretation
- e) Fit

■ Methodological / Computational Statistics (?)

- 'Robustness' / 'Scalability'
- Algorithmic Design
- 'Sufficiency'

■ Big Data...

■ What is it?

- a) xxxBytes
- b) Storage
- c) Manipulation
- d) Interpretation
- e) Fit

■ Methodological / Computational Statistics (?)

- 'Robustness' / 'Scalability'
- Algorithmic Design
- 'Sufficiency'

■ Big Data...

■ What is it?

- a) xxxBytes
- b) Storage
- c) Manipulation
- d) Interpretation
- e) Fit

■ Methodological / Computational Statistics (?)

- 'Robustness' / 'Scalability'
- Algorithmic Design
- 'Sufficiency'

- Big Data. . .
 - What is it?
 - a) xxxBytes
 - b) Storage
 - c) Manipulation
 - d) Interpretation
 - e) Fit
 - Methodological / Computational Statistics (?)
 - 'Robustness' / 'Scalability'
 - Algorithmic Design
 - 'Sufficiency'

- Big Data. . .
 - What is it?
 - a) xxxBytes
 - b) Storage
 - c) Manipulation
 - d) Interpretation
 - e) Fit
 - Methodological / Computational Statistics (?)
 - 'Robustness' / 'Scalability'
 - Algorithmic Design
 - 'Sufficiency'

- Big Data. . .
 - What is it?
 - a) xxxBytes
 - b) Storage
 - c) Manipulation
 - d) Interpretation
 - e) Fit
 - Methodological / Computational Statistics (?)
 - 'Robustness' / 'Scalability'
 - Algorithmic Design
 - 'Sufficiency'

- Big Data. . .
 - What is it?
 - a) xxxBytes
 - b) Storage
 - c) Manipulation
 - d) Interpretation
 - e) Fit
 - Methodological / Computational Statistics (?)
 - 'Robustness' / 'Scalability'
 - Algorithmic Design
 - 'Sufficiency'

“Mine Is Bigger Than Yours”



“I think you’ll find that mine is bigger...”

Typically have (with x as data if we are in the Bayesian context ...)

$$\pi(x) \propto \prod_{i=1}^N f_i(x)$$

Want to **avoid** calculating $\pi(x)$ at every iteration of an MCMC.

■ Multi-Core Methods

- Break data into K pieces / kernels
- Compute posteriors
- Recombine
- Recombination Approaches: **Averaging** (Xing / Scott / Dunson); KDE (Xing / Dunson)

■ Single-Core Methods

- Know something about your posterior – Firefly MCMC
- Pseudo-Marginal – Use subsampling to estimate likelihood...
- **Employ gradient based MCMC algorithms...**

Typically have (with x as data if we are in the Bayesian context ...)

$$\pi(x) \propto \prod_{i=1}^N f_i(x)$$

Want to **avoid** calculating $\pi(x)$ at every iteration of an MCMC.

■ Multi-Core Methods

- Break data into K pieces / kernels
- Compute posteriors
- Recombine
- Recombination Approaches: **Averaging** (Xing / Scott / Dunson); KDE (Xing / Dunson)

■ Single-Core Methods

- Know something about your posterior – Firefly MCMC
- Pseudo-Marginal – Use subsampling to estimate likelihood...
- **Employ gradient based MCMC algorithms...**

Typically have (with x as data if we are in the Bayesian context ...)

$$\pi(x) \propto \prod_{i=1}^N f_i(x)$$

Want to **avoid** calculating $\pi(x)$ at every iteration of an MCMC.

■ Multi-Core Methods

- Break data into K pieces / kernels
- Compute posteriors
- Recombine
- Recombination Approaches: **Averaging** (Xing / Scott / Dunson); KDE (Xing / Dunson)

■ Single-Core Methods

- Know something about your posterior – Firefly MCMC
- Pseudo-Marginal – Use subsampling to estimate likelihood...
- **Employ gradient based MCMC algorithms...**

Typically have (with x as data if we are in the Bayesian context ...)

$$\pi(x) \propto \prod_{i=1}^N f_i(x)$$

Want to **avoid** calculating $\pi(x)$ at every iteration of an MCMC.

■ Multi-Core Methods

- Break data into K pieces / kernels
- Compute posteriors
- Recombine
- Recombination Approaches: **Averaging** (Xing / Scott / Dunson); KDE (Xing / Dunson)

■ Single-Core Methods

- Know something about your posterior – Firefly MCMC
- Pseudo-Marginal – Use subsampling to estimate likelihood...
- **Employ gradient based MCMC algorithms...**

Typically have (with x as data if we are in the Bayesian context ...)

$$\pi(x) \propto \prod_{i=1}^N f_i(x)$$

Want to **avoid** calculating $\pi(x)$ at every iteration of an MCMC.

■ Multi-Core Methods

- Break data into K pieces / kernels
- Compute posteriors
- Recombine
- Recombination Approaches: **Averaging** (Xing / Scott / Dunson); KDE (Xing / Dunson)

■ Single-Core Methods

- Know something about your posterior – Firefly MCMC
- Pseudo-Marginal – Use subsampling to estimate likelihood...
- **Employ gradient based MCMC algorithms...**

Typically have (with x as data if we are in the Bayesian context ...)

$$\pi(x) \propto \prod_{i=1}^N f_i(x)$$

Want to **avoid** calculating $\pi(x)$ at every iteration of an MCMC.

■ Multi-Core Methods

- Break data into **K pieces** / kernels
- Compute posteriors
- Recombine
- Recombination Approaches: **Averaging** (Xing / Scott / Dunson); **KDE** (Xing / Dunson)

■ Single-Core Methods

- Know something about your posterior – Firefly MCMC
- Pseudo-Marginal – Use subsampling to estimate likelihood...
- **Employ gradient based MCMC algorithms...**

Typically have (with x as data if we are in the Bayesian context ...)

$$\pi(x) \propto \prod_{i=1}^N f_i(x)$$

Want to **avoid** calculating $\pi(x)$ at every iteration of an MCMC.

■ Multi-Core Methods

- Break data into **K pieces** / kernels
- Compute posteriors
- Recombine
- Recombination Approaches: **Averaging** (Xing / Scott / Dunson); **KDE** (Xing / Dunson)

■ Single-Core Methods

- Know something about your posterior – Firefly MCMC
- Pseudo-Marginal – Use subsampling to estimate likelihood...
- **Employ gradient based MCMC algorithms...**

Typically have (with x as data if we are in the Bayesian context ...)

$$\pi(x) \propto \prod_{i=1}^N f_i(x)$$

Want to **avoid** calculating $\pi(x)$ at every iteration of an MCMC.

■ Multi-Core Methods

- Break data into **K pieces** / kernels
- Compute posteriors
- Recombine
- Recombination Approaches: **Averaging** (Xing / Scott / Dunson); **KDE** (Xing / Dunson)

■ Single-Core Methods

- **Know something** about your posterior – Firefly MCMC
- Pseudo-Marginal – Use subsampling to estimate likelihood...
- **Employ gradient based MCMC algorithms...**

Typically have (with x as data if we are in the Bayesian context ...)

$$\pi(x) \propto \prod_{i=1}^N f_i(x)$$

Want to **avoid** calculating $\pi(x)$ at every iteration of an MCMC.

■ Multi-Core Methods

- Break data into **K pieces** / kernels
- Compute posteriors
- Recombine
- Recombination Approaches: **Averaging** (Xing / Scott / Dunson); **KDE** (Xing / Dunson)

■ Single-Core Methods

- **Know something** about your posterior – Firefly MCMC
- **Pseudo-Marginal** – Use subsampling to estimate likelihood...
- **Employ gradient based MCMC algorithms...**

Typically have (with x as data if we are in the Bayesian context ...)

$$\pi(x) \propto \prod_{i=1}^N f_i(x)$$

Want to **avoid** calculating $\pi(x)$ at every iteration of an MCMC.

■ Multi-Core Methods

- Break data into **K pieces** / kernels
- Compute posteriors
- Recombine
- Recombination Approaches: **Averaging** (Xing / Scott / Dunson); **KDE** (Xing / Dunson)

■ Single-Core Methods

- **Know something** about your posterior – Firefly MCMC
- **Pseudo-Marginal** – Use subsampling to estimate likelihood...
- **Employ gradient based MCMC algorithms...**

Eg for the **Metropolis** algorithm, need to accept a proposed move from θ to ϕ with probability

$$\min \left\{ 1, \frac{\pi(\phi)}{\pi(\theta)} \right\}$$

Pseudo-marginal MCMC (Andrieu + R, 2009, Ann Stat) allows us to instead use unbiased positive estimators of $\pi(\theta)$ and $\pi(\phi)$, accepting instead with probability

$$\min \left\{ 1, \frac{\hat{\pi}(\phi)}{\hat{\pi}(\theta)} \right\}.$$

There is **no systematic bias** induced by this: the cost comes in the mixing of the chain.

Eg for the **Metropolis** algorithm, need to accept a proposed move from θ to ϕ with probability

$$\min \left\{ 1, \frac{\pi(\phi)}{\pi(\theta)} \right\}$$

Pseudo-marginal MCMC (Andrieu + R, 2009, Ann Stat) allows us to instead use unbiased positive estimators of $\pi(\theta)$ and $\pi(\phi)$, accepting instead with probability

$$\min \left\{ 1, \frac{\hat{\pi}(\phi)}{\hat{\pi}(\theta)} \right\} .$$

There is **no systematic bias** induced by this: the cost comes in the mixing of the chain.

Can we have positive unbiased estimators for $\prod_{i=1}^N f_i(x)$ which

- 1 cost $o(N)$ to compute;
- 2 have variance which is $o(N)$?

Positive estimators of this type do exist but the answer to the above question appears to be **no**.

Estimating products unbiasedly is much more expensive than estimating sums unbiasedly.

Can we have positive unbiased estimators for $\prod_{i=1}^N f_i(x)$ which

- 1 cost $o(N)$ to compute;
- 2 have variance which is $o(N)$?

Positive estimators of this type do exist but the answer to the above question appears to be **no**.

Estimating products unbiasedly is much more expensive than estimating sums unbiasedly.

Can we have positive unbiased estimators for $\prod_{i=1}^N f_i(x)$ which

- 1 cost $o(N)$ to compute;
- 2 have variance which is $o(N)$?

Positive estimators of this type do exist but the answer to the above question appears to be **no**.

Estimating products unbiasedly is much more expensive than estimating sums unbiasedly.

Can we have positive unbiased estimators for $\prod_{i=1}^N f_i(x)$ which

- 1 cost $o(N)$ to compute;
- 2 have variance which is $o(N)$?

Positive estimators of this type do exist but the answer to the above question appears to be **no**.

Estimating products unbiasedly is much more expensive than estimating sums unbiasedly.

Traditionally used where π - high dimensional / intractable target

Our context: $\pi(x) = p(x) \prod_{i=1}^N f_i(x)$.

Langevin Diffusion: $dX_t = \frac{1}{2} \nabla \log \pi(X_t) dt + dB_t$ has invariant distribution π .

Nice structure: the diffusion drift is a sum.

$$\nabla \log \pi(x) = \nabla \log p(x) + \sum_{i=1}^N \nabla \log f_i(x)$$

Traditionally used where π - high dimensional / intractable target

Our context: $\pi(x) = p(x) \prod_{i=1}^N f_i(x)$.

Langevin Diffusion: $dX_t = \frac{1}{2} \nabla \log \pi(X_t) dt + dB_t$ has invariant distribution π .

Nice structure: the diffusion drift is a sum.

$$\nabla \log \pi(x) = \nabla \log p(x) + \sum_{i=1}^N \nabla \log f_i(x)$$

Traditionally used where π - high dimensional / intractable target

Our context: $\pi(x) = p(x) \prod_{i=1}^N f_i(x)$.

Langevin Diffusion: $dX_t = \frac{1}{2} \nabla \log \pi(X_t) dt + dB_t$ has invariant distribution π .

Nice structure: the diffusion drift is a sum.

$$\nabla \log \pi(x) = \nabla \log p(x) + \sum_{i=1}^N \nabla \log f_i(x)$$

Traditionally used where π - high dimensional / intractable target

Our context: $\pi(x) = p(x) \prod_{i=1}^N f_i(x)$.

Langevin Diffusion: $dX_t = \frac{1}{2} \nabla \log \pi(X_t) dt + dB_t$ has invariant distribution π .

Nice structure: the diffusion drift is a sum.

$$\nabla \log \pi(x) = \nabla \log p(x) + \sum_{i=1}^N \nabla \log f_i(x)$$

Two major problems with implementation of using this to sample π .

- exactness
- infinite time horizon

How can we deal with this?

- Discretise: Langevin increments $\approx N\left(\frac{1}{2}\nabla \log \pi(X_t)\Delta t, \Delta t\right)\dots$
- Euler-Maruyama: $X_{t+\Delta t} = X_t + \frac{1}{2}\nabla \log \pi(X_t)\Delta t + \xi$ where $\xi \sim N(0, \Delta t)$
- More problems...
 - Computational cost
 - Target
 - Metropolis correction (MALA)...

Two major problems with implementation of using this to sample π .

- exactness
- infinite time horizon

How can we deal with this?

- Discretise: Langevin increments $\approx N\left(\frac{1}{2}\nabla \log \pi(X_t)\Delta t, \Delta t\right)\dots$
- Euler-Maruyama: $X_{t+\Delta t} = X_t + \frac{1}{2}\nabla \log \pi(X_t)\Delta t + \xi$ where $\xi \sim N(0, \Delta t)$
- More problems...
 - Computational cost
 - Target
 - Metropolis correction (MALA)...

Two major problems with implementation of using this to sample π .

- exactness
- infinite time horizon

How can we deal with this?

- Discretise: Langevin increments $\approx N\left(\frac{1}{2}\nabla \log \pi(X_t)\Delta t, \Delta t\right)\dots$
- Euler-Maruyama: $X_{t+\Delta t} = X_t + \frac{1}{2}\nabla \log \pi(X_t)\Delta t + \xi$ where $\xi \sim N(0, \Delta t)$
- More problems...
 - Computational cost
 - Target
 - Metropolis correction (MALA)...

Two major problems with implementation of using this to sample π .

- exactness
- infinite time horizon

How can we deal with this?

- **Discretise**: Langevin increments $\approx N\left(\frac{1}{2}\nabla \log \pi(X_t)\Delta t, \Delta t\right) \dots$
- Euler-Maruyama: $X_{t+\Delta t} = X_t + \frac{1}{2}\nabla \log \pi(X_t)\Delta t + \xi$ where $\xi \sim N(0, \Delta t)$
- More problems...
 - Computational cost
 - Target
 - Metropolis correction (MALA)...

Two major problems with implementation of using this to sample π .

- exactness
- infinite time horizon

How can we deal with this?

- **Discretise**: Langevin increments $\approx N\left(\frac{1}{2}\nabla \log \pi(X_t)\Delta t, \Delta t\right) \dots$
- **Euler-Maruyama**: $X_{t+\Delta t} = X_t + \frac{1}{2}\nabla \log \pi(X_t)\Delta t + \xi$ where $\xi \sim N(0, \Delta t)$
- More problems...
 - Computational cost
 - Target
 - Metropolis correction (MALA)...

Two major problems with implementation of using this to sample π .

- exactness
- infinite time horizon

How can we deal with this?

- **Discretise**: Langevin increments $\approx N\left(\frac{1}{2}\nabla \log \pi(X_t)\Delta t, \Delta t\right) \dots$
- **Euler-Maruyama**: $X_{t+\Delta t} = X_t + \frac{1}{2}\nabla \log \pi(X_t)\Delta t + \xi$ where $\xi \sim N(0, \Delta t)$
- **More problems...**
 - Computational cost
 - Target
 - Metropolis correction (MALA)...

Two major problems with implementation of using this to sample π .

- exactness
- infinite time horizon

How can we deal with this?

- **Discretise**: Langevin increments $\approx N\left(\frac{1}{2}\nabla \log \pi(X_t)\Delta t, \Delta t\right) \dots$
- **Euler-Maruyama**: $X_{t+\Delta t} = X_t + \frac{1}{2}\nabla \log \pi(X_t)\Delta t + \xi$ where $\xi \sim N(0, \Delta t)$
- **More problems...**
 - Computational cost
 - Target
 - Metropolis correction (MALA)...

Two major problems with implementation of using this to sample π .

- exactness
- infinite time horizon

How can we deal with this?

- **Discretise**: Langevin increments $\approx N\left(\frac{1}{2}\nabla \log \pi(X_t)\Delta t, \Delta t\right) \dots$
- **Euler-Maruyama**: $X_{t+\Delta t} = X_t + \frac{1}{2}\nabla \log \pi(X_t)\Delta t + \xi$ where $\xi \sim N(0, \Delta t)$
- **More problems...**
 - Computational cost
 - Target
 - Metropolis correction (MALA)...

Two major problems with implementation of using this to sample π .

- exactness
- infinite time horizon

How can we deal with this?

- **Discretise**: Langevin increments $\approx N\left(\frac{1}{2}\nabla \log \pi(X_t)\Delta t, \Delta t\right) \dots$
- **Euler-Maruyama**: $X_{t+\Delta t} = X_t + \frac{1}{2}\nabla \log \pi(X_t)\Delta t + \xi$ where $\xi \sim N(0, \Delta t)$
- **More problems...**
 - Computational cost
 - Target
 - Metropolis correction (MALA)...

The **Exact Algorithm** for diffusion simulation (Beskos, Papaspiliopoulos and R, 2006, Bernoulli and 2008, MCAP) allows in principle to simulate **exactly** from Langevin diffusion on a fixed finite time interval.

Avoids the need for an accept/reject step!

Big problem. . .

$$\text{Step 1: } h(X_t) \propto (\pi(X_t))^{1/2} \exp \left\{ -\frac{(X_t - X_0)^2}{2t} \right\}$$

The **Exact Algorithm** for diffusion simulation (Beskos, Papaspiliopoulos and R, 2006, Bernoulli and 2008, MCAP) allows in principle to simulate **exactly** from Langevin diffusion on a fixed finite time interval.

Avoids the need for an accept/reject step!

Big problem. . .

$$\text{Step 1: } h(X_t) \propto (\pi(X_t))^{1/2} \exp \left\{ -\frac{(X_t - X_0)^2}{2t} \right\}$$

The **Exact Algorithm** for diffusion simulation (Beskos, Papaspiliopoulos and R, 2006, Bernoulli and 2008, MCAP) allows in principle to simulate **exactly** from Langevin diffusion on a fixed finite time interval.

Avoids the need for an accept/reject step!

Big problem. . .

$$\text{Step 1: } h(X_t) \propto (\pi(X_t))^{1/2} \exp \left\{ -\frac{(X_t - X_0)^2}{2t} \right\}$$

ScaLable Langevin Exact Algorithm

Continuous time, multi-level splitting, retrospective sequential sampler

The methods involves subsampling from the big data set

Requires: $(\log f_i)'$, $(\log f_i)''$, $(\log p)'$, $(\log p)''$, N , (\hat{x})

Parallelisable (Non-Trivially) (not to be discussed in this talk)

Sca_lable L_angevin E_xact Algorithm

Continuous time, multi-level splitting, retrospective sequential sampler

The method involves subsampling from the big data set

Requires: $(\log f_i)'$, $(\log f_i)''$, $(\log p)'$, $(\log p)''$, N , (\hat{x})

Parallelisable (Non-Trivially) (not to be discussed in this talk)

ScaLe Langevin Exact Algorithm

Continuous time, multi-level splitting, retrospective sequential sampler

The method involves subsampling from the big data set

Requires: $(\log f_i)'$, $(\log f_i)''$, $(\log p)'$, $(\log p)''$, N , (\hat{x})

Parallelisable (Non-Trivially) (not to be discussed in this talk)

Sca^lable L^angevin E^xact Algorithm

Continuous time, multi-level splitting, retrospective sequential sampler

The methods involves **subsampling** from the big data set

Requires: $(\log f_i)'$, $(\log f_i)''$, $(\log p)'$, $(\log p)''$, N , (\hat{x})

Parallelisable (Non-Trivially) (not to be discussed in this talk)

ScaLable Langevin Exact Algorithm

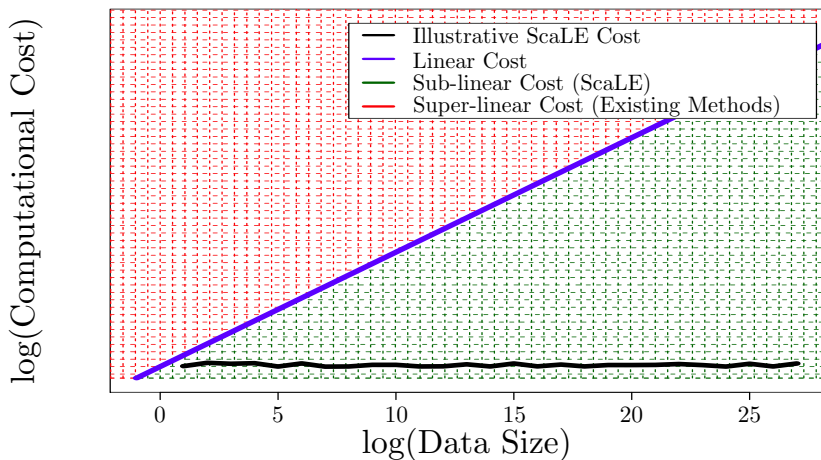
Continuous time, multi-level splitting, retrospective sequential sampler

The methods involves **subsampling** from the big data set

Requires: $(\log f_i)'$, $(\log f_i)''$, $(\log p)'$, $(\log p)''$, N , (\hat{x})

Parallelisable (Non-Trivially) **(not to be discussed in this talk)**

Computational Cost vs. Data Size



Consider a diffusion given by a d -dimensional diffusion process

$$d\mathbf{X}_s = \alpha(\mathbf{X}_s) ds + d\mathbf{B}_s, \quad s \in [0, t]. \quad (1)$$

Assume

- 1 The diffusion in (1) is non-explosive.
- 2 α is continuously differentiable in all its arguments.
- 3 There exists $l > -\infty$ such that $\phi(\mathbf{u}) := (\|\alpha(\mathbf{u})\|^2 + \nabla^2 A(\mathbf{u}))/2 - l \geq 0$.
- 4 There exists a function $A : \mathcal{X}^d \rightarrow \mathbf{R}$ such that $\alpha(\mathbf{u}) = \nabla A(\mathbf{u})$.

The transition density of (1) is typically intractable but we have the **Dacunha-Castelle** formula

$$p_t(\mathbf{x} \mid \mathbf{x}_0) = \mathcal{N}_t(\mathbf{x} - \mathbf{x}_0) \exp\{A(\mathbf{x}) - A(\mathbf{x}_0) - It\} \mathbb{E}_{\mathbf{x}_0, \mathbf{x}} \left[\exp \left\{ - \int_0^t \phi(\mathbf{X}_s) ds \right\} \right] \quad (2)$$

where $\mathcal{N}_t(\mathbf{u})$ denotes the density of the d -dimensional normal distribution with mean $\mathbf{0}$ and variance $t\mathbf{I}_d$ evaluated at $\mathbf{u} \in \mathbf{R}^d$.

The expectation is taken w.r.t. a Brownian bridge, \mathbf{x}_s , $s \in [0, t]$, with $\mathbf{X}_0 = \mathbf{x}_0$ and $\mathbf{X}_t = \mathbf{x}_t$.

The transition density of (1) is typically intractable but we have the **Dacunha-Castelle** formula

$$p_t(\mathbf{x} \mid \mathbf{x}_0) = \mathcal{N}_t(\mathbf{x} - \mathbf{x}_0) \exp\{A(\mathbf{x}) - A(\mathbf{x}_0) - It\} \mathbb{E}_{\mathbf{x}_0, \mathbf{x}} \left[\exp \left\{ - \int_0^t \phi(\mathbf{X}_s) ds \right\} \right] \quad (2)$$

where $\mathcal{N}_t(\mathbf{u})$ denotes the density of the d -dimensional normal distribution with mean $\mathbf{0}$ and variance $t\mathbf{I}_d$ evaluated at $\mathbf{u} \in \mathbf{R}^d$.

The expectation is taken w.r.t. a Brownian bridge, \mathbf{x}_s , $s \in [0, t]$, with $\mathbf{X}_0 = \mathbf{x}_0$ and $\mathbf{X}_t = \mathbf{x}_t$.

The diffusion's limiting distribution (if it exists) is more tractable.

Theorem

The diffusion in (1) is positive recurrent if and only if

$$\int_{\mathbf{R}^d} e^{2A(\mathbf{z})} d\mathbf{z} < \infty .$$

If either condition holds, then the diffusion admits a unique invariant probability measure with Lebesgue density given by

$$\nu(d\mathbf{x}) = \frac{e^{2A(\mathbf{x})} d\mathbf{x}}{\int_{\mathbf{R}^d} e^{2A(\mathbf{z})} d\mathbf{z}} := \nu(\mathbf{x}) d\mathbf{x} \quad (3)$$

and

$$p_t(\mathbf{x} \mid \mathbf{x}_0) \rightarrow \nu(\mathbf{x}) \quad (4)$$

with this convergence of densities holding for all $\mathbf{x} \in \mathbf{R}^d$, and also in L^1 .

The diffusion's limiting distribution (if it exists) is more tractable.

Theorem

The diffusion in (1) is positive recurrent if and only if

$$\int_{\mathbf{R}^d} e^{2A(\mathbf{z})} d\mathbf{z} < \infty .$$

If either condition holds, then the diffusion admits a unique invariant probability measure with Lebesgue density given by

$$\nu(d\mathbf{x}) = \frac{e^{2A(\mathbf{x})} d\mathbf{x}}{\int_{\mathbf{R}^d} e^{2A(\mathbf{z})} d\mathbf{z}} := \nu(\mathbf{x}) d\mathbf{x} \quad (3)$$

and

$$p_t(\mathbf{x} \mid \mathbf{x}_0) \rightarrow \nu(\mathbf{x}) \quad (4)$$

with this convergence of densities holding for all $\mathbf{x} \in \mathbf{R}^d$, and also in L^1 .

Maybe we can try Rejection Sampling on diffusion path space.

Let $\mathbb{Q} (= \mathbb{Q}_{0,T}^x)$ be the law of our diffusion (1), which is absolutely continuous with respect to \mathbb{W} (Brownian motion started at x) with Radon-Nikodym derivative given by Girsanov's formula:

$$\begin{aligned}\frac{d\mathbb{Q}}{d\mathbb{W}}(X) &= \exp \left\{ \int_0^T \alpha(X_s) dW_s - \frac{1}{2} \int_0^T \alpha^2(X_s) ds \right\} \\ &= \exp \left\{ A(X_T) - A(X_0) - \int_0^T \phi(X_s) ds \right\}\end{aligned}$$

(Recall, $\phi = (\alpha^2 + \alpha')/2$.)

Maybe we can try Rejection Sampling on diffusion path space.

Let $\mathbb{Q} (= \mathbb{Q}_{0,T}^x)$ be the law of our diffusion (1), which is absolutely continuous with respect to \mathbb{W} (Brownian motion started at x) with Radon-Nikodym derivative given by Girsanov's formula:

$$\begin{aligned}\frac{d\mathbb{Q}}{d\mathbb{W}}(X) &= \exp \left\{ \int_0^T \alpha(X_s) dW_s - \frac{1}{2} \int_0^T \alpha^2(X_s) ds \right\} \\ &= \exp \left\{ A(X_T) - A(X_0) - \int_0^T \phi(X_s) ds \right\}\end{aligned}$$

(Recall, $\phi = (\alpha^2 + \alpha')/2$.)

Maybe we can try Rejection Sampling on diffusion path space.

Let $\mathbb{Q} (= \mathbb{Q}_{0,T}^x)$ be the law of our diffusion (1), which is absolutely continuous with respect to \mathbb{W} (Brownian motion started at x) with Radon-Nikodym derivative given by Girsanov's formula:

$$\begin{aligned}\frac{d\mathbb{Q}}{d\mathbb{W}}(X) &= \exp \left\{ \int_0^T \alpha(X_s) dW_s - \frac{1}{2} \int_0^T \alpha^2(X_s) ds \right\} \\ &= \exp \left\{ A(X_T) - A(X_0) - \int_0^T \phi(X_s) ds \right\}\end{aligned}$$

(Recall, $\phi = (\alpha^2 + \alpha')/2$.)

Set $d\mathbb{W}$ to be probability measure proportional to $e^{A(X_T)} \cdot d\mathbb{W}$ so that

$$\frac{d\mathbb{Q}}{d\mathbb{Z}}(X) \propto \exp \left\{ - \int_0^T \phi(X_s) ds \right\}$$

Typically ϕ bounded below so this RN derivative is **bounded**.

- 1 Simulate $X \sim \mathbb{Z}$
- 2 With probability $P_{\text{W}}(X) := \frac{1}{M} \frac{d\mathbb{Q}}{d\mathbb{Z}}(X)$ set $(I = 1)$, else $(I = 0)$

$$X|I = 1 \sim \mathbb{Q}.$$

But how do we carry out rejection step?

- 1 Simulate $X \sim \mathbb{Z}$
- 2 With probability $P_{\text{W}}(X) := \frac{1}{M} \frac{d\mathbb{Q}}{d\mathbb{Z}}(X)$ set $(I = 1)$, else $(I = 0)$

$$X|I = 1 \sim \mathbb{Q}.$$

But how do we carry out rejection step?

- 1 Simulate $X \sim \mathbb{Z}$
- 2 With probability $P_{\text{W}}(X) := \frac{1}{M} \frac{d\mathbb{Q}}{d\mathbb{Z}}(X)$ set $(I = 1)$, else $(I = 0)$

$$X|(I = 1) \sim \mathbb{Q}.$$

But how do we carry out rejection step?

- 1 Simulate $X \sim \mathbb{Z}$
- 2 With probability $P_{\text{W}}(X) := \frac{1}{M} \frac{d\mathbb{Q}}{d\mathbb{Z}}(X)$ set $(I = 1)$, else $(I = 0)$

$$X|(I = 1) \sim \mathbb{Q}.$$

But how do we carry out rejection step?

How can we simulate, store and calculate integrals from $X \sim \mathbb{Z}$?

Simulation of finite skeletons of biased Brownian motion \mathbb{Z} is straightforward.

Acceptance probability can be written as

$$P = \exp \left\{ - \int_0^T (\phi(X_s) - \ell) \, ds \right\}$$

where $\phi(X_s) - \ell$ is non-negative.

P is just the probability that an event of hazard rate $\phi(X_s) - \ell$ has not occurred by time T .

Can achieve this event by Poisson thinning (sometimes quite complicated) from a constant hazard rate.

How can we simulate, store and calculate integrals from $X \sim \mathbb{Z}$?

Simulation of finite skeletons of biased Brownian motion \mathbb{Z} is straightforward.

Acceptance probability can be written as

$$P = \exp \left\{ - \int_0^T (\phi(X_s) - \ell) \, ds \right\}$$

where $\phi(X_s) - \ell$ is non-negative.

P is just the probability that an event of hazard rate $\phi(X_s) - \ell$ has not occurred by time T .

Can achieve this event by Poisson thinning (sometimes quite complicated) from a constant hazard rate.

How can we **simulate, store and calculate integrals** from $X \sim \mathbb{Z}$?

Simulation of finite skeletons of **biased Brownian motion** \mathbb{Z} is straightforward.

Acceptance probability can be written as

$$P = \exp \left\{ - \int_0^T (\phi(X_s) - \ell) \, ds \right\}$$

where $\phi(X_s) - \ell$ is non-negative.

P is just the probability that an event of **hazard rate** $\phi(X_s) - \ell$ has not occurred by time T .

Can achieve this event by Poisson thinning (sometimes quite complicated) from a constant hazard rate.

How can we **simulate, store and calculate integrals** from $X \sim \mathbb{Z}$?

Simulation of finite skeletons of **biased Brownian motion** \mathbb{Z} is straightforward.

Acceptance probability can be written as

$$P = \exp \left\{ - \int_0^T (\phi(X_s) - \ell) \, ds \right\}$$

where $\phi(X_s) - \ell$ is non-negative.

P is just the probability that an event of **hazard rate** $\phi(X_s) - \ell$ has not occurred by time T .

Can achieve this event by Poisson thinning (sometimes quite complicated) from a constant hazard rate.

How can we **simulate, store and calculate integrals** from $X \sim \mathbb{Z}$?

Simulation of finite skeletons of **biased Brownian motion** \mathbb{Z} is straightforward.

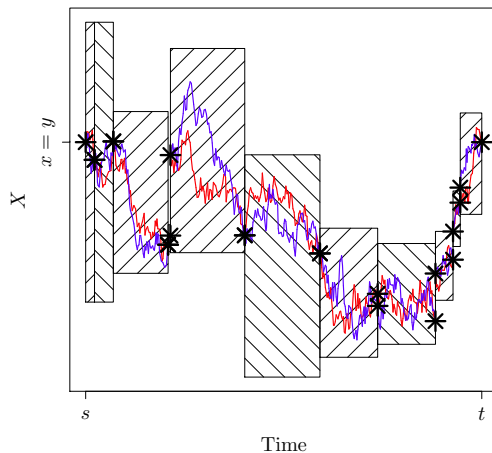
Acceptance probability can be written as

$$P = \exp \left\{ - \int_0^T (\phi(X_s) - \ell) \, ds \right\}$$

where $\phi(X_s) - \ell$ is non-negative.

P is just the probability that an event of **hazard rate** $\phi(X_s) - \ell$ has not occurred by time T .

Can achieve this event by Poisson thinning (sometimes quite complicated) from a constant hazard rate.



Many extensions of these ideas in the literature: EA0, EA1, EA2, EA3, JEA, CIS ... Relaxations of smoothness conditions, multi-dimensional, time-inhomogeneous versions of these algorithms

Methods are surprisingly efficient. There is no intrinsic cost of exactness.

Methods are genuinely multi-dimensional, but will scale at least linearly with dimension.

But existing methods do rely on being able to identify ϕ .

Many extensions of these ideas in the literature: EA0, EA1, EA2, EA3, JEA, CIS ... Relaxations of smoothness conditions, multi-dimensional, time-inhomogeneous versions of these algorithms

Methods are surprisingly efficient. There is no intrinsic cost of exactness.

Methods are genuinely multi-dimensional, but will scale at least linearly with dimension.

But existing methods do rely on being able to identify ϕ .

Many extensions of these ideas in the literature: EA0, EA1, EA2, EA3, JEA, CIS ... Relaxations of smoothness conditions, multi-dimensional, time-inhomogeneous versions of these algorithms

Methods are surprisingly efficient. **There is no intrinsic cost of exactness.**

Methods are genuinely multi-dimensional, but will scale at least linearly with dimension.

But existing methods do rely on being able to identify ϕ .

Many extensions of these ideas in the literature: EA0, EA1, EA2, EA3, JEA, CIS ... Relaxations of smoothness conditions, multi-dimensional, time-inhomogeneous versions of these algorithms

Methods are surprisingly efficient. **There is no intrinsic cost of exactness.**

Methods are genuinely multi-dimensional, but will scale at least linearly with dimension.

But existing methods do rely on being able to identify ϕ .

Many extensions of these ideas in the literature: EA0, EA1, EA2, EA3, JEA, CIS ... Relaxations of smoothness conditions, multi-dimensional, time-inhomogeneous versions of these algorithms

Methods are surprisingly efficient. **There is no intrinsic cost of exactness.**

Methods are genuinely multi-dimensional, but will scale at least linearly with dimension.

But existing methods do rely on being able to identify ϕ .

Can we use the EA framework for big data?

The idea would be to completely avoid a Metropolis-Hastings accept/reject step, which would be $O(N)$ expensive.

Recall: $\alpha(X_t) := \frac{1}{2} \nabla \log \pi(X_t)$

Two big problems:

- 1 Simulating from A is $O(N)$.
- 2 Calculating α, α' is $O(N)$

Can we use the EA framework for big data?

The idea would be to completely **avoid** a Metropolis-Hastings accept/reject step, which would be $O(N)$ expensive.

Recall: $\alpha(X_t) := \frac{1}{2} \nabla \log \pi(X_t)$

Two big problems:

- 1 Simulating from A is $O(N)$.
- 2 Calculating α, α' is $O(N)$

Can we use the EA framework for big data?

The idea would be to completely **avoid** a Metropolis-Hastings accept/reject step, which would be $O(N)$ expensive.

$$\text{Recall: } \alpha(X_t) := \frac{1}{2} \nabla \log \pi(X_t)$$

Two big problems:

- 1 Simulating from A is $O(N)$.
- 2 Calculating α, α' is $O(N)$

Can we use the EA framework for big data?

The idea would be to completely **avoid** a Metropolis-Hastings accept/reject step, which would be $O(N)$ expensive.

$$\text{Recall: } \alpha(X_t) := \frac{1}{2} \nabla \log \pi(X_t)$$

Two big problems:

- 1 Simulating from A is $O(N)$.
- 2 Calculating α, α' is $O(N)$

$$p_T(0, x) \propto \underbrace{\exp\left\{-\frac{x^2}{2T}\right\}}_{\text{Gaussian}} \cdot \underbrace{\exp\left\{\int_0^x \alpha(u) du\right\}}_{(\pi(x))^{1/2}} \cdot \underbrace{\mathbb{P}[\text{"Surv"}]}_{\text{"PE" / Aux RV F}} \rightarrow \pi(x)$$

If we ignore the middle term, we should bias $p_T(0, x)$ by the ratio $\pi(x)^{-1/2}$. Therefore expect that we have convergence of this modified continuum of distributions to $\pi(x)^{1/2}$.

So we solve problem 1 above, only to converge to the wrong distribution!

$$p_T(0, x) \propto \underbrace{\exp\left\{-\frac{x^2}{2T}\right\}}_{\text{Gaussian}} \cdot \underbrace{\exp\left\{\int_0^x \alpha(u) du\right\}}_{(\pi(x))^{1/2}} \cdot \underbrace{\mathbb{P}[\text{"Surv"}]}_{\text{"PE" / Aux RV F}} \rightarrow \pi(x)$$

If we **ignore** the middle term, we should bias $p_T(0, x)$ by the ratio $\pi(x)^{-1/2}$. Therefore expect that we have convergence of this modified continuum of distributions to $\pi(x)^{1/2}$.

So we solve problem 1 above, only to converge to the wrong distribution!

$$p_T(0, x) \propto \underbrace{\exp\left\{-\frac{x^2}{2T}\right\}}_{\text{Gaussian}} \cdot \underbrace{\exp\left\{\int_0^x \alpha(u) du\right\}}_{(\pi(x))^{1/2}} \cdot \underbrace{\mathbb{P}[\text{"Surv"}]}_{\text{"PE" / Aux RV F}} \rightarrow \pi(x)$$

If we **ignore** the middle term, we should bias $p_T(0, x)$ by the ratio $\pi(x)^{-1/2}$. Therefore expect that we have convergence of this modified continuum of distributions to $\pi(x)^{1/2}$.

So we solve **problem 1** above, only to converge to the wrong distribution!

Double Drift (!): $\alpha(X_t) := \nabla \log \pi(X_t)$

$$p_T(0, x) \propto \underbrace{\exp\left\{-\frac{x^2}{2T}\right\}}_{\text{Gaussian}} \cdot \underbrace{\exp\left\{\int_0^x \alpha(u) du\right\}}_{\pi(x)} \cdot \underbrace{\mathbb{P}[\text{"Surv"}]}_{\text{"PE" / Aux RV F}} \rightarrow (\pi(x))^2$$

$$\tilde{p}_T(0, x) \propto \underbrace{\exp\left\{-\frac{x^2}{2T}\right\}}_{\text{Gaussian}} \cdot \underbrace{\mathbb{P}[\text{"Surv"}]}_{\text{"PE" / Aux RV F}} \rightarrow \pi(x)$$

Type I **quasi-stationary distribution**

But does it converge??

Double Drift (!): $\alpha(X_t) := \nabla \log \pi(X_t)$

$$p_T(0, x) \propto \underbrace{\exp\left\{-\frac{x^2}{2T}\right\}}_{\text{Gaussian}} \cdot \underbrace{\exp\left\{\int_0^x \alpha(u) du\right\}}_{\pi(x)} \cdot \underbrace{\mathbb{P}[\text{"Surv"}]}_{\text{"PE" / Aux RV F}} \rightarrow (\pi(x))^2$$

$$\tilde{p}_T(0, x) \propto \underbrace{\exp\left\{-\frac{x^2}{2T}\right\}}_{\text{Gaussian}} \cdot \underbrace{\mathbb{P}[\text{"Surv"}]}_{\text{"PE" / Aux RV F}} \rightarrow \pi(x)$$

Type I **quasi-stationary distribution**

But does it converge??

Double Drift (!): $\alpha(X_t) := \nabla \log \pi(X_t)$

$$p_T(0, x) \propto \underbrace{\exp\left\{-\frac{x^2}{2T}\right\}}_{\text{Gaussian}} \cdot \underbrace{\exp\left\{\int_0^x \alpha(u) du\right\}}_{\pi(x)} \cdot \underbrace{\mathbb{P}[\text{"Surv"}]}_{\text{"PE" / Aux RV F}} \rightarrow (\pi(x))^2$$

$$\tilde{p}_T(0, x) \propto \underbrace{\exp\left\{-\frac{x^2}{2T}\right\}}_{\text{Gaussian}} \cdot \underbrace{\mathbb{P}[\text{"Surv"}]}_{\text{"PE" / Aux RV F}} \rightarrow \pi(x)$$

Type I quasi-stationary distribution

But does it converge??

Double Drift (!): $\alpha(X_t) := \nabla \log \pi(X_t)$

$$p_T(0, x) \propto \underbrace{\exp\left\{-\frac{x^2}{2T}\right\}}_{\text{Gaussian}} \cdot \underbrace{\exp\left\{\int_0^x \alpha(u) du\right\}}_{\pi(x)} \cdot \underbrace{\mathbb{P}[\text{"Surv"}]}_{\text{"PE" / Aux RV F}} \rightarrow (\pi(x))^2$$

$$\tilde{p}_T(0, x) \propto \underbrace{\exp\left\{-\frac{x^2}{2T}\right\}}_{\text{Gaussian}} \cdot \underbrace{\mathbb{P}[\text{"Surv"}]}_{\text{"PE" / Aux RV F}} \rightarrow \pi(x)$$

Type I **quasi-stationary distribution**

But does it converge??

Double Drift (!): $\alpha(X_t) := \nabla \log \pi(X_t)$

$$p_T(0, x) \propto \underbrace{\exp\left\{-\frac{x^2}{2T}\right\}}_{\text{Gaussian}} \cdot \underbrace{\exp\left\{\int_0^x \alpha(u) du\right\}}_{\pi(x)} \cdot \underbrace{\mathbb{P}[\text{"Surv"}]}_{\text{"PE" / Aux RV F}} \rightarrow (\pi(x))^2$$

$$\tilde{p}_T(0, x) \propto \underbrace{\exp\left\{-\frac{x^2}{2T}\right\}}_{\text{Gaussian}} \cdot \underbrace{\mathbb{P}[\text{"Surv"}]}_{\text{"PE" / Aux RV F}} \rightarrow \pi(x)$$

Type I **quasi-stationary distribution**

But does it converge??

We **only** have L^1 convergence of $p_T(0, x)$ to π as $T \rightarrow \infty$.

We **require** a stronger f -norm result:

$$p_T(0, \cdot) \rightarrow \pi$$

in $L^1(f)$ where $f(x) = e^{-A(x)} = \pi^{-1/2}(x)$ where the f -norm is given by

$$\|g\|_f = \sup_{h; |h| \leq f} \int |h(x)g(x)|dx$$

It turns out that we get this f -norm convergence (essentially) when the Langevin diffusion has invariant density ν such that $\int \nu(x)^{1/2}dx < \infty$

But this is immediate when we use $\nu = \pi^2$ (Fort and R, 2005).

We **only** have L^1 convergence of $p_T(0, x)$ to π as $T \rightarrow \infty$.

We **require** a stronger f -norm result:

$$p_T(0, \cdot) \rightarrow \pi$$

in $L^1(f)$ where $f(x) = e^{-A(x)} = \pi^{-1/2}(x)$ where the f -norm is given by

$$\|g\|_f = \sup_{h; |h| \leq f} \int |h(x)g(x)|dx$$

It turns out that we get this f -norm convergence (essentially) when the Langevin diffusion has invariant density ν such that $\int \nu(x)^{1/2}dx < \infty$

But this is immediate when we use $\nu = \pi^2$ (Fort and R, 2005).

We **only** have L^1 convergence of $p_T(0, x)$ to π as $T \rightarrow \infty$.

We **require** a stronger f -norm result:

$$p_T(0, \cdot) \rightarrow \pi$$

in $L^1(f)$ where $f(x) = e^{-A(x)} = \pi^{-1/2}(x)$ where the f -norm is given by

$$\|g\|_f = \sup_{h; |h| \leq f} \int |h(x)g(x)|dx$$

It turns out that we get this f -norm convergence (essentially) when the Langevin diffusion has invariant density ν such that $\int \nu(x)^{1/2}dx < \infty$

But this is immediate when we use $\nu = \pi^2$ (Fort and R, 2005).

We **only** have L^1 convergence of $p_T(0, x)$ to π as $T \rightarrow \infty$.

We **require** a stronger f -norm result:

$$p_T(0, \cdot) \rightarrow \pi$$

in $L^1(f)$ where $f(x) = e^{-A(x)} = \pi^{-1/2}(x)$ where the f -norm is given by

$$\|g\|_f = \sup_{h; |h| \leq f} \int |h(x)g(x)|dx$$

It turns out that we get this f -norm convergence (essentially) when the Langevin diffusion has invariant density ν such that $\int \nu(x)^{1/2}dx < \infty$

But this is immediate when we use $\nu = \pi^2$ (Fort and R, 2005).

le how do we overcome the fact that we cannot evaluate ϕ pointwise without incurring an $O(N)$ cost?

Use a **retrospective sampling** idea.

The EA construction requires (in thinning Poisson process argument) to kill a proposed path with probability

$$k = \frac{\phi(X_s)}{M}.$$

Actually we can sample an event of this probability by instead sampling from an event of probability K where K is an unbiased estimator of k taking values in $[0, 1]$.

Can do this without any loss of efficiency, unlike the pseudo-marginal MCMC methodology.

le how do we overcome the fact that we cannot evaluate ϕ pointwise without incurring an $O(N)$ cost?

Use a **retrospective sampling** idea.

The EA construction requires (in thinning Poisson process argument) to kill a proposed path with probability

$$k = \frac{\phi(X_s)}{M}.$$

Actually we can sample an event of this probability by instead sampling from an event of probability K where K is an unbiased estimator of k taking values in $[0, 1]$.

Can do this without any loss of efficiency, unlike the pseudo-marginal MCMC methodology.

le how do we overcome the fact that we cannot evaluate ϕ pointwise without incurring an $O(N)$ cost?

Use a **retrospective sampling** idea.

The EA construction requires (in thinning Poisson process argument) to kill a proposed path with probability

$$k = \frac{\phi(X_s)}{M}.$$

Actually we can sample an event of this probability by instead sampling from an event of probability K where K is an unbiased estimator of k taking values in $[0, 1]$.

Can do this without any loss of efficiency, unlike the pseudo-marginal MCMC methodology.

le how do we overcome the fact that we cannot evaluate ϕ pointwise without incurring an $O(N)$ cost?

Use a **retrospective sampling** idea.

The EA construction requires (in thinning Poisson process argument) to kill a proposed path with probability

$$k = \frac{\phi(X_s)}{M}.$$

Actually we can sample an event of this probability by instead sampling from an event of probability K where K is an unbiased estimator of k taking values in $[0, 1]$.

Can do this without any loss of efficiency, unlike the pseudo-marginal MCMC methodology.

le how do we overcome the fact that we cannot evaluate ϕ pointwise without incurring an $O(N)$ cost?

Use a **retrospective sampling** idea.

The EA construction requires (in thinning Poisson process argument) to kill a proposed path with probability

$$k = \frac{\phi(X_s)}{M}.$$

Actually we can sample an event of this probability by instead sampling from an event of probability K where K is an unbiased estimator of k taking values in $[0, 1]$.

Can do this without any loss of efficiency, unlike the pseudo-marginal MCMC methodology.

Implementation through continuous-time sequential monte Carlo methodology. **Resampling** needed to make the method robust over long time periods.

Simultaneously project a population of particles. Trajectories die according to the prescribed hazard rate, and are replaced by resampling from currently alive population.

Many important details about how to make algorithm efficient, eg by not permitting poisson rate to be $O(N)$ are omitted.

Implementation through continuous-time sequential monte Carlo methodology. **Resampling** needed to make the method robust over long time periods.

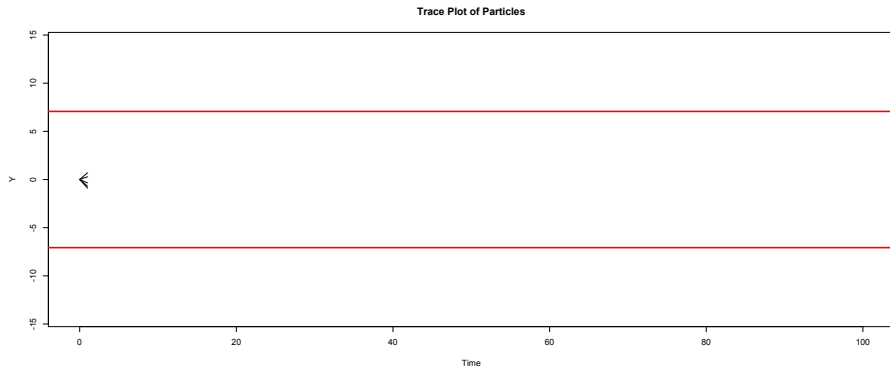
Simultaneously project a population of particles. Trajectories die according to the prescribed hazard rate, and are replaced by resampling from currently alive population.

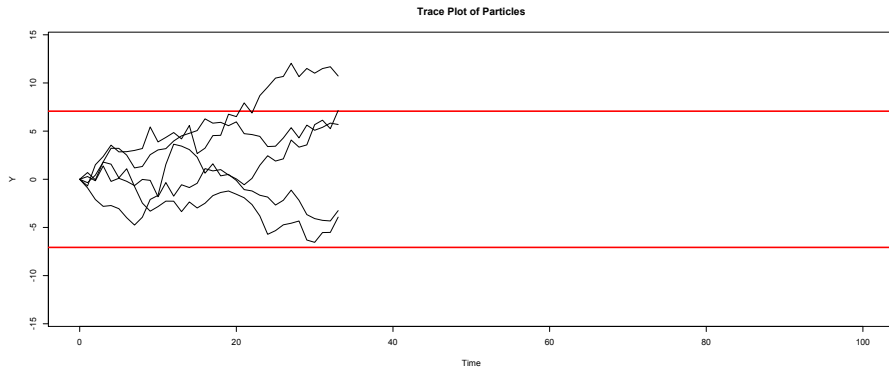
Many important details about how to make algorithm efficient, eg by not permitting poisson rate to be $O(N)$ are omitted.

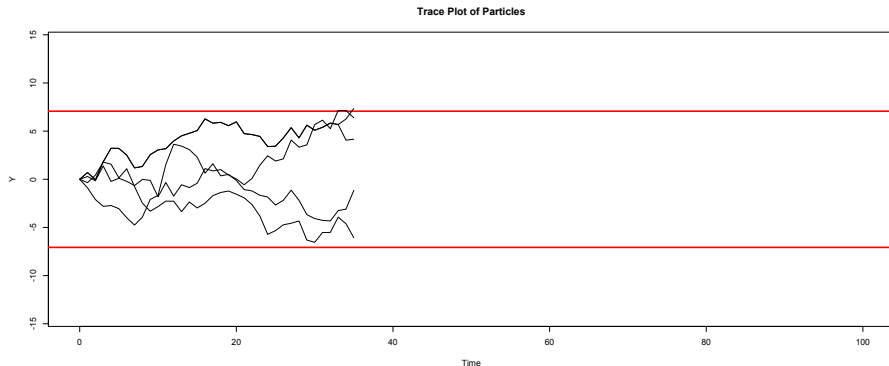
Implementation through continuous-time sequential monte Carlo methodology. [Resampling](#) needed to make the method robust over long time periods.

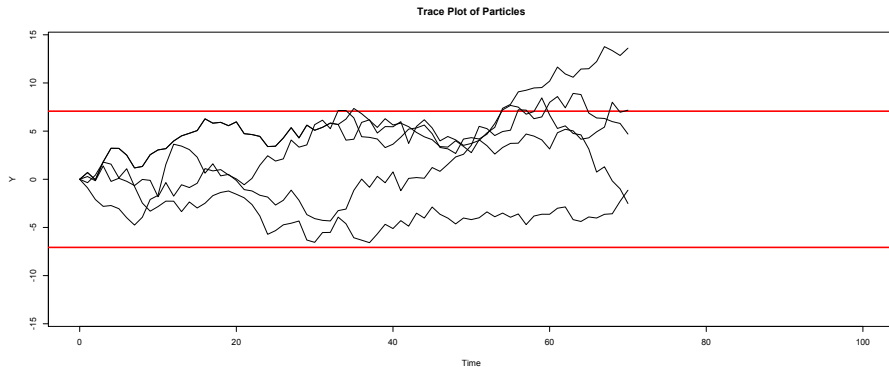
Simultaneously project a population of particles. Trajectories die according to the prescribed hazard rate, and are replaced by resampling from currently alive population.

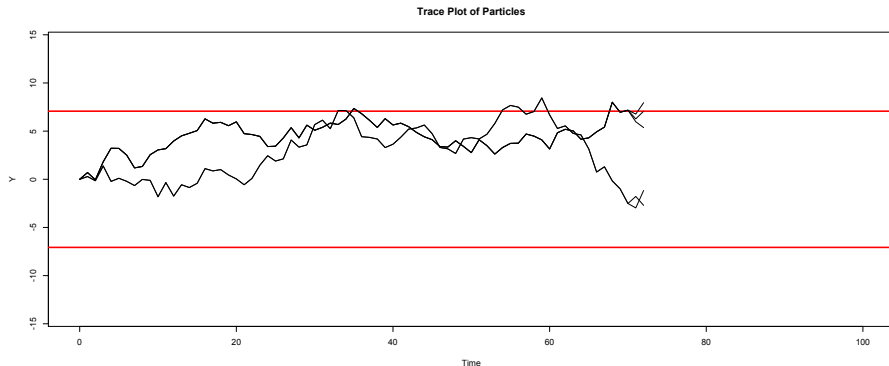
Many important details about how to make algorithm efficient, eg by not permitting poisson rate to be $O(N)$ are omitted.

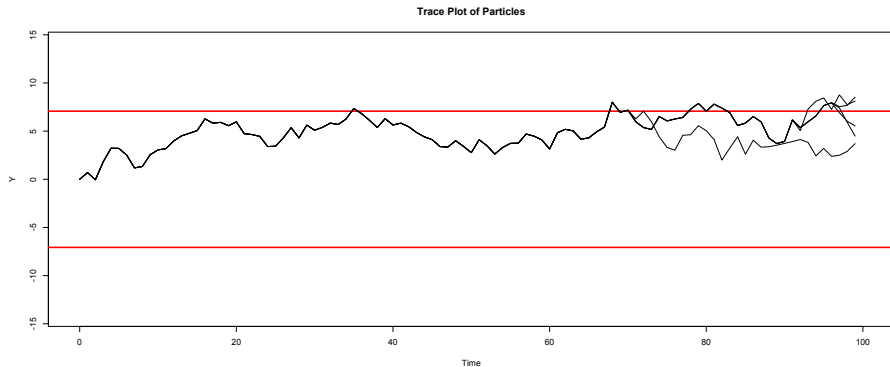


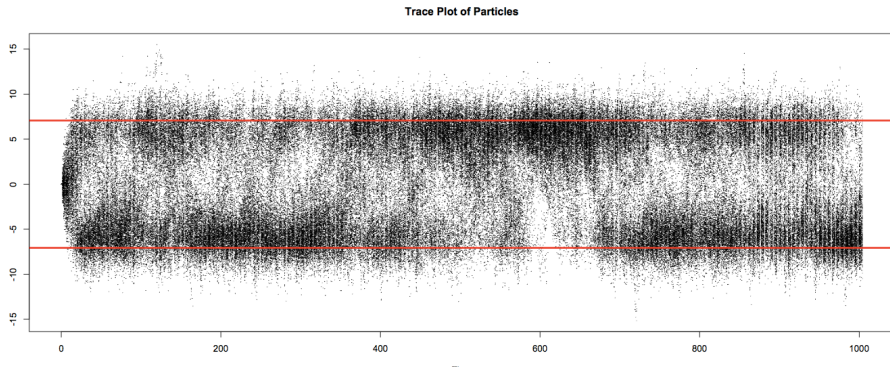


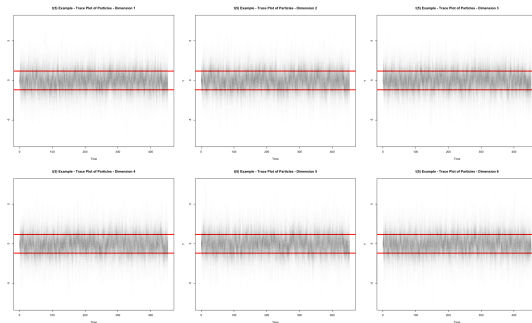




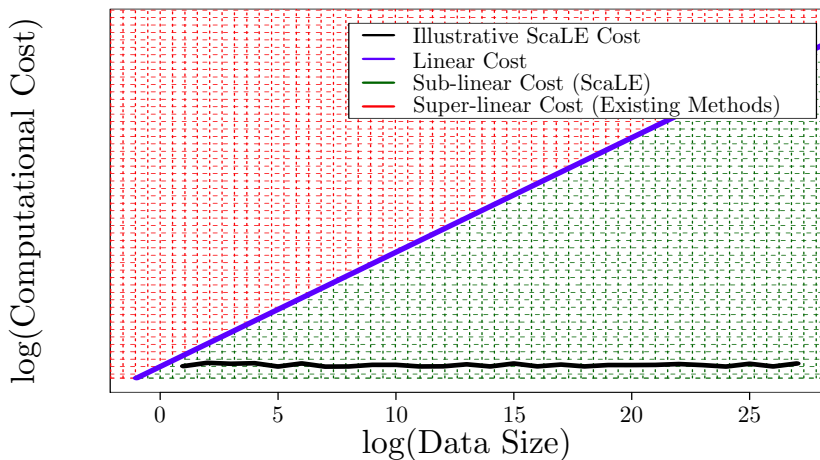








Computational Cost vs. Data Size



- This method provides provide “exact” simulation from posterior distributions from Bayesian statistical analyses, for “arbitrarily large” data sets.
- High-dimensional parameter spaces will be difficult, though not necessarily impossible to deal with.
- It is always important to bear in mind that **exactness** may not be needed or worthwhile.
- However there is **no intrinsic cost for exactness**.
- Current applications on Bayesian analysis for massive data sets: eg logistic regressions, contaminated regression models....
- Scales **extremely well** in size of data. Scaling in dimensionality of parameter space is less clear...

- This method provides provide “exact” simulation from posterior distributions from Bayesian statistical analyses, for “arbitrarily large” data sets.
- High-dimensional parameter spaces will be difficult, though not necessarily impossible to deal with.
- It is always important to bear in mind that **exactness** may not be needed or worthwhile.
- However there is **no intrinsic cost for exactness**.
- Current applications on Bayesian analysis for massive data sets: eg logistic regressions, contaminated regression models....
- Scales **extremely well** in size of data. Scaling in dimensionality of parameter space is less clear...

- This method provides provide “exact” simulation from posterior distributions from Bayesian statistical analyses, for “arbitrarily large” data sets.
- High-dimensional parameter spaces will be difficult, though not necessarily impossible to deal with.
- It is always important to bear in mind that **exactness** may not be needed or worthwhile.
- However there is **no intrinsic cost for exactness**.
- Current applications on Bayesian analysis for massive data sets: eg logistic regressions, contaminated regression models....
- Scales **extremely well** in size of data. Scaling in dimensionality of parameter space is less clear...

- This method provides provide “exact” simulation from posterior distributions from Bayesian statistical analyses, for “arbitrarily large” data sets.
- High-dimensional parameter spaces will be difficult, though not necessarily impossible to deal with.
- It is always important to bear in mind that **exactness** may not be needed or worthwhile.
- However there is **no intrinsic cost for exactness**.
- Current applications on Bayesian analysis for massive data sets: eg logistic regressions, contaminated regression models....
- Scales **extremely well** in size of data. Scaling in dimensionality of parameter space is less clear...

- This method provides provide “exact” simulation from posterior distributions from Bayesian statistical analyses, for “arbitrarily large” data sets.
- High-dimensional parameter spaces will be difficult, though not necessarily impossible to deal with.
- It is always important to bear in mind that **exactness** may not be needed or worthwhile.
- However there is **no intrinsic cost for exactness**.
- Current applications on Bayesian analysis for massive data sets: eg logistic regressions, contaminated regression models....
- Scales **extremely well** in size of data. Scaling in dimensionality of parameter space is less clear...

- This method provides provide “exact” simulation from posterior distributions from Bayesian statistical analyses, for “arbitrarily large” data sets.
- High-dimensional parameter spaces will be difficult, though not necessarily impossible to deal with.
- It is always important to bear in mind that **exactness** may not be needed or worthwhile.
- However there is **no intrinsic cost for exactness**.
- Current applications on Bayesian analysis for massive data sets: eg logistic regressions, contaminated regression models....
- Scales **extremely well** in size of data. Scaling in dimensionality of parameter space is less clear...

Consider the collection of probability measures $\{\mathbb{K}_{t,\mathbf{x}_0}, t \geq 0\}$ with $\mathbb{K}_{t,\mathbf{x}_0}$ describing a probability law on $C[0, t]$ such that $\mathbb{K}_{t,\mathbf{x}_0}(\mathbf{X}_0 = \mathbf{x}_0) = 1$ and

$$\frac{d\mathbb{K}_{t,\mathbf{x}_0}}{d\mathbb{W}_{\mathbf{x}_0}}(\mathbf{X}) = \kappa_{t,\mathbf{x}_0}^{-1} \exp \left\{ - \int_0^t \phi(\mathbf{X}_s) ds \right\} \quad (5)$$

where

$$\kappa_{t,\mathbf{x}_0} = \mathbb{E}_{\mathbb{W}_{\mathbf{x}_0}} \left[\exp \left\{ - \int_0^t \phi(\mathbf{X}_s) ds \right\} \right]. \quad (6)$$

$\mathbb{K}_{t,\mathbf{x}_0}$ can be interpreted as normalised Brownian motion killed instantaneously at a state-dependent rate $\phi(X_s)$.

Let $\mathbb{M}_{t,\mathbf{x}_0}$ be the marginal distribution of $\mathbb{K}_{t,\mathbf{x}_0}$ evaluated at time t . From (5),

$$\frac{\mathbb{M}_{t,\mathbf{x}_0}(d\mathbf{x})}{d\mathbf{x}} := m(\mathbf{x}) = \kappa_{t,\mathbf{x}_0}^{-1} \mathbb{E}_{\mathbf{x}_0,\mathbf{x}} \left[\exp \left\{ - \int_0^t \phi(\mathbf{X}_s) ds \right\} \right] \mathcal{N}_t(\mathbf{x} - \mathbf{x}_0) \quad (7)$$

which from (2) can be written

$$m(\mathbf{x}) = \kappa_{t,\mathbf{x}_0}^{-1} \exp\{-A(\mathbf{x}) + A(\mathbf{x}_0) + It\} p_t(\mathbf{x} | \mathbf{x}_0) \quad (8)$$

Since $e^{-A(\mathbf{x})}$ is unbounded, we therefore need a little more than L^1 convergence of $p_t(\mathbf{x} | \mathbf{x}_0)$ to ensure L^1 convergence of m to a probability density proportional to $e^{A(\mathbf{x})}$.

Let $\mathbb{M}_{t,\mathbf{x}_0}$ be the marginal distribution of $\mathbb{K}_{t,\mathbf{x}_0}$ evaluated at time t . From (5),

$$\frac{\mathbb{M}_{t,\mathbf{x}_0}(d\mathbf{x})}{d\mathbf{x}} := m(\mathbf{x}) = \kappa_{t,\mathbf{x}_0}^{-1} \mathbb{E}_{\mathbf{x}_0,\mathbf{x}} \left[\exp \left\{ - \int_0^t \phi(\mathbf{X}_s) ds \right\} \right] \mathcal{N}_t(\mathbf{x} - \mathbf{x}_0) \quad (7)$$

which from (2) can be written

$$m(\mathbf{x}) = \kappa_{t,\mathbf{x}_0}^{-1} \exp\{-A(\mathbf{x}) + A(\mathbf{x}_0) + It\} p_t(\mathbf{x} | \mathbf{x}_0) \quad (8)$$

Since $e^{-A(\mathbf{x})}$ is unbounded, we therefore need a little more than L^1 convergence of $p_t(\mathbf{x} | \mathbf{x}_0)$ to ensure L^1 convergence of m to a probability density proportional to $e^{A(\mathbf{x})}$.

Let $\mathbb{M}_{t,\mathbf{x}_0}$ be the marginal distribution of $\mathbb{K}_{t,\mathbf{x}_0}$ evaluated at time t . From (5),

$$\frac{\mathbb{M}_{t,\mathbf{x}_0}(d\mathbf{x})}{d\mathbf{x}} := m(\mathbf{x}) = \kappa_{t,\mathbf{x}_0}^{-1} \mathbb{E}_{\mathbf{x}_0,\mathbf{x}} \left[\exp \left\{ - \int_0^t \phi(\mathbf{X}_s) ds \right\} \right] \mathcal{N}_t(\mathbf{x} - \mathbf{x}_0) \quad (7)$$

which from (2) can be written

$$m(\mathbf{x}) = \kappa_{t,\mathbf{x}_0}^{-1} \exp\{-A(\mathbf{x}) + A(\mathbf{x}_0) + It\} p_t(\mathbf{x} | \mathbf{x}_0) \quad (8)$$

Since $e^{-A(\mathbf{x})}$ is unbounded, we therefore need a little more than L^1 convergence of $p_t(\mathbf{x} | \mathbf{x}_0)$ to ensure L^1 convergence of m to a probability density proportional to $e^{A(\mathbf{x})}$.

Fortunately stronger results exist.

Define the f norm of a signed measure ξ to be

$$\|\nu\|_f = \sup\{\xi(g); |g| \leq f\} \quad (9)$$

eg $f = 1$ is usual total variation distance.

We need an f -norm convergence result for \mathbb{M}_{t,x_0} with $f \propto e^{-A(x)}$.

The easiest theory is for the case of **geometrically ergodic Markov processes**. But here we give the more general **polynomially ergodic case**.

Fortunately stronger results exist.

Define the f norm of a signed measure ξ to be

$$\|\nu\|_f = \sup\{\xi(g); |g| \leq f\} \quad (9)$$

eg $f = 1$ is usual total variation distance.

We need an f -norm convergence result for $\mathbb{M}_{t, \mathbf{x}_0}$ with $f \propto e^{-A(\mathbf{x})}$.

The easiest theory is for the case of geometrically ergodic Markov processes. But here we give the more general polynomially ergodic case.

Fortunately stronger results exist.

Define the f norm of a signed measure ξ to be

$$\|\nu\|_f = \sup\{\xi(g); |g| \leq f\} \quad (9)$$

eg $f = 1$ is usual total variation distance.

We need an f -norm convergence result for $\mathbb{M}_{t, \mathbf{x}_0}$ with $f \propto e^{-A(\mathbf{x})}$.

The easiest theory is for the case of **geometrically ergodic Markov processes**. But here we give the more general **polynomially ergodic case**.

Theorem

Fort and R (2005)

Let $1 \leq V < \infty$ be a Borel function and $0 < \alpha \leq 1$. Assume that

- (i) some skeleton chain P^m is irreducible.*
- (ii) there exists a closed petite set C such that $\sup_C V < \infty$ and for all $\alpha \leq \eta \leq 1$, $t \mapsto V^{\eta-\alpha}(X_t)$ is integrable \mathbf{P} -a.s. and*

$$\mathcal{A}V^\eta \leq -c_\eta V^{\eta-\alpha} + b\mathbf{1}_C, \quad 0 \leq b < \infty, 0 < c_\eta < \infty. \quad (10)$$

Then there exists an unique invariant distribution π , $\pi(V^{1-\alpha}) < \infty$ and for all $0 < p < 1$ and $b \in \mathbb{R}$ or $p = 1$ and $b \geq 0$ or $p = 0$ and $b \leq 0$,

$$\lim_{t \rightarrow +\infty} (1+t)^{(1-p)(1-\alpha)/\alpha} (\log t)^b \|P^t(x, \cdot) - \pi(\cdot)\|_{V^{(1-\alpha)p} (\ln V)^{-b} \vee 1} = 0 \quad x \in \mathcal{X}.$$

Consider the simplest case - that's all we need later, although the theory is much more general.

$$d\mathbf{X}_s = \alpha(\mathbf{X}_s) ds + d\mathbf{B}_s, \quad s \in [0, t]. \quad (11)$$

$$\text{where } \alpha = \frac{\nabla \log v(\mathbf{x})}{2}$$

Very suitable for Lyapunov function methods by taking $V(\mathbf{x}) \propto \pi(\mathbf{x})^{-r}$ for some $0 < r < 1$.

Direct application as in Fort and R, 2005:

Consider the simplest case - that's all we need later, although the theory is much more general.

$$d\mathbf{X}_s = \alpha(\mathbf{X}_s) ds + d\mathbf{B}_s, \quad s \in [0, t]. \quad (11)$$

$$\text{where } \alpha = \frac{\nabla \log v(x)}{2}$$

Very suitable for Lyapunov function methods by taking $V(\mathbf{x}) \propto \pi(\mathbf{x})^{-r}$ for some $0 < r < 1$.

Direct application as in Fort and R, 2005:

Theorem

Consider ν is a positive, d -dimensional, C^2 , invariant density of \mathbf{X} . Suppose there exists some $0 < \beta < d^{-1}$ with

$$0 < \liminf_{|x| \rightarrow +\infty} \frac{|\nabla \log \nu(x)|}{\nu^\beta(x)} \leq \limsup_{|x| \rightarrow +\infty} \frac{|\nabla \log \nu(x)|}{\nu^\beta(x)} < \infty, \quad (12)$$

$$2\beta - 1 < \gamma := \liminf_{|x| \rightarrow +\infty} \frac{\text{Tr}(\nabla^2 \log \nu(x))}{|\nabla \log \nu(x)|^2} \leq \limsup_{|x| \rightarrow +\infty} \frac{\text{Tr}(\nabla^2 \log \nu(x))}{|\nabla \log \nu(x)|^2} < \infty. \quad (13)$$

For all $0 \leq \kappa < 1 + \gamma - 2\beta$,

$$\lim_{t \rightarrow +\infty} (t+1)^\tau \|P^t(x, \cdot) - \nu(\cdot)\|_{1+\nu^{-\kappa}} = 0 \quad \tau < \frac{1 + \gamma - 2\beta - \kappa}{2\beta}. \quad (14)$$

Under some regularity conditions, a density ν with tail that recede at least as quickly as

$$\|\mathbf{x}\|^{-d+k}$$

requires that $k > d$ for the conditions of the theorem to be satisfied.

In other words, we require that ν be a density such that $\nu^{1/2}$ is integrable.