

# Irreversible Langevin samplers and variance reduction: a large deviations approach

Konstantinos Spiliopoulos

Department of Mathematics & Statistics, Boston University  
Partially supported by NSF-DMS 1312124  
Joint work with Luc Rey-Bellet

# Outline

- 1 Motivation: Accelerating Monte Carlo
- 2 Investigation of convergence criteria
- 3 What can large deviations theory say?
- 4 What about variance reduction?
- 5 Simulation results
- 6 Summary
- 7 References

# Part I

## Motivation: Accelerating Monte Carlo

# Motivation

- Sampling from a given high dimensional distribution is a classical problem.
- One knows the target distributions only up to normalizing constants. Hence approximations are necessary.
- Often, such approximations are based on constructing Markov processes that have the target distribution as their target distribution, e.g., MCMC.

# Motivation

- Sampling from a given high dimensional distribution is a classical problem.
- One knows the target distributions only up to normalizing constants. Hence approximations are necessary.
- Often, such approximations are based on constructing Markov processes that have the target distribution as their target distribution, e.g., MCMC.
- The degree of how good the approximation is depends on
  - 1 the approximating Markov process, and
  - 2 on the criterion used for comparison.

## Problem formulation-Steady state simulation

Let us assume that target distribution is of Gibbs type

$$\bar{\pi}(dx) = \frac{e^{-2U(x)} dx}{\int_E e^{-2U(x)} dx}$$

Often one is interested in quantities of the form

$$\bar{f} \equiv \int_E f(x) \bar{\pi}(dx)$$

One may consider a Markov process  $X_t$  which has  $\bar{\pi}$  as its invariant distribution and under the assumption that  $X_t$  is positive recurrent, the ergodic theorem gives

$$\frac{1}{t} \int_0^t f(X_s) ds \rightarrow \int_E f(x) \bar{\pi}(dx), \text{ a.s. as } t \rightarrow \infty, \quad (1)$$

for all  $f \in L^1(\bar{\pi})$ . Hence the estimator  $f_t \equiv \frac{1}{t} \int_0^t f(X_s) ds$  can be used to approximate the expectation  $\bar{f}$ .

## Detailed Balance Condition

Most of the times, a Markov chain is constructed that is time-reversible or in other words satisfies the detailed balance condition (DBC).

If for example that target stationary distribution is  $\pi = (\pi_1, \dots, \pi_N)$ , then a **sufficient** condition to guarantee that

$$\lim_{t \rightarrow \infty} P_i(t) = \pi_i$$

is the detailed balance condition

$$\pi_i p_{ij} = \pi_j p_{ji} \text{ OR } \pi(x)P(x, y) = \pi(y)P(y, x)$$

But DBC is **only sufficient and not necessary!!**

# Questions.

- 1 What is the best Markov process  $X_t$  ?
- 2 What would be a reasonable criterion of optimality?



# Questions.

- 1 What is the best Markov process  $X_t$  ?
- 2 What would be a reasonable criterion of optimality?
- *What does large deviations theory have to say? Connections with variance reduction?*

## Part II

### Investigation of convergence criteria

## Overdamped Langevin equation.

To sample the Gibbs measure  $\bar{\pi}$  on the set  $E$

$$\bar{\pi}(dx) = \frac{e^{-2U(x)} dx}{\int_E e^{-2U(x)} dx}$$

one can consider the (time-reversible) Langevin equation

$$dX_t = -\nabla U(X_t)dt + dW_t. \quad (2)$$

## Overdamped Langevin equation.

To sample the Gibbs measure  $\bar{\pi}$  on the set  $E$

$$\bar{\pi}(dx) = \frac{e^{-2U(x)} dx}{\int_E e^{-2U(x)} dx}$$

one can consider the (time-reversible) Langevin equation

$$dX_t = -\nabla U(X_t) dt + dW_t. \quad (2)$$

How should we describe the rate of convergence

$$\frac{1}{t} \int_0^t f(X_s) ds \rightarrow \int_E f(x) \bar{\pi}(dx), ? \quad (3)$$

**Important remark:** Need to use a performance measure that works directly with the empirical measure, which is what is used in practice!!!!

## Standard measures of performance and their problems.

**Spectral gap and 2nd eigenvalue.** Consider the transition kernel

$$p(t, x_0, dx) = \mathbb{P}[X_t \in dx | X_0 = x_0]$$

Under the appropriate conditions, we have that the limit

$$p(t, x_0, dx) \rightarrow \pi(dx)$$

is determined by the spectral gap of the generator corresponding to  $X$  and the rate is exponential. In other words:

$$\|E.f(X_t) - \bar{f}\|_{L^2(\bar{\pi})} \leq C_0 \|f - \bar{f}\| e^{-\lambda t}$$

where

$\lambda = \inf\{\text{real part of non-zero eigenvalues in the spectrum of the operator}\}.$

## Standard measures of performance and their problems.

**Spectral gap and 2nd eigenvalue.** Consider the transition kernel

$$p(t, x_0, dx) = \mathbb{P}[X_t \in dx | X_0 = x_0]$$

Under the appropriate conditions, we have that the limit

$$p(t, x_0, dx) \rightarrow \pi(dx)$$

is determined by the spectral gap of the generator corresponding to  $X$  and the rate is exponential. In other words:

$$\|E.f(X_t) - \bar{f}\|_{L^2(\bar{\pi})} \leq C_0 \|f - \bar{f}\| e^{-\lambda t}$$

where

$\lambda = \inf\{\text{real part of non-zero eigenvalues in the spectrum of the operator}\}.$

Is spectral gap the most appropriate measure to characterize convergence?

## Another standard measures of performance.

**Asymptotic variance:** Hard to compute and it is a property of the algorithm only when we are at equilibrium. This is ok if we are interested in steady-state simulation.

$$t^{1/2} \left( \frac{1}{t} \int_0^t f(X_s) ds - \int f d\bar{\pi} \right) \Rightarrow N(0, \sigma_f^2)$$

and the asymptotic variance  $\sigma_f^2$  is given in terms of the integrated autocorrelation function,

$$\sigma_f^2 = 2 \int_0^\infty \mathbb{E}_{\bar{\pi}} [(f(X_0) - \bar{f})(f(X_t) - \bar{f})] dt$$

Hard to compute.

**Question:** Are there other possible Markov processes to use? How to compare their performance?

## Overdamped Langevin equation.

To sample the Gibbs measure  $\bar{\pi}$  on the set  $E$

$$\bar{\pi}(dx) = \frac{e^{-2U(x)} dx}{\int_E e^{-2U(x)} dx}$$

one can consider the (time-reversible) Langevin equation

$$dX_t = -\nabla U(X_t) dt + dW_t. \quad (4)$$



## Overdamped Langevin equation.

To sample the Gibbs measure  $\bar{\pi}$  on the set  $E$

$$\bar{\pi}(dx) = \frac{e^{-2U(x)} dx}{\int_E e^{-2U(x)} dx}$$

one can consider the (time-reversible) Langevin equation

$$dX_t = -\nabla U(X_t) dt + dW_t. \quad (4)$$

There are however many other stochastic differential equations with the same invariant measure and we may consider instead the family of equations

$$dX_t = [-\nabla U(X_t) + C(X_t)] dt + dW_t$$

where the vector field  $C(x)$  satisfies the condition

$$\operatorname{div}(Ce^{-2U}) = 0$$

In this case the Markov process is not time-reversible!!

## Overdamped Langevin equation.

There are many such  $C$ , indeed since  $\operatorname{div}(Ce^{-2U}) = 0$  is equivalent to

$$\operatorname{div}(C) = 2C\nabla U,$$

so that we, for example, can choose  $C$  to be both divergence free and orthogonal to  $\nabla U$ .

In particular, it is proved in [Barbarosie, 2011] that in dimension  $d$  any divergence free vector field can be written, locally, as the exterior product  $C = \nabla V_1 \wedge \cdots \wedge \nabla V_{n-1}$  for some for  $V_i \in \mathcal{C}^1(E; \mathbb{R})$ .

## Overdamped Langevin equation.

There are many such  $C$ , indeed since  $\operatorname{div}(Ce^{-2U}) = 0$  is equivalent to

$$\operatorname{div}(C) = 2C\nabla U,$$

so that we, for example, can choose  $C$  to be both divergence free and orthogonal to  $\nabla U$ .

In particular, it is proved in [Barbarosie, 2011] that in dimension  $d$  any divergence free vector field can be written, locally, as the exterior product  $C = \nabla V_1 \wedge \cdots \wedge \nabla V_{n-1}$  for some  $V_i \in \mathcal{C}^1(E; \mathbb{R})$ .

Therefore we can **pick**  $C$  of the form

$$C = \nabla U \wedge \nabla V_2 \cdots \nabla V_{n-1}.$$

for arbitrary  $V_2, \cdots, V_{n-1}$ . For  $d = 2$ ,  $C$  has the general form  $C = S\nabla U$  for any antisymmetric  $S$  and for  $d = 3$   $C$  has the general form  $\nabla U \times \nabla V$  for some  $V \in \mathcal{C}^1(E; \mathbb{R})$ .

## Donsker-Varadhan and Gärtner large deviations theory.

From a practical Monte-Carlo point of view one is interested in the distribution of the ergodic average  $t^{-1} \int_0^t f(X_s) ds$  and how likely it is that this average differs from  $\int f d\bar{\pi}$ .

## Donsker-Varadhan and Gärtner large deviations theory.

From a practical Monte-Carlo point of view one is interested in the distribution of the ergodic average  $t^{-1} \int_0^t f(X_s) ds$  and how likely it is that this average differs from  $\int f d\bar{\pi}$ .

Define the empirical measure

$$\pi_t \equiv \frac{1}{t} \int_0^t \delta_{X_s} ds$$

which converges to  $\bar{\pi}$  almost surely. If we have a large deviation for the family of measures  $\pi_t$ , which we write, symbolically as

$$\mathbb{P} \{ \pi_t \approx \mu \} \asymp e^{-tI_C(\mu)}$$

## Donsker-Varadhan and Gärtner large deviations theory.

From a practical Monte-Carlo point of view one is interested in the distribution of the ergodic average  $t^{-1} \int_0^t f(X_s) ds$  and how likely it is that this average differs from  $\int f d\bar{\pi}$ .

Define the empirical measure

$$\pi_t \equiv \frac{1}{t} \int_0^t \delta_{X_s} ds$$

which converges to  $\bar{\pi}$  almost surely. If we have a large deviation for the family of measures  $\pi_t$ , which we write, symbolically as

$$\mathbb{P} \{ \pi_t \approx \mu \} \asymp e^{-tI_C(\mu)}$$

Note that rate function  $I_C(\mu)$  quantifies the exponential rate at which the random measure  $\pi_t$  converges to  $\bar{\pi}$ . **Clearly, the larger  $I_C$  is, the faster the convergence occurs.**

## Standard measures of performance and their problems.

How should we describe the rate of convergence

$$\frac{1}{t} \int_0^t f(X_s) ds \rightarrow \int_E f(x) \bar{\pi}(dx), ? \quad (5)$$

**Important remark:** Standard measures of performance do not work directly with the empirical measure, which is what is used in practice!!!!

**Spectral gap and 2nd eigenvalue.** Consider the transition kernel

$$p(t, x_0, dx) = \mathbb{P}[X_t \in dx | X_0 = x_0]$$

Under the appropriate conditions, we have that the limit

$$p(t, x_0, dx) \rightarrow \pi(dx)$$

is determined by the spectral gap of the generator corresponding to  $X$  and the rate is exponential.

## Standard measures of performance and their problems.

How should we describe the rate of convergence

$$\frac{1}{t} \int_0^t f(X_s) ds \rightarrow \int_E f(x) \bar{\pi}(dx), ? \quad (5)$$

**Important remark:** Standard measures of performance do not work directly with the empirical measure, which is what is used in practice!!!!

**Spectral gap and 2nd eigenvalue.** Consider the transition kernel

$$p(t, x_0, dx) = \mathbb{P}[X_t \in dx | X_0 = x_0]$$

Under the appropriate conditions, we have that the limit

$$p(t, x_0, dx) \rightarrow \pi(dx)$$

is determined by the spectral gap of the generator corresponding to  $X$  and the rate is exponential.

**Is spectral gap the most appropriate measure to characterize convergence?**



## Other standard measures of performance.

- **Asymptotic variance:** Hard to compute and it is a property of the algorithm only when we are at equilibrium. This is ok if we are interested in steady-state simulation.
- **Large deviations rate function.** It quantifies the exponential rate at which the empirical measure  $\pi_t$  converges to the Gibbs measure  $\bar{\pi}$ . The larger the rate function is, the fastest the convergence is.

Moreover, it turns out that we can characterize asymptotic variance, using the large deviations rate function!

## Standard measures of performance.

Is spectral gap the most appropriate measure to characterize convergence?

The **problem** with spectral gap is that the information is on the density at fixed time  $t$  and not on the empirical measure. But empirical measure depends on sample path. Hence, spectral gap neglects potentially significant effects of time averaging in empirical measure (see also Dupuis, Liu, Plattner, and J. D. Doll (2012)).

**Counter example.** Let us consider the family of diffusions

$$dX_t = \delta dt + dW_t$$

on the circle  $S^1$  with generator

$$\mathcal{L}_\delta = \Delta + \delta \nabla$$

For any  $\delta \in \mathbb{R}$  the Lebesgue measure on  $S^1$  is invariant but the diffusion is reversible only if  $\delta = 0$ .

## Standard measures of performance.

The eigenvalues and eigenfunctions of  $\mathcal{L}_\delta$  are

$$e_n = e^{inx}, \quad \lambda_n = -n^2 + in\delta, \quad n \in \mathbb{Z}.$$

The spectral gap is  $-1$  for all  $\delta \in \mathbb{R}$ , i.e. the spectral gap does not move. However the asymptotic variance does decrease. For any real-valued function  $f$  with  $\int_{S^1} f dx = 0$  we have for the asymptotic variance of the estimator

$$\sigma_f^2(\delta) = \int_0^\infty \langle e^{t\mathcal{L}} f(x), f(x) \rangle_{L^2(dx)} dt = \langle \mathcal{L}_\delta^{-1} f, f \rangle_{L^2(dx)}$$

where  $\mathcal{L}_\delta^{-1}$  is the inverse of  $\mathcal{L}_\delta$  on the orthogonal complement of the eigenfunction 1. Expanding in the eigenfunctions we find

$$\sigma_f^2(\delta) = \sum_{n \in \mathbb{Z}, n \neq 0} \frac{|c_n|^2}{n^2 + in\delta} = \sum_{n=1}^\infty \frac{2|c_n|^2}{n^2 + \delta^2}.$$

Even though the spectral gap does not decrease at all, the variance not only decreases, but it can be made as small as we want by increasing  $\delta^2$ .

## What is known in the literature in terms of spectral gap?

- Spectral gap decreases under a natural non-degeneracy condition on adding some irreversibility The corresponding eigenspace should not be invariant under the action of the added drift  $C$  (Hwang, Hwang-Ma and Sheu, (2005)).
- Let  $U = 0$  and consider a one-parameter family of perturbations  $C = \delta C_0$  for  $\delta \in \mathbb{R}$  and  $C_0$  is some divergence vector field. If the flow is weak-mixing then the second largest eigenvalue tends to 0 as  $\delta \rightarrow \infty$  (Constantin-Kiselev-Ryshik-Zlatoš, (2008)).
- Detailed analysis of linear diffusion processes with  $U(x) = Ax$  and  $C = JAx$  for a antisymmetric  $J$  can be found in Hwang-Ma-Sheu (1993) and Lelievre-Nier-Pavliotis (2012) where the optimal choice of  $J$  is determined.
- Evidence that violation of detailed balance accelerates relaxation in recent physics literature, Ichiki-Ohzeki (2013).

## Part III

What can large deviations theory say?

## Approximation via diffusions

To sample the Gibbs measure  $\bar{\pi}$  on the set  $E$

$$\bar{\pi}(dx) = \frac{e^{-2U(x)} dx}{\int_E e^{-2U(x)} dx}$$

one can consider the (time-reversible) Langevin equation

$$dX_t = -\nabla U(X_t)dt + dW_t. \quad (6)$$

## Approximation via diffusions

To sample the Gibbs measure  $\bar{\pi}$  on the set  $E$

$$\bar{\pi}(dx) = \frac{e^{-2U(x)} dx}{\int_E e^{-2U(x)} dx}$$

one can consider the (time-reversible) Langevin equation

$$dX_t = -\nabla U(X_t) dt + dW_t. \quad (6)$$

There are however many other stochastic differential equations with the same invariant measure and we may consider instead the family of equations

$$dX_t = [-\nabla U(X_t) + C(X_t)] dt + dW_t$$

where the vector field  $C(x)$  satisfies the condition

$$\operatorname{div}(Ce^{-2U}) = 0$$

In this case the Markov process is not time-reversible!! Can we somehow optimize by choosing  $C$ ?

## Donsker-Varadhan and Gärtner theory

Define the empirical measure

$$\pi_t \equiv \frac{1}{t} \int_0^t \delta_{X_s} ds$$

which converges to  $\bar{\pi}$  almost surely. If we have a large deviation for the family of measures  $\pi_t$ , which we write, symbolically as

$$\mathbb{P} \{ \pi_t \approx \mu \} \asymp e^{-tI_C(\mu)}$$

The information in  $I_C(\mu)$  can be used to study observable: we have for  $f \in \mathcal{C}(E; \mathbb{R})$  the large deviation principle

$$\mathbb{P} \left\{ \frac{1}{t} \int_0^t f(X_s) ds \approx \ell \right\} \asymp e^{-t\tilde{I}_{f,C}(\ell)}$$

where

$$\tilde{I}_{f,C}(\ell) = \inf_{\mu \in \mathcal{P}(E)} \{ I_C(\mu) : \langle f, \mu \rangle = \ell \} ,$$



## Donsker-Varadhan and Gärtner theory

In particular, if  $\mathcal{A}$  is the generator of the Markov process and  $\mathcal{D}$  its domain, then the Donsker-Varadhan functional takes the form

$$I(\mu) = - \inf_{u \in \{u \in \mathcal{D}, u > 0\}} \int_E \frac{\mathcal{A}u}{u} d\mu$$

## Donsker-Varadhan and Gärtner theory

In particular, if  $\mathcal{A}$  is the generator of the Markov process and  $\mathcal{D}$  its domain, then the Donsker-Varadhan functional takes the form

$$I(\mu) = - \inf_{u \in \{u \in \mathcal{D}, u > 0\}} \int_E \frac{\mathcal{A}u}{u} d\mu$$

A more explicit formula due to Gärtner:

Theorem (Gärtner).

Consider the SDE

$$dX_t = b(X_t) + dW_t$$

on  $E = \mathbb{T}^d$  with  $b \in \mathcal{C}^1(E; \mathbb{R}^d)$ . The Donsker-Varadhan rate function  $I(\mu)$  takes the form

$$I(\mu) = \frac{1}{2} \int_E |\nabla \phi(x)|^2 d\mu(x) \quad (7)$$

where  $\phi$  is the unique (up to constant) solution of the equation

$$\Delta \phi + \frac{1}{\rho} (\nabla \rho, \nabla \phi) = \frac{1}{\rho} \mathcal{L}^* \rho \quad (8)$$

# Simplifications

In the special case where  $b = -\nabla U$  is a gradient, then  $\phi(x) = \frac{1}{2} \log p(x) + U(x) + \text{constant}$  and we get

$$I(\mu) = \frac{1}{2} \int_E \left| \frac{1}{2} \frac{\nabla p(x)}{p(x)} + \nabla U(x) \right|^2 d\mu(x) \quad (9)$$

which is the usual explicit formula for the rate function in the reversible case.

## Simplifications

In the special case where  $b = -\nabla U$  is a gradient, then  $\phi(x) = \frac{1}{2} \log p(x) + U(x) + \text{constant}$  and we get

$$I(\mu) = \frac{1}{2} \int_E \left| \frac{1}{2} \frac{\nabla p(x)}{p(x)} + \nabla U(x) \right|^2 d\mu(x) \quad (9)$$

which is the usual explicit formula for the rate function in the reversible case.

Motivated, by this if we set  $\phi(x) = \frac{1}{2} \log p(x) + \psi(x)$ , then we get the following representation.

### Lemma

We have

$$I(\mu) = \frac{1}{8} \int_E \left| \frac{\nabla p(x)}{p(x)} \right|^2 d\mu(x) + \frac{1}{2} \int_E |\nabla \psi(x)|^2 d\mu(x) - \frac{1}{2} \int_E \frac{b \nabla p}{p} d\mu(x)$$

where  $\psi$  is the unique (up to constant) solution of the equation

$$\operatorname{div} [p(b + \nabla \psi)] = 0.$$

## Behavior of rate function

Recall that we are comparing

$$dX_t = [-\nabla U(X_t)] dt + dW_t$$

$$dX_t = [-\nabla U(X_t) + C(X_t)] dt + dW_t$$

Same invariant measure but reversible versus irreversible!

### Theorem

Assume that  $C \neq 0$  such that  $\operatorname{div} C = 2C \nabla U$ . For any  $\mu \in \mathcal{P}(E)$  we have  $I_C(\mu) \geq I_0(\mu)$ . If  $\mu(dx) = p(x)dx$  is a measure with positive density  $p \in \mathcal{C}^{(2+\alpha)}(E)$  for some  $\alpha > 0$  and  $\mu \neq \bar{\pi}$  then we have

$$I_C(\mu) = I_0(\mu) + \frac{1}{2} \int_E |\nabla \psi_C(x) - \nabla U(x)|^2 d\mu(x).$$

where  $\psi_C$  is the unique solution (up to a constant) of the equation

$$\operatorname{div} [p(-\nabla U + C + \nabla \psi_C)] = 0.$$

Moreover we have  $I_C(\mu) = I_0(\mu)$  if and only if the positive density  $p(x)$  satisfies  $\operatorname{div}(p(x)C(x)) = 0$ . Equivalently such  $p$  have the form  $p(x) = e^{2G(x)}$  where  $G$  is such that  $G + U$  is an invariant for the vector field  $C$  (i.e.,  $C \nabla(G + U) = 0$ ).

## Behavior of rate function

To obtain a slightly more quantitative result let us consider a one-parameter family  $C(x) = \delta C_0(x)$  where  $\delta \in \mathbb{R}$  and  $C_0 \neq 0$  such that  $\operatorname{div} C_0 = 2C_0 \nabla U$ .

### Theorem

Assume that  $C_0 \neq 0$  such that  $\operatorname{div} C_0 = 2C_0 \nabla U$ . Consider the measure  $\mu(dx) = p(x)dx$  with positive density  $p \in \mathcal{C}^{(2+\alpha)}(E)$  for some  $\alpha > 0$ . Then we

$$I_{\delta C_0}(\mu) = I_0(\mu) + \delta^2 K(\mu).$$

where the functional  $K(\mu)$  is strictly positive if and only if  $\operatorname{div}(p(x)C(x)) \neq 0$ .

Namely, rate function is quadratic in  $\delta$ !

## What about observables?

For  $f \in \mathcal{C}(E, \mathbb{R})$ , by the contraction principle,

$$\mathbb{P} \left\{ \frac{1}{t} \int_0^t f(X_s) ds \approx \ell \right\} \asymp e^{-t\tilde{I}_{f,C}(\ell)}$$

where

$$\tilde{I}_{f,C}(\ell) = \inf_{\mu \in \mathcal{P}(E)} \{I_C(\mu) : \langle f, \mu \rangle = \ell\}$$

### Theorem

Consider  $f \in \mathcal{C}^{(\alpha)}(E)$  and  $\ell \in (\min_x f(x), \max_x f(x))$  with  $\ell \neq \int f d\bar{\pi}$ . Fix a vector field  $C$  as in assumption **(H)**. Then we have

$$\tilde{I}_{f,C}(\ell) \geq \tilde{I}_{f,0}(\ell).$$

Moreover if there exists  $\ell_0$  such that for this particular field  $C$ ,  $\tilde{I}_{f,C}(\ell_0) = \tilde{I}_{f,0}(\ell_0)$  then we must have

$$\widehat{\beta}(\ell_0)f = \frac{1}{2}\Delta(G + U) + \frac{1}{2}|\nabla G|^2 - \frac{1}{2}|\nabla U|^2, \quad (10)$$

where  $G$  is such that  $G + U$  is invariant under the particular vector field  $C$ .

## What about observables?

- Letting  $\mathcal{L}_0$  denote the infinitesimal generator of the reversible process  $X_t$  (i.e., when  $C = 0$ ), we get that (10) can be rewritten as a nonlinear Poisson equation of the form

$$\mathcal{H}(G + U) = \widehat{\beta}(\ell_0)f, \quad (11)$$

where

$$\mathcal{H}(G + U) = e^{-(G+U)}\mathcal{L}_0e^{G+U} = \frac{1}{2}\Delta(G + U) + \frac{1}{2}|\nabla G|^2 - \frac{1}{2}|\nabla U|^2.$$

- This result does not necessarily mean that there are observables  $f$ , for which variance reduction cannot be attained. It only means, that for a given observable, one should choose a vector field  $C$ , such that there is no  $G$  that satisfies both  $C\nabla(G + U) = 0$  and (10).



# Sketch of the proof 1/5.

Overview of the proof:

- 1 Recall that we already know that  $I_C(\mu) > I_0(\mu)$ .
- 2 Since  $\tilde{I}_{f,C}(\ell) = \inf_{\mu \in \mathcal{P}(E)} \{I_C(\mu) : \langle f, \mu \rangle = \ell\}$ , is the minimizer  $\mu$  achieved?
- 3 The mimimizer  $\mu$  is achieved for “good” functions  $f$ .
- 4 Exploiting these properties, a contradiction argument gives us that  $\tilde{I}_{f,C}(\ell) \geq \tilde{I}_{f,0}(\ell)$  and the condition under which the inequality is strict.

## Sketch of the proof 2/5.

The proof of the **existence and characterization of the minimizer**  $\mu$  is based on the representation of the rate function in terms of a Legendre transform

$$\tilde{I}_{f,C}(\ell) = \sup_{\beta \in \mathbb{R}} \{ \beta \ell - \lambda(\beta f) \} .$$

where the eigenvalue  $\lambda(\beta f)$  is a smooth strictly convex function of  $\beta$

$$\lambda(\beta f) = \lim_{t \rightarrow \infty} \frac{1}{t} \log \mathbb{E} \left[ e^{\int_0^t \beta f(X_s) ds} \right] .$$

If  $\ell$  belongs to the range of  $f$  we have

$$\tilde{I}_{f,C}(\ell) = \hat{\beta} \ell - \lambda(\hat{\beta} f), \quad \text{with } \hat{\beta} \text{ given by } \ell = \frac{d}{d\beta} \lambda(\hat{\beta} f).$$

## Sketch of the proof 3/5.

Since  $f \in \mathcal{C}^{(\alpha)}$ ,  $\lambda(\beta f)$  is the maximal eigenvalue of  $\mathcal{L}_C + \beta f$

$$(\mathcal{L}_C + \beta f)u(\beta f) = \lambda(\beta f)u(\beta f), \quad (12)$$

and is a smooth convex function of  $\beta$ . Here  $u(\beta f)$  is the corresponding eigenfunction.

With  $u(\beta f) = e^{\phi(\beta f)}$ , the eigenvalue equation can be equivalently written as

$$\mathcal{L}_C \phi(\beta f) + \frac{1}{2} |\nabla \phi(\beta f)|^2 = \lambda(\beta f) - \beta f \quad (13)$$

Differentiating with respect to  $\beta$  and setting  $\psi(\beta f) = \frac{\partial \phi}{\partial \beta}(\beta f)$  we see that  $\psi(\beta f)$  satisfies the equation

$$\mathcal{L}_C \psi(\beta f) + (\nabla \phi(\beta f), \nabla \psi(\beta f)) = \lambda'_f(\beta) - f$$

## Sketch of the proof 4/5.

Equivalently

$$\mathcal{L}_{C+\nabla\phi(\beta f)}\psi = \lambda'_f(\beta) - f$$

Thus, the constraint  $\langle f, \mu \rangle = \ell$ , implies that in order to have  $\ell = \lambda'_f(\hat{\beta})$  for some  $\hat{\beta}$ ,  $\mu_{\hat{\beta}}$  should be the invariant measure for the process with generator  $\mathcal{L}_{C+\nabla\phi(\hat{\beta}f)}$ .

Since  $\nabla\phi \in \mathcal{C}^{(1+\alpha)}$  the corresponding invariant measure  $\mu_{\hat{\beta}}$  is strictly positive and has a density  $p(x) \in \mathcal{C}^{(2+\alpha)}$ .

To conclude the proof, by Gärtner's result we have

$I_C(\mu_{\hat{\beta}}) = \mu(\hat{\beta}f) - \lambda(\hat{\beta}f)$ . But since  $\mu(f) = \ell$  this is also equal to  $I_{f,C}(\ell)$ .

## Sketch of the proof 5/5.

Hence, we have obtained

$$\tilde{I}_{f,C}(\ell) = I_C(\mu_{C,\hat{\beta}}) \text{ with } \mu_{C,\hat{\beta}}(dx) = p_{C,\hat{\beta}}(x)dx.$$

The rest is standard.

- 1 Let  $\text{div}(Cp_{C,\hat{\beta}}) \neq 0$ . Assume that the rate functions with  $C = 0$  and  $C \neq 0$  are equal and get a contradiction.
- 2 Let  $\text{div}(Cp_{C,\hat{\beta}}) = 0$ . Let us write  $p_{C,\hat{\beta}} = e^{-2G}$ , so we must have  $C \cdot \nabla G = 2\text{div}(C)$ . Keeping in mind that  $p_{C,\hat{\beta}}(x)$  is invariant density corresponding to a known operator gives us that we must have  $p_{C,\hat{\beta}} = p_{0,\hat{\beta}} = e^{2(\phi(\beta f) - U) + \text{const}}$ . Thus  $\phi(\beta f) = G + U$  and  $C \cdot \nabla \phi(\beta f) = 2\text{div}(C)$  and (13) reduces to  $(\mathcal{L}_0 + \hat{\beta}f)e^\phi = \lambda(\hat{\beta}f)e^\phi$ . Solving for  $f$  gives the nonlinear Poisson equation (11):

$$\hat{\beta}f = e^{-\phi} \mathcal{L}_0 e^\phi + \text{const.} \quad (14)$$

## Part IV

What about variance reduction?

## Asymptotic variance

Under our assumptions the central limit theorem holds for the ergodic average  $f_t$  and we have

$$t^{1/2} \left( \frac{1}{t} \int_0^t f(X_s) ds - \int f d\bar{\pi} \right) \Rightarrow N(0, \sigma_f^2)$$

and the asymptotic variance  $\sigma_f^2$  is given in terms of the integrated autocorrelation function,

$$\sigma_f^2 = 2 \int_0^\infty \mathbb{E}_{\bar{\pi}} [(f(X_0) - \bar{f})(f(X_t) - \bar{f})] dt$$

## Asymptotic variance

This is a convenient quantity from a practical point of view since there exists easily implementable estimators for  $\sigma_{\bar{f}}^2$ . On the other hand the asymptotic variance  $\sigma_{\bar{f}}^2$  is related to the curvature of the rate function  $I_f(\ell)$  around the mean  $\bar{f}$  we have

$$\tilde{I}_f''(\bar{f}) = \frac{1}{2\sigma_{\bar{f}}^2}.$$

From previous theorem it follows immediately that  $\sigma_{\bar{f},C}^2 \leq \sigma_{\bar{f},0}^2$  but in fact the addition of an irreversible drift strictly generically decreases the asymptotic variance.

### Theorem

Assume that  $C \neq 0$  is a vector field such that  $\operatorname{div} C = 2C \nabla U$ . Let  $f \in \mathcal{C}^{(\alpha)}(E)$  such that for some  $\epsilon > 0$  and  $\ell \in (\bar{f} - \epsilon, \bar{f} + \epsilon) \setminus \{\bar{f}\}$  we have  $\tilde{I}_{f,C}(\ell) > \tilde{I}_{f,0}(\ell)$ . Then we have

$$\sigma_{\bar{f},C}^2 < \sigma_{\bar{f},0}^2.$$



## Sketch of the proof

It is clear that the relation  $\sigma_{\bar{f}}^2 = \frac{1}{2\tilde{I}_{\bar{f}}''(\bar{f})}$  implies that it is enough to prove that for  $C \neq 0$  and  $f \in \mathcal{C}^{(\alpha)}(E)$

$$\tilde{I}_{\bar{f},C}''(\bar{f}) - \tilde{I}_{\bar{f},0}''(\bar{f}) > 0$$

The proof of this statement follows by precise computation of first and then second order Gâteaux derivatives.

## Sketch of the proof

It is clear that the relation  $\sigma_{\bar{f}}^2 = \frac{1}{2\tilde{I}_{\bar{f}}''(\bar{f})}$  implies that it is enough to prove that for  $C \neq 0$  and  $f \in \mathcal{C}^{(\alpha)}(E)$

$$\tilde{I}_{\bar{f},C}''(\bar{f}) - \tilde{I}_{\bar{f},0}''(\bar{f}) > 0$$

The proof of this statement follows by precise computation of first and then second order Gâteaux derivatives.

The formula of the second order derivative  $\tilde{I}_{\bar{f},C}''(\bar{f})$  that is derived, provides a natural optimization problem for the choice of  $C$ .

# Part V

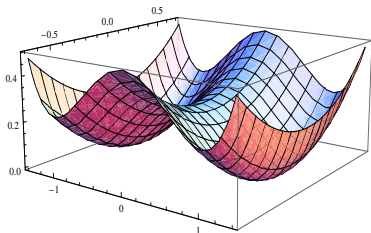
## Simulation results

## Example.

Consider that we want to sample from the stationary distribution

$$\bar{\pi}(dxdy) = \frac{e^{-\frac{U(x,y)}{D}}}{\int_{\mathbb{R}^2} e^{-\frac{U(x,y)}{D}} dxdy} dxdy$$

where  $D$  is some constant and  $U(x, y) = \frac{1}{4}(x^2 - 1)^2 + \frac{1}{2}y^2$ ,



## Example.

- Consider the Markov process

$$dZ_t = [-\nabla U(Z_t) + C(Z_t)] dt + \sqrt{2D} dW_t, \quad Z_0 = 0$$

where for  $z = (x, y)$ ,  $U(x, y) = \frac{1}{4}(x^2 - 1)^2 + \frac{1}{2}y^2$ .

- Let  $D = 0.1$  and  $C(x, y) = \delta C_0(x, y)$  with  $C_0(x, y) = J \nabla U(x, y)$ . Here,  $\delta \in \mathbb{R}$ ,  $I$  is the  $2 \times 2$  identity matrix and  $J$  is the standard  $2 \times 2$  antisymmetric matrix, i.e.,  $J_{12} = 1$  and  $J_{21} = -1$ .
- Notice that for any  $\delta \in \mathbb{R}$ , the invariant measure is

$$\bar{\pi}(dx dy) = \frac{e^{-\frac{U(x,y)}{D}}}{\int_{\mathbb{R}^2} e^{-\frac{U(x,y)}{D}} dx dy} dx dy$$

## Example.

Let us suppose that we want to compute the following observables

$$\bar{f}_i = \int_{\mathbb{R}^2} f_i(x, y) \bar{\pi}(dx dy), \quad i = 1, 2$$

where

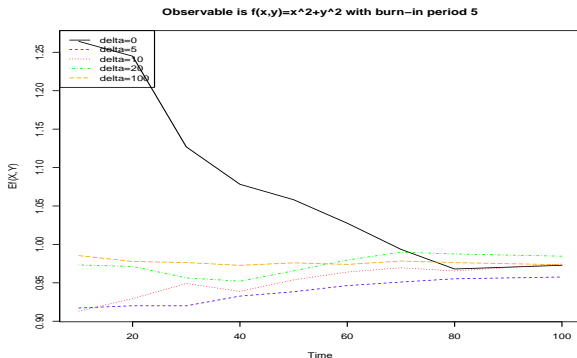
$$f_1(x, y) = x^2 + y^2, \quad f_2(x, y) = U(x, y) = \frac{1}{4}(x^2 - 1)^2 + \frac{1}{2}y^2$$

## Speed of convergence

It is known that an estimator for  $\bar{f}_i$  is given by

$$\hat{\bar{f}}_i(t) = \frac{1}{t - \nu} \int_{\nu}^t f_i(X_s, Y_s) ds$$

where  $\nu$  is some burn-in period that is used with the hope that the bias has been significantly reduced by time  $\nu$ .



## Variance reduction

In general, a central limit theorem holds and takes the following form

$$t^{1/2} \left( \hat{f}(t) - \bar{f} \right) \Rightarrow N(0, \sigma_f^2)$$

In order to estimate  $\sigma_f^2$ , we use the well established method of batch means  
Then for  $\kappa = 1, \dots, m$  ( $m$  is number of batches) we define

$$\hat{f}(t; \kappa) = \frac{1}{t/m} \int_{(\kappa-1)t/m}^{\kappa t/m} f(X_s, Y_s) ds,$$

$$\hat{f}(t) = \frac{1}{m} \sum_{\kappa=1}^m \hat{f}(t; \kappa)$$

and

$$s_m^2(t) = \frac{1}{m-1} \sum_{\kappa=1}^m \left( \hat{f}(t; \kappa) - \hat{f}(t) \right)^2$$



## Variance reduction

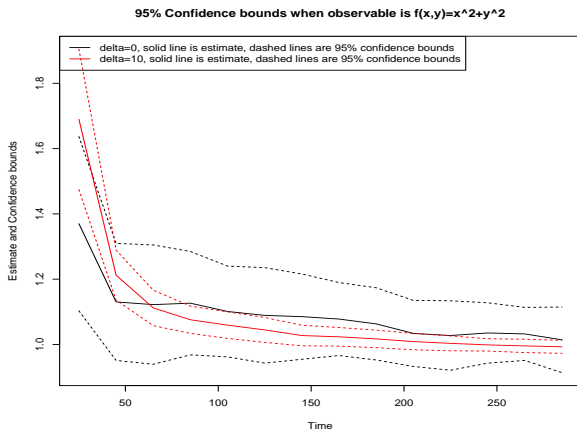
Then, we have in distribution

$$\sqrt{m} \frac{\hat{f}(t) - \bar{f}}{s_m(t)} \Rightarrow T_{m-1}, \quad \text{as } t \rightarrow \infty$$

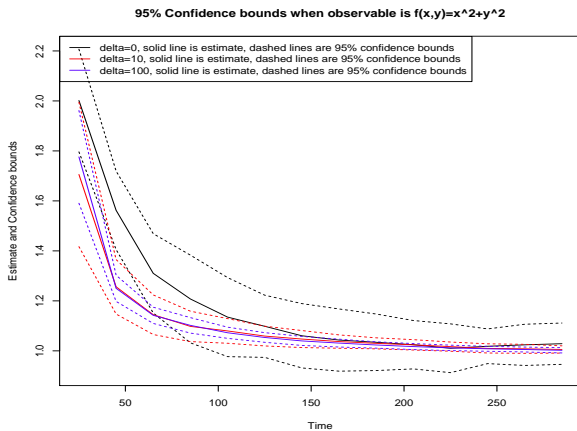
where  $T_{m-1}$  is the Student's T distribution with  $m - 1$  degrees of freedom. So, a  $(1 - \alpha)\%$  confidence interval is given by

$$\left( \hat{f}(t) - t_{\alpha/2, m-1} s_m(t) / \sqrt{m}, \hat{f}(t) + t_{\alpha/2, m-1} s_m(t) / \sqrt{m} \right)$$

# Variance reduction



# Variance reduction



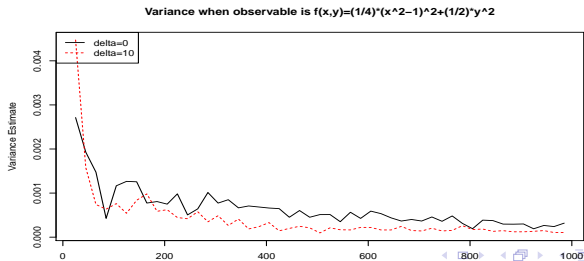
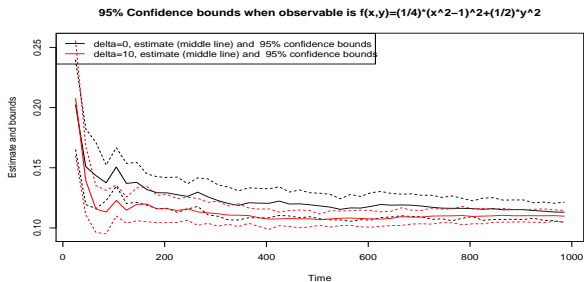
## Variance reduction

$\delta \mid t$	25	100	160	220	295
0	0.22	0.08	0.038	0.029	0.011
10	0.19	0.01	0.007	0.005	0.002
100	0.09	0.001	$3e - 04$	$2.8e - 04$	$1.3e - 04$

**Table:** Estimated variance values for different pairs  $(\delta, t)$ .

- Variance reduction of about two orders of magnitude when  $\delta = 100$  versus  $\delta = 0$ !!
- However, for accuracy purposes, larger  $\delta$  requires smaller discretization step in the numerical algorithm. Hence, there is trade-off to consider here.

# Variance reduction



# Part VI

## Summary

To summarize, we have developed a systematic approach to the problem.

- 1 Large deviations action functional is directly related to the empirical measure.
  - ▶ Can be used a measure of performance. The bigger it is the better.
- 2 Asymptotic variance is inversely proportional to the second derivative of the rate function!
- 3 A natural optimization problem for the choice of the optimal perturbation is being defined.
- 4 Introducing irreversibility speeds up convergence, and reduces significantly the variance of the estimator.

# Part VII

## References



# References. I

- 1 K.A. Athreya, H. Doss and J. Sethuraman, On the convergence of the Markov chain simulation method, *Annals of Statistics*, Vol. 24, (1996), pp. 69-100.
- 2 C. Barbarosie, Representation of divergence-free vector fields, *Quarterly of Applied Mathematics*, Vol. 69, (2011), pp. 309–316.
- 3 M. Bedard and J.S. Rosenthal, Optimal Scaling of Metropolis Algorithms: Heading Towards General Target Distributions, *Canadian Journal of Statistics*, Vol. 36, Issue 4, (2008), pp. 483-503.
- 4 P. Constantin, A. Kiselev, L. Ryshik and A. Zlatoš, Diffusion and mixing in fluid flow *Annals of Mathematics*, Vo. 168 (2008), pp. 643-674

## References. II

- 5 M.D. Donsker and S.R.S. Varadhan, Asymptotic evaluation of certain Markov process expectations for large times, I, *Communications Pure in Applied Mathematics*, Vol. 28, (1975), pp. 1-47, II, *Communications on Pure in Applied Mathematics*, Vol. 28, (1975), pp. 279–301, and III, *Communications on Pure in Applied Mathematics*, Vol. 29, (1976), pp. 389-461.
- 6 P. Dupuis, Y. Liu, N. Plattner, and J. D. Doll, On the Infinite Swapping Limit for Parallel Tempering. *SIAM Multiscale Modeling and Simulation*, Vol. 10, Issue 3, (2012), pp. 986-1022.
- 7 B. Franke, C.-R. Hwang, H.-M. Pai, and S.-J. Sheu, The behavior of the spectral gap under growing drift, *Transactions of the American Mathematical Society*, Vol 362, No. 3 (2010), pp. 1325-1350.
- 8 A. Frigessi, C.R. Hwang and L. Younes, Optimal spectral structures of reversible stochastic matrices, Monte Carlo methods and the simulation of Markov random fields, *Annals of Applied Probability*, Vol. 2, (1992), pp. 610-628.

## References. III

- 9 A. Frigessi, C.R. Hwang, S.J. Sheu and P. Di Stefano, Convergence rates of the Gibbs sampler, the Metropolis algorithm, and their single-site updating dynamics, *Journal of Royal Statistical Society Series B, Statistical Methodology*, Vol. 55, (1993), pp. 205-219.
- 10 A. Ichiki and M. Ohzeki, Violation of detailed balance condition accelerates relaxation, *Physical Review E*, Vol. 88, (2013), pp. 020101(R).
- 11 Jürgen Gärtner, On large deviations from the invariant measure, *Theory of probability and its applications*, Vol. XXII, No. 1, (1977), pp. 24-39.
- 12 W.R. Gilks and G.O. Roberts, Strategies for improving MCMC, *Monte Carlo Markov Chain in practice*, Chapman and Hall, Boca Raton, FL, (1996), pp. 89-114.
- 13 C.R. Hwang, S.Y. Hwang-Ma and S.J. Sheu, Accelerating Gaussian diffusions. *The Annals of Applied Probability* Vol. 3, (1993) pp. 897-913.
- 14 C.R. Hwang, S.Y. Hwang-Ma and S.J. Sheu, Accelerating diffusions, *The Annals of Applied Probability*, Vol 15, No. 2, (2005), pp. 1433-1444.

## References. IV

- 15 T. Lelièvre, F. Nier and G.A. Pavliotis, Optimal non-reversible linear drift for the convergence to equilibrium of a diffusion, *Journal of Statistical Physics*, 152(2), 237-274, (2013)
- 16 K.L. Mergessen and R.L. Tweedie, Rates of convergence of the Hastings and Metropolis algorithms, *Annals of Statistics*, Vol. 24, (1996), pp. 101-121.
- 17 R. Pinsky, The I-function for diffusion processes with boundaries, *The Annals of Probability*, Vol. 13, No. 3, (1985), pp. 676-692.
- 18 G.O. Roberts and J.S. Rosenthal, General state space Markov Chain and MCMC algorithms, *Probability Surveys*, Vol. 1, (2004), pp. 20-71.
- 19 Luc Rey-Bellet and K. Spiliopoulos, Irreversible Langevin samplers and variance reduction: a large deviations approach, (2014), submitted.
- 20 Luc Rey-Bellet and K. Spiliopoulos, Variance reduction for irreversible samplers, (2014), submitted.

**Thank You!!!!**