

Borne de convergence explicite pour l'algorithme Langevin Monte Carlo cinétique.

Aurélien Enfroy¹

Joint work with: Alain Durmus², Eric Moulines²,

¹Télécom SudParis

²ENS Paris-Saclay

³Ecole Polytechnique

Outline

- 1 Motivation and setting
- 2 The Kinetic Langevin algorithm
- 3 Bound distance between semi-group and Discretization
- 4 Semi-group convergence
- 5 Convergence

Introduction

- Sampling distribution over high-dimensional state-space has recently attracted a lot of research efforts in computational statistics and machine learning community...
- **Applications** (non-exhaustive)
 1. Bayesian inference for high-dimensional models,
 2. Bayesian inverse problems (e.g., image restoration and deblurring),
 3. Aggregation of estimators and experts,
 4. Bayesian non-parametrics.
- Most of the sampling techniques known so far **do not scale** to high-dimension... Challenges are numerous in this area...

Bayesian setting

- A Bayesian model is specified by
 1. the sampling distribution of the observed data conditional on its parameters, often termed **likelihood**: $Y \sim L(\cdot|\theta)$
 2. a **prior distribution** π_0 on the parameter space $\theta \in \mathbb{R}^d$
- The inference is based on the **posterior distribution**:

$$\pi(d\theta) = \frac{\pi_0(d\theta)L(Y|\theta)}{\int L(Y|u)\pi_0(du)}.$$

- In most cases the normalizing constant is **not tractable**:

$$\pi(d\theta) \propto \pi_0(d\theta)L(Y|\theta).$$

Bayesian setting

- Bayesian decision theory relies on computing expectations:

$$\pi(f) = \int_{\mathbb{R}^d} f(x) d\pi(x) = \int_{\mathbb{R}^d} f(x) \pi(x) dx$$

Generic problem: estimation of an integral $\pi(f)$, where

- π is known up to a multiplicative factor ;
- Sampling directly from π is not an option;
- A solution is to approximate $\pi(f)$ by

$$n^{-1} \sum_{i=1}^n f(X_i) ,$$

where $(X_i)_{i \geq 0}$ is a Markov chain associated with a Markov kernel P with invariant distribution π .

Markov chain theory

Let $(X_k)_{k \geq 0}$ be a Markov chain on $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$.

- **P Markov kernel** associated with $(X_k)_{k \geq 0}$ if
 - for any ν , νP is the distribution of X_1 starting from $X_0 \sim \nu$
 - νP^k the distribution of X_k for $k \geq 0$
- **Invariant probability measure**: π is said to be an invariant probability measure for the Markov kernel P if

$$X_0 \sim \pi \text{ then } X_1 \sim \pi \quad \text{equivalent to } \pi P = \pi$$

- **Ergodic Theorem** (Meyn and Tweedie, 2003): If π is invariant, With some conditions on P , we have for any $f \in L^1(\pi)$,

$$\frac{1}{n} \sum_{i=1}^n f(X_i) \xrightarrow{\pi\text{-a.s.}} \int f(x) \pi(x) dx.$$

Convergence of Markov chains

- A measure of efficiency of MCMC to target π associated to a Markov kernel P :

$$\|P^k(x, \cdot) - \pi\|_{\text{TV}} \leq C(x)v(k) ,$$

1. The total variation distance defined for μ, ν two probability measures on \mathbb{R}^d by

$$\|\mu - \nu\|_{\text{TV}} = \sup_{|f| \leq 1} |\mu(f) - \nu(f)| .$$

2. $C(x) \geq 0$: dependence on the initial condition.
3. Ideally $\lim_{k \rightarrow +\infty} v(k) = 0$ (or close to 0) with the better possible rate.

We answer to the following questions:

- For a target precision $\varepsilon > 0$, we can find $N \geq 0$ such that

$$\|\delta_x P^n - \pi\|_{\text{TV}} \leq \varepsilon \text{ for all } n \geq N .$$

- In general N is not explicit.

Outline

- 1 Motivation and setting
- 2 The Kinetic Langevin algorithm
- 3 Bound distance between semi-group and Discretization
- 4 Semi-group convergence
- 5 Convergence

Framework

- Denote by π a target density w.r.t. the Lebesgue measure on \mathbb{R}^d , known up to a normalisation factor

$$x \mapsto e^{-U(x)} / \int_{\mathbb{R}^d} e^{-U(y)} dy ,$$

- Assume for the moment that U is L -smooth : continuously differentiable and there exists a constant L such that for all $x, y \in \mathbb{R}^d$,

$$\|\nabla U(x) - \nabla U(y)\| \leq L\|x - y\| .$$

Kinetic Langevin diffusion

- Kinetic Langevin SDE:

$$d\mathbf{X}_t = \mathbf{V}_t dt ,$$

$$d\mathbf{V}_t = -(\kappa_1 \mathbf{V}_t + \kappa_2 \nabla U(\mathbf{X}_t))dt + \sqrt{2\kappa_1\kappa_2}dB_t ,$$

where $(B_t)_{t \geq 0}$ is a d -dimensional Brownian Motion and κ_1, κ_2 are positive constants.

- **Notation:** $(P_t)_{t \geq 0}$ the Markov semigroup associated to the Kinetic Langevin diffusion:

$$P_t(z, A) = \mathbb{P}((\mathbf{X}_t, \mathbf{V}_t) \in A | (\mathbf{X}_0, \mathbf{V}_0) = z) , \quad x \in \mathbb{R}^{2d}, A \in \mathcal{B}(\mathbb{R}^{2d}) .$$

- $\mu(x, v) \propto \exp(-U(x)) \exp(-\|v\|^2/2) \propto \pi(x) \exp(-\|v\|^2/2)$ is the unique **invariant probability** measure.

Discretized Kinetic Langevin diffusion

- **Idea:** Sample the diffusion paths:

$$\begin{aligned}d\tilde{\mathbf{X}}_t &= \tilde{\mathbf{V}}_t dt, \\d\tilde{\mathbf{V}}_t &= -(\kappa_1 \tilde{\mathbf{V}}_t + \kappa_2 \nabla U(\tilde{\mathbf{X}}_{\Gamma_k}))dt + \sqrt{2\kappa_1\kappa_2}dB_t,\end{aligned}$$

where

- $(\gamma_k)_{k \geq 1}$ is a sequence of stepsizes, which can either be held constant or be chosen to decrease to 0 at a certain rate
- For any $k \in \mathbb{N}$, $\Gamma_k = \sum_{i=1}^k \gamma_i$
- **Noation:** For any $k \in \mathbb{N}$, $(X_k, V_k) = (\mathbf{X}_{\Gamma_k}, \mathbf{V}_{\Gamma_k})$
- This algorithm is referred to as the **Kinetic Langevin Algorithm**.

Explicit form of the Kinetic Langevin Algorithm

- Explicit form:

$$\begin{aligned} X_{k+1} &= X_k - \frac{1}{\kappa_1} (e^{-\kappa_1 \gamma_{k+1}} - 1) V_k - \frac{\kappa_2}{\kappa_1} \left(\gamma_{k+1} + \frac{e^{-\kappa_1 \gamma_{k+1}} - 1}{\kappa_1} \right) \nabla U(X_k) \\ &\quad - p_{\mathbb{R}^d \times \{0\}} \left(\Sigma_{\gamma_{k+1}}^{1/2} G_{k+1} \right), \\ V_{k+1} &= e^{-\kappa_1 \gamma_{k+1}} V_k + \frac{\kappa_2}{\kappa_1} (e^{-\kappa_1 \gamma_{k+1}} - 1) \nabla U(X_k) + p_{\{0\} \times \mathbb{R}^d} \left(\Sigma_{\gamma_{k+1}}^{1/2} G_{k+1} \right), \end{aligned}$$

where

- $(G_k)_{k \geq 1}$ is i.i.d. $\mathcal{N}(0, I_{2d})$
- $(\gamma_k)_{k \geq 1}$ is a sequence of stepsizes, which can either be held constant or be chosen to decrease to 0 at a certain rate
-

$$\Sigma_\gamma = 2\kappa_1^{-1}\kappa_2 \begin{pmatrix} \frac{\kappa_1 \gamma - 2^{-1}(1 - e^{-\kappa_1 \gamma})^2 - (1 - e^{-\kappa_1 \gamma})}{\kappa_1} I_d & -2^{-1}(1 - e^{-\kappa_1 \gamma})^2 I_d \\ -2^{-1}(1 - e^{-\kappa_1 \gamma})^2 I_d & 2^{-1}\kappa_1(1 - e^{-2\kappa_1 \gamma}) I_d \end{pmatrix}.$$

Discretized Kinetic Langevin diffusion: constant stepsize

- When the stepsize is held **constant**, i.e. $\gamma_k = \gamma$, then $(X_k, V_k)_{k \geq 1}$ is an **homogeneous Markov chain** with Markov kernel R_γ
- Under some appropriate conditions,

$$R_\gamma$$

is irreducible, positive recurrent \rightsquigarrow unique invariant distribution μ_γ which **does not coincide** with the target distribution μ .

- **Questions:**
 - For a given precision $\epsilon > 0$, how should I choose the stepsize $\gamma > 0$ and the number of iterations n so that : $\|\delta_x R_\gamma^n - \mu\|_{TV} \leq \epsilon$
 - quantify the distance between μ_γ and μ .

Discretized Kinetic Langevin diffusion: decreasing stepsize

- When $(\gamma_k)_{k \geq 1}$ is nonincreasing and non constant, $(X_k, V_k)_{k \geq 1}$ is an **inhomogeneous Markov chain** associated with the kernels $(R_{\gamma_k})_{k \geq 1}$.
- **Notation:** Q_γ^p is the composition of Markov kernels

$$Q_\gamma^p = R_{\gamma_1} R_{\gamma_2} \dots R_{\gamma_p}$$

With this notation, $\mathbb{E}_x[f(X_p, V_p)] = \delta_x Q_\gamma^p f$.

- **Questions:**
 - **Convergence** : is there a way to choose the step sizes so that $\|\delta_x Q_\gamma^p - \mu\|_{TV} \rightarrow 0$?

Outline

- 1 Motivation and setting
- 2 The Kinetic Langevin algorithm
- 3 Bound distance between semi-group and Discretization**
- 4 Semi-group convergence
- 5 Convergence

Kullback-Leibler divergence

- For $f : \mathbb{R}^d \rightarrow \mathbb{R}$, such that for any $x \in \mathbb{R}^d$, $f(x) > 0$, $\int_{\mathbb{R}^d} f(x) d\mu(x) < +\infty$, define the entropy of f with respect to μ by

$$\text{Ent}_\mu(f) = \int_{\mathbb{R}^d} f(x) \log(f(x)) d\mu(x) - \log\left(\int_{\mathbb{R}^d} f(x) d\mu(x)\right) \int_{\mathbb{R}^d} f(x) d\mu(x) .$$

- The Kullback-Leibler divergence between μ and ν is defined by

$$\text{KL}(\nu, \mu) = \text{Ent}_\mu(d\nu/d\mu) ,$$

if $\nu \ll \mu$ and $\text{KL}(\mu, \nu) = +\infty$ otherwise.

Main result

Theorem 1

For any $k \in \mathbb{N}^*$ and $z \in \mathbb{R}^{2d}$,

$$\text{KL}(\delta_z P_{\Gamma_k} \mid \delta_z Q_\gamma^k) \leq (L^2 \kappa_2 / (2\kappa_1)) \sum_{j=0}^{k-1} \gamma_{j+1}^3 D(\gamma_{j+1}, \delta_z Q_\gamma^j),$$

where for any distribution ν on $\mathcal{B}(\mathbb{R}^d)$ and $\gamma \in \mathbb{R}_+^*$,

$$D(\gamma, \nu) = \int_{\mathbb{R}^{2d}} \{A(\gamma) \|x'\|^2 + B(\gamma) \|v'\|^2\} d\nu(dx' dv') + d\kappa_1 \kappa_2 \gamma / 4,$$
$$A(\gamma) = L^2 \left(\frac{\kappa_2^2 \gamma^2}{15} + \frac{\kappa_2 \gamma}{8} \right), \quad B(\gamma) = \frac{1}{3} + \frac{\kappa_2 \gamma}{8}.$$

Main result

Corollary 2

For any $z \in \mathbb{R}^{2d}$, $k \in \mathbb{N}^*$, $p \in \mathbb{N}$, $p < k$,

$$\|\delta_z Q_\gamma^k - \mu\|_{TV} \leq L \sqrt{\frac{\kappa_2}{\kappa_1} \left(\sum_{j=p}^{k-1} \gamma_{j+1}^3 D(\gamma_{j+1}, \delta_z Q_\gamma^j) \right)} + \|\delta_z Q_\gamma^p P_{\Gamma_{p,k}} - \mu\|_{TV} .$$

Outline

- 1 Motivation and setting
- 2 The Kinetic Langevin algorithm
- 3 Bound distance between semi-group and Discretization
- 4 Semi-group convergence**
- 5 Convergence

Poincaré

- **Objective** compute bound for $\|\delta_z Q_\gamma^p P_{\Gamma_{p,k}} - \mu\|_{TV}$
- **Assumption:** π satisfies a Poincaré inequality i.e there exists C_P such that, for any $f \in C^2(\mathbb{R}^d)$ satisfying $\int_{\mathbb{R}^d} f d\pi = 0$, $\int_{\mathbb{R}^d} (\nabla f)^2 d\pi \geq C_P \int_{\mathbb{R}^d} f^2 d\pi$.
- **Remark:** If U is such that $\langle \nabla U(x), x \rangle \geq L\chi_2 \|x\| - \tau_2$ then π satisfies a Poincaré inequality.

Theorem 3

For any $t \in \mathbb{R}_+$ and initial distribution ν_0 such that $\nu_0 \ll \mu$,

$$\|\nu_0 P_t - \mu\|_{TV} \leq A e^{-\alpha t} \text{Var}_\mu(d\nu_0/d\mu)$$

where A and α are explicit constants.

- Assumption:

- π satisfies a log-Sobolev inequality i.e. there exists $C_{LS} > 0$ such that, for any continuously differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$ satisfying $f(x) > 0$ for any $x \in \mathbb{R}^d$, $\int_{\mathbb{R}^d} f(x) dx < +\infty$

$$\text{Ent}_{\pi}(f) \leq C_{LS} \int_{\mathbb{R}^d} \{\|\nabla f(x)\|^2 / f(x)\} d\pi(x).$$

- The potential U is infinitely continuously differentiable on \mathbb{R}^d and for any $\alpha \in \mathbb{N}^d$ multi-index, $\sup_{x \in \mathbb{R}^d} |\partial_x^\alpha U(x)| < +\infty$.

- Remark: If U is strongly convex then π satisfies a log-Sobolev inequality.

Theorem 4

For any distribution ν_0 such that $d\nu_0/d\mu \in C_b^{2,+,*}(\mathbb{R}^{2d})$ and $t \in \mathbb{R}_+$,

$$\text{KL}(\nu_0 P_t, \mu) \leq \exp(-\alpha t) \left(\text{KL}(\nu_0, \mu) + \frac{1}{\beta} \int \Phi_2 \left(\frac{d\nu_0}{d\mu} \right) d\mu \right),$$

where α is an explicit constant and Φ_2, β are defined by,

$$\Phi_2(f) = \left(\|(\nabla_x + \kappa_1 \nabla_v) f\|^2 + \|\nabla_v f\|^2 \right) / f, \quad \beta = 2(1 + \kappa_2^2 L^2 + \kappa_1^2/2) / (\kappa_1^2 \kappa_2).$$

Outline

- 1 Motivation and setting
- 2 The Kinetic Langevin algorithm
- 3 Bound distance between semi-group and Discretization
- 4 Semi-group convergence
- 5 Convergence**

Convergence

- By the previous section there exist $\varrho \in (0, 1)$, such that for any initial distribution ν_0 , there exists $C(\nu_0) < +\infty$ such that,

$$\|\nu_0 P_t - \mu\|_{TV} \leq C(\nu_0) \varrho^t$$

Theorem 5

Assume that $\lim_{k \rightarrow +\infty} \gamma_k = 0$, $\lim_{k \rightarrow +\infty} \Gamma_k = +\infty$ and some technical conditions. Then for any $z \in \mathbb{R}^{2d}$

$$\lim_{k \rightarrow +\infty} \|\delta_z Q_\gamma^k - \mu\|_{TV} = 0$$

Algorithm complexity

Theorem 6

For all $\varepsilon > 0$, we get $\|\delta_z R_\gamma^k - \mu\|_{TV} \leq \varepsilon$ if

$$k > T\gamma^{-1} + 1 \quad \text{and} \quad \gamma \leq \sqrt{\frac{\kappa_1 \varepsilon^2}{4\kappa_2 \bar{D}(z) L^2(T + \bar{\gamma})}} \wedge \bar{\gamma},$$

where

$$T = (\log(\bar{C}(z)) - \log(\varepsilon/2)) / (-\log(\varrho)),$$

for some explicit constants $\bar{C}(z)$ and $\bar{D}(z)$.

Distance between μ_γ and μ

- **Assumption:** There exist $C < +\infty$, $\nu > 0$ such that, for any initial distribution ν_0, μ_0 ,

$$\|\nu_0 P_t - \mu_0 P_t\|_{TV} \leq C e^{-\nu t} \|\nu_0 - \mu_0\|_{TV} .$$

Theorem 7

For any $\gamma < \bar{\gamma}$,

$$\|\mu - \mu_\gamma\|_{TV} \leq \gamma E$$

for some constant E .

Thank you for your attention.