

Langevin Monte Carlo for Log-Concave Densities

Sampling Algorithms and Geometries
on Probability Distributions

Arnak Dalalyan

(joint with A. Karagulyan, L. Riou-Durand)

CREST, ENSAE Paris, IP Paris



The problem of sampling

- **The target density**

$$\pi(\boldsymbol{\theta}) \propto \exp(-f(\boldsymbol{\theta})), \quad f: \mathbb{R}^p \rightarrow \mathbb{R},$$

$$\mu_q(\pi) = \left(\int_{\mathbb{R}^p} \|\boldsymbol{\theta}\|_2^q \pi(\boldsymbol{\theta}) d\boldsymbol{\theta} \right)^{1/q} < +\infty, \quad q = 2.$$

- **Conditions on f :** gradient Lipschitz + convex

$$\|\nabla f(\mathbf{u}) - \nabla f(\mathbf{v})\|_2 \leq M \|\mathbf{u} - \mathbf{v}\|_2 \quad \& \quad \nabla^2 f(\mathbf{u}) \succeq 0.$$

- **Example:** posterior of the multivariate logistic regression.
- **Goal:** find a distr. ν easy to sample and s.t. $W_q(\nu, \pi)$ is small.
- **Constant sampling:** if $\nu = \delta_0$, then $W_q(\nu, \pi) \leq \mu_q(\pi)$.
- **Equiv. of moments:** there is A_q s.t. for all log-concave π ,

$$\mu_q(\pi) \leq A_q \mu_2(\pi)$$

(explicit expression leading to $A_3 \leq 3.5, A_4 \leq 4.6$).

First-order MCMC methods

Gradient oracle: assume that at any $\theta \in \mathbb{R}^p$, we can evaluate $\nabla f(\theta)$.

- Langevin Monte Carlo (LMC)

$$\vartheta_{k+1}^{\text{LMC}} = \vartheta_k^{\text{LMC}} - h\nabla f(\vartheta_k^{\text{LMC}}) + \sqrt{2h} \xi_{k+1}; \quad k = 0, 1, \dots \quad (\text{LMC})$$

where $\{\xi_k\}$ is iid $\mathcal{N}_p(0, I_p)$ indep. of ϑ_0 . Set $\nu_k^{\text{LMC}} = \mathcal{L}(\vartheta_k^{\text{LMC}})$.

- Langevin Monte Carlo with averaging (LMCa)

$$\vartheta_k^{\text{aLMC}} = \vartheta_\tau^{\text{LMC}}; \quad \tau \sim \text{Unif}(1, \dots, k), \quad k = 0, 1, \dots \quad (\text{LMCa})$$

and set $\nu_k^{\text{aLMC}} = \mathcal{L}(\vartheta_k^{\text{aLMC}})$.

- Kinetic Langevin Monte Carlo (KLMC, aka underdamped LMC)

$$\begin{bmatrix} \mathbf{v}_{k+1} \\ \vartheta_{k+1} \end{bmatrix} = \begin{bmatrix} \psi_0(h)\mathbf{v}_k - \psi_1(h)\nabla f(\vartheta_k) \\ \vartheta_k + \psi_1(h)\mathbf{v}_k - \psi_2(h)\nabla f(\vartheta_k) \end{bmatrix} + \sqrt{2\gamma} \begin{bmatrix} \xi_{k+1}^{(1)} \\ \xi_{k+1}^{(2)} \end{bmatrix} \quad (\text{KLMC})$$

- Metropolis adjusted Langevin Algorithm (MALA)

First-order MCMC methods

Gradient oracle: assume that at any $\theta \in \mathbb{R}^p$, we can evaluate $\nabla f(\theta)$.

- Langevin Monte Carlo (LMC)

$$\vartheta_{k+1}^{\text{LMC}} = \vartheta_k^{\text{LMC}} - h\nabla f(\vartheta_k^{\text{LMC}}) + \sqrt{2h} \xi_{k+1}; \quad k = 0, 1, \dots \quad (\text{LMC})$$

where $\{\xi_k\}$ is iid $\mathcal{N}_p(0, I_p)$ indep. of ϑ_0 . Set $\nu_k^{\text{LMC}} = \mathcal{L}(\vartheta_k^{\text{LMC}})$.

- Langevin Monte Carlo with averaging (LMCa)

$$\vartheta_k^{\text{aLMC}} = \vartheta_\tau^{\text{LMC}}; \quad \tau \sim \text{Unif}(1, \dots, k), \quad k = 0, 1, \dots \quad (\text{LMCa})$$

and set $\nu_k^{\text{aLMC}} = \mathcal{L}(\vartheta_k^{\text{aLMC}})$.

- Kinetic Langevin Monte Carlo (KLMC, aka underdamped LMC)

$$\begin{bmatrix} \mathbf{v}_{k+1} \\ \vartheta_{k+1} \end{bmatrix} = \begin{bmatrix} \psi_0(h)\mathbf{v}_k - \psi_1(h)\nabla f(\vartheta_k) \\ \vartheta_k + \psi_1(h)\mathbf{v}_k - \psi_2(h)\nabla f(\vartheta_k) \end{bmatrix} + \sqrt{2\gamma} \begin{bmatrix} \xi_{k+1}^{(1)} \\ \xi_{k+1}^{(2)} \end{bmatrix} \quad (\text{KLMC})$$

- Metropolis adjusted Langevin Algorithm (MALA)

An illustration and the objective

Sampling guarantees

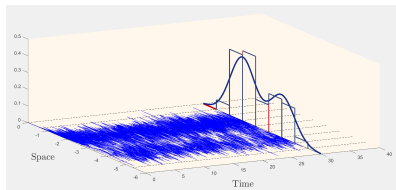


Figure: The blue lines represent different paths of a discretized Langevin process. We see that the histogram of the state at time $t = 30$ is close to the target density (the dark blue line).

Main goal: number of gradient evaluations that are sufficient to get ε -accuracy in W_q (especially for $q = 1$ and $q = 2$).

Mixing time of an approximate sampling algorithm Alg :

$$K_{\text{Alg}, W_q}(p, \varepsilon) = \min \{ k \in \mathbb{N} : W_q(\nu_k^{\text{Alg}}, \pi) \leq \varepsilon \mu_q(\pi), \forall \pi \in \mathcal{P} \}.$$

Quick overview

Order of magnitude of the mixing time of various first-order samplers.

$\kappa = M\mu_2^2/p$	LMCa	MALA	α -LMC	α -KLMC
W_2	—	—	$\kappa p^2/\varepsilon^6$	$\kappa^{1.5} p^2/\varepsilon^5$
W_1	—	—	$\kappa p^2/\varepsilon^4$	$\kappa^{1.5} p^2/\varepsilon^3$
d_{TV}	$\kappa p^2/\varepsilon^4 \triangle$	$p^3(\kappa/\varepsilon)^{3/2} \diamond$	$\kappa^2 p^3/\varepsilon^4 \square$	—

- \triangle behavior of the LMC with averaging [Durmus et al., 2019].
- \diamond derived from [Dwivedi et al., 2018]
- \square behavior of the LMC [Dalalyan, 2017].

First approach

Wasserstein from TV

- **Proposition** For any pair of probability measures (ν, ν') , and for any $q \geq 1$, we have:

$$W_q(\nu, \nu') \leq \inf_{r \geq q} \left\{ (\mu_r(\nu) + \mu_r(\nu')) d_{\text{TV}}(\nu, \nu')^{\frac{1}{q} - \frac{1}{r}} \right\}.$$

- **Proof** optimal coupling $X \sim \nu$ and $Y \sim \nu'$ for the TV-distance:
 $d_{\text{TV}}(\nu, \nu') = P(X \neq Y)$ and $W_q(\nu, \nu') \leq E[\|X - Y\|_2^q]^{1/q}$

- If π log-concave, $\mu_1(\nu) \lesssim \mu_2(\pi)$ and $d_{\text{TV}}(\nu, \pi) \leq (\varepsilon/A_r)^{1+1/(r-1)}$
then

$$W_1(\nu, \pi) \leq \varepsilon \mu_2(\pi).$$

- If π log-concave, $\mu_2(\nu) \lesssim \mu_2(\pi)$ and $d_{\text{TV}}(\nu, \pi) \leq (\varepsilon/A_r)^{2+4/(r-2)}$
then

$$W_2(\nu, \pi) \leq \varepsilon \mu_2(\pi).$$

First approach

Wasserstein from TV

- **Proposition** For any pair of probability measures (ν, ν') , and for any $q \geq 1$, we have:

$$W_q(\nu, \nu') \leq \inf_{r \geq q} \left\{ (\mu_r(\nu) + \mu_r(\nu')) d_{\text{TV}}(\nu, \nu')^{\frac{1}{q} - \frac{1}{r}} \right\}.$$

- **Proof** optimal coupling $X \sim \nu$ and $Y \sim \nu'$ for the TV-distance: $d_{\text{TV}}(\nu, \nu') = P(X \neq Y)$ and $W_q(\nu, \nu') \leq E[\|X - Y\|_2^q \mathbf{1}_{X \neq Y}]^{1/q}$. Use the Hölder inequality to conclude.

- If π log-concave, $\mu_1(\nu) \lesssim \mu_2(\pi)$ and $d_{\text{TV}}(\nu, \pi) \leq (\varepsilon/A_r)^{1+1/(r-1)}$ then

$$W_1(\nu, \pi) \leq \varepsilon \mu_2(\pi).$$

- If π log-concave, $\mu_2(\nu) \lesssim \mu_2(\pi)$ and $d_{\text{TV}}(\nu, \pi) \leq (\varepsilon/A_r)^{2+4/(r-2)}$ then

$$W_2(\nu, \pi) \leq \varepsilon \mu_2(\pi).$$

First approach

Wasserstein from TV

- **Proposition** For any pair of probability measures (ν, ν') , and for any $q \geq 1$, we have:

$$W_q(\nu, \nu') \leq \inf_{r \geq q} \left\{ (\mu_r(\nu) + \mu_r(\nu')) d_{\text{TV}}(\nu, \nu')^{\frac{1}{q} - \frac{1}{r}} \right\}.$$

- **Proof** optimal coupling $X \sim \nu$ and $Y \sim \nu'$ for the TV-distance: $d_{\text{TV}}(\nu, \nu') = P(X \neq Y)$ and
- If π log-concave, $\mu_1(\nu) \lesssim \mu_2(\pi)$ and $d_{\text{TV}}(\nu, \pi) \leq (\varepsilon/A_r)^{1+1/(r-1)}$ then

$$W_1(\nu, \pi) \leq \varepsilon \mu_2(\pi).$$

- If π log-concave, $\mu_2(\nu) \lesssim \mu_2(\pi)$ and $d_{\text{TV}}(\nu, \pi) \leq (\varepsilon/A_r)^{2+4/(r-2)}$ then

$$W_2(\nu, \pi) \leq \varepsilon \mu_2(\pi).$$

Quick overview

Order of magnitude of the mixing time of various first-order samplers.

$\kappa = M\mu_2^2/p$	LMCa	MALA	α -LMC	α -KLMC
W_2	$A_r^8 p^2 / \varepsilon^{8+8/r}$	$A_r^3 p^3 / \varepsilon^{3+3/r}$	p^2 / ε^6	p^2 / ε^5
W_1	$A_r^4 p^2 / \varepsilon^{4+4/r}$	$A_r^{3/2} p^3 / \varepsilon^{3/2+1/r}$	p^2 / ε^4	p^2 / ε^3
d_{TV}	$p^2 / \varepsilon^4 \triangle$	$p^3 / \varepsilon^{3/2} \diamond$	$p^3 / \varepsilon^4 \square$	—

- dependence on p is not better than for the penalized LMC.
- The results for MALA involve very large constants.
- Good dependence on ε requires large r , but then the constants A_r blow up

Second approach: Poincaré inequality

- π satisfies the Poincaré inequality if

$$\text{Var}_\pi[h] \leq C_P \int \|\nabla h(\boldsymbol{\theta})\|_2^2 \pi(\boldsymbol{\theta}) d\boldsymbol{\theta}.$$

We call C_P the Poincaré constant.

- For any log-concave π , $C_P < \infty$.
- The Langevin diffusion satisfies ([Chewi et al., 2020, Lehec, 2021, Vempala and Wibisono, 2019])

$$W_2(\nu_t, \pi)^2 \leq 2C_P \chi^2(\nu_0, \pi) e^{-t/C_P}, \quad t \geq 0.$$

Theorem [Lehec, 2021, Thm 3] If π is log-concave and f is L -Lipschitz, then LMC with step-size $h \leq p/L^2$ satisfies

$$W_2(\nu_k, \pi) \leq 4C_P \chi^2(\nu_0, \pi) e^{-kh/C_P} + 3Lkh^{3/2} \sqrt{p}.$$

If $\boldsymbol{\vartheta}_0 \sim \mathcal{N}(x_*, (p/L^2)I)$ then $\log \chi^2(\nu_0, \pi) \leq p + p \log(L^2 C_P / p)$.

Quick overview

Order of magnitude of the mixing time of various first-order samplers.

$\kappa = M\mu_2^2/p$	LMCa	MALA	α -LMC	α -KLMC
W_2	$A_r^8 p^2 / \varepsilon^{8+8/r}$	$A_r^3 p^3 / \varepsilon^{3+3/r}$	p^2 / ε^6	p^2 / ε^5
W_1	$A_r^4 p^2 / \varepsilon^{4+4/r}$	$A_r^{3/2} p^3 / \varepsilon^{3/2+1/r}$	p^2 / ε^4	p^2 / ε^3
d_{TV}	$p^2 / \varepsilon^4 \triangle$	$p^3 / \varepsilon^{3/2} \diamond$	$p^3 / \varepsilon^4 \square$	—

- Lehec's result leads to

$$K_{LMC, W_2} = \Theta\left(\frac{C_P^3 L^2 p^2}{\varepsilon^4}\right)$$

- Mathematically elegant result, but dependence on C_P is annoying.
- f global-Lipschitz assumption might be violated.

Third approach

Adding a quadratic penalty

- If f is α -strongly log-concave, then one has $\mathcal{C}_P \leq 1/\alpha$ and $W_2(\nu_t, \pi) \leq e^{-\alpha t} W_2(\nu_0, \pi)$.
- Define the str. convex surrogate $f_\alpha(\theta) := f(\theta) + \alpha \|\theta\|_2^2/2$ and

$$\pi_\alpha(\theta) := \frac{e^{-f_\alpha(\theta)}}{\int_{\mathbb{R}^p} e^{-f_\alpha(\mathbf{v})} d\mathbf{v}}.$$

- **Proposition** We have

$$d_{\text{TV}}(\pi, \pi_\alpha) \leq \alpha \mu_2^2(\pi) \quad W_q^q(\pi, \pi_\alpha) \leq C_q \alpha \mu_2(\pi)^{q+2}.$$

In particular, $C_1 \leq 22$ and $C_2 \leq 111$.

- Define α -LMC as LMC for f_α .
- Use the triangle inequality

$$\text{dist}(\nu_{k,\alpha}^{\text{Alg}}, \pi) \leq \text{dist}(\nu_{k,\alpha}^{\text{Alg}}, \pi_\alpha) + \text{dist}(\pi_\alpha, \pi). \quad (1)$$

Main result for α -LMC

Theorem

Suppose that the potential f is convex and M -Lipschitz. Let $q \in [1, 2]$. Then, for every $\alpha \leq M/20$ and $h \leq 1/(M + \alpha)$, we have

$$W_q(\nu_K^{\alpha\text{-LMC}}, \pi) \leq \underbrace{\sqrt{\mu_2}(1 - \alpha h)^{K/2}}_{\text{error due to the time finiteness}} + \underbrace{(2.1hMp/\alpha)^{1/2}}_{\text{discretization error}} + \underbrace{\left(C_q \alpha \mu_2^{q+2}\right)^{1/q}}_{\text{error due to the lack of strong-convexity}}.$$

- **FAQ:** why α is in the discretization error as well ?
- Optimizing wrt to α and h , we get

$$K_{\alpha\text{-LMC}, W_1} \leq 5 \times 10^4 M \frac{\mu_2^2 p}{\varepsilon^4} \log(100/\varepsilon)$$

$$K_{\alpha\text{-LMC}, W_2} \leq 4 \times 10^6 M \frac{\mu_2^2 p}{\varepsilon^6} \log(100/\varepsilon).$$

Main result for α -KLMC

Theorem

Suppose f is convex and M -Lipschitz. Let $q \in [1, 2]$. Then for every $\alpha \leq M/20$, $\gamma \geq \sqrt{M + 2\alpha}$ and $h \leq \alpha/(4\gamma(M + \alpha))$, we have

$$W_q(\nu_K^{\alpha\text{-KLMC}}, \pi) \leq \underbrace{\sqrt{2\mu_2} \left(1 - \frac{3\alpha h}{4\gamma}\right)^K}_{\text{error finite time}} + \underbrace{1.5Mp^{1/2}(h/\alpha)}_{\text{discretization error}} + \underbrace{\left(C_q \alpha \mu_2^{q+2}\right)^{1/q}}_{\text{error due to the lack of strong-convexity}}.$$

Optimizing wrt to α and h , we get

$$K_{\alpha\text{-KLMC}, w_1}(p, \varepsilon) \leq 9.2 \times 10^3 (M\mu_2^2)^{3/2} (p^{1/2}/\varepsilon^3) \log(150/\varepsilon)$$

$$K_{\alpha\text{-KLMC}, w_2}(p, \varepsilon) \leq 4.4 \times 10^5 (M\mu_2^2)^{3/2} (p^{1/2}/\varepsilon^5) \log(150/\varepsilon).$$

Conclusions and outlook

- Non-asymptotic sampling guarantees for convex (but not strongly convex) and gradient-Lipschitz potentials.
- The simple convexification trick is still “competitive”.
- Faster rates are obtained under additional smoothness (Hessian Lipschitz) assumptions.
- Current work: variable step-size h_k + variable penalty α_k + randomized mid-point discretization [Shen and Lee, 2019].
- Time-continuous bound in [Karagulyan and Dalalyan, 2020]: if $\alpha(t) = 1/(t + \mu_2^2(\pi))$ then

$$W_2(\nu_t, \pi) \leq \frac{10\mu_2^2(1 + \log(1 + t/\mu_2^2))}{\sqrt{t + \mu_2^2}}.$$

Conclusions and outlook

- Non-asymptotic sampling guarantees for convex (but not strongly convex) and gradient-Lipschitz potentials.
- The simple convexification trick is still “competitive”.
- Faster rates are obtained under additional smoothness (Hessian Lipschitz) assumptions.
- Current work: variable step-size h_k + variable penalty α_k + randomized mid-point discretization [Shen and Lee, 2019].
- Time-continuous bound in [Karagulyan and Dalalyan, 2020]: if $\alpha(t) = 1/(t + \mu_2^2(\pi))$ then

$$W_2(\nu_t, \pi) \leq \frac{10\mu_2^2(1 + \log(1 + t/\mu_2^2))}{\sqrt{t + \mu_2^2}}.$$

Interested in a postdoc in Paris: send me an email.

References I

- S. Chewi, T. L. Gouic, C. Lu, T. Maunu, P. Rigollet, and A. Stromme. Exponential ergodicity of mirror-Langevin diffusions, 2020.
- A. S. Dalalyan. Theoretical guarantees for approximate sampling from a smooth and log-concave density. *J. R. Stat. Soc. B*, 79:651 – 676, 2017.
- A. Durmus, S. Majewski, and B. Miasojedow. Analysis of Langevin Monte-Carlo via convex optimization. *J. Mach. Learn. Res.*, 20:73–1, 2019.
- R. Dwivedi, Y. Chen, M. J. Wainwright, and B. Yu. Log-concave sampling: Metropolis-hastings algorithms are fast! In *Conference on Learning Theory*, pages 793–797. PMLR, 2018.

References II

- A. Karagulyan and A. Dalalyan. Penalized Langevin dynamics with vanishing penalty for smooth and log-concave targets. *Advances in Neural Information Processing Systems*, 33, 2020.
- J. Lehec. The langevin monte carlo algorithm in the non-smooth log-concave case. *arXiv preprint arXiv:2101.10695*, 2021.
- R. Shen and Y. T. Lee. The randomized midpoint method for log-concave sampling. In *Advances in Neural Information Processing Systems*, pages 2098–2109, 2019.
- S. Vempala and A. Wibisono. Rapid convergence of the unadjusted Langevin algorithm: Isoperimetry suffices. In *Advances in Neural Information Processing Systems 32*, pages 8094–8106. 2019.