

# Penalized Langevin dynamics with vanishing penalty for smooth and log-concave targets

Avetik Karagulyan

MMS 2020, “Luminy”, France



## General setting of the problem

- **Problem:** sample from a given target distribution  $\pi$  defined in  $\mathbb{R}^p$  with a large value of  $p$ .
- More precisely, for a given precision level  $\epsilon$ , generate a random vector  $\theta$  with values in  $\mathbb{R}^p$  such that its distribution  $\mu$  satisfies

$$\text{distance}(\mu, \pi) \leq \epsilon. \quad (1)$$

- Important particular case:  $\pi$  has a density (w.r.t. the Lebesgue measure) given by

$$\pi(\theta) \propto \exp(-f(\theta)), \quad (2)$$

with a “potential”  $f : \mathbb{R}^p \rightarrow \mathbb{R}$ .

## General setting of the problem

### Main assumption $\mathbf{A}(m, M)$ :

- $f \in C^2(\mathbb{R}^p)$
- $m$ -strong convexity:  $\nabla^2 f(\boldsymbol{\theta}) \succeq mI_p$ , with  $m \geq 0$ ;
- $M$ -Lipschitz gradients:  $\nabla^2 f(\boldsymbol{\theta}) \preceq MI_p$ , with  $M > 0$ .

### Wasserstein distance:

$$W_2(\nu, \nu') = \inf \left\{ \mathbf{E}[\|\boldsymbol{\vartheta} - \boldsymbol{\vartheta}'\|_2^2]^{1/2} : \boldsymbol{\vartheta} \sim \nu \text{ and } \boldsymbol{\vartheta}' \sim \nu' \right\}, \quad (3)$$

where the infimum is over all joint distributions having  $\nu$  and  $\nu'$  as the first and the second marginal distributions.

# Langevin Diffusion

- **Vanilla Langevin diffusion:**

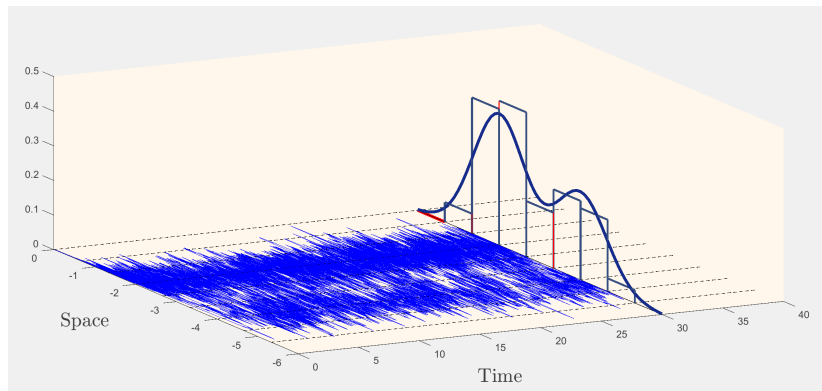
$$d\mathbf{L}_t^{\text{LD}} = -\nabla f(\mathbf{L}_t^{\text{LD}})dt + \sqrt{2}d\mathbf{W}_t. \quad (\text{LD})$$

The solution of this equation is a Markov process having  $\pi$  as an invariant distribution.

- When the potential function  $f$  is  $m$ -strongly convex, the Markov process is ergodic and it converges to  $\pi$  exponentially (Villani 2008):

$$W_2(\nu_t^{\text{LD}}, \pi) \leq e^{-mt} W_2(\nu_0^{\text{LD}}, \pi). \quad (4)$$

## Illustration of LD



The blue lines represent different paths of the Langevin process. We see that the histogram of the state at time  $t = 30$  is close to the target density.

## Langevin Monte-Carlo

Langevin Monte-Carlo (LMC) is defined as:

$$\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - h\nabla f(\boldsymbol{\theta}_k) + \sqrt{2h} \boldsymbol{\xi}_{k+1}; \quad k = 0, 1, 2, \dots \quad (5)$$

$(\boldsymbol{\xi}_k)_{k \in \mathbb{N}}$  is a sequence of i.i.d standard Gaussians ind. from  $\boldsymbol{\theta}_k$ .

### Theorem (Durmus et al, 2018)

*Suppose  $f$  satisfies the assumption  $A(m, M)$  with  $m > 0$ . If  $h < 1/M$ , then the following upper bound is satisfied:*

$$W_2(\mu_K, \pi) \leq (1 - mh)^{K/2} W_2(\mu_0, \pi) + \{2hp\kappa\}^{1/2} \quad (6)$$

*where  $\mu_K$  is the law of  $\boldsymbol{\theta}_K$  and  $\kappa = M/m$  is the condition number.*

## Remarks on the theorem:

- The term  $(1 - mh)^{K/2} W_2(\mu_0, \pi)$  can be viewed as an approximation error.
- The term  $\{2hp\kappa\}^{1/2}$  is the error due to the discretization.
- We have at our disposal the step-size  $h$  and the time horizon  $T = Kh$ . Therefore choosing

$$h \leq \frac{\varepsilon^2}{p\kappa} \quad \text{and} \quad K = \frac{p\kappa}{m\varepsilon^2} \log(2/\varepsilon) \quad (7)$$

we have  $W_2(\mu_K, \pi) < \varepsilon$ .

# Kinetic Langevin Dynamics

## Kinetic (underdamped) Langevin diffusion:

$$\begin{aligned}d\mathbf{L}_t^{\text{KLD}} &= \mathbf{V}_t^{\text{KLD}} dt; \\d\mathbf{V}_t^{\text{KLD}} &= -(\eta\mathbf{V}_t^{\text{KLD}} + \nabla f(\mathbf{L}_t^{\text{KLD}}))dt + \sqrt{2\eta}\mathbf{W}_t.\end{aligned}$$

- Here  $W_t$  is a Brownian motion and  $\eta$  is the friction parameter. (LD) is the limit of the rescaled kinetic diffusion  $\bar{\mathbf{L}}_t = \mathbf{L}_{\eta t}^{\text{KLD}}$  when the friction coefficient  $\eta$  tends to infinity (Nelson et al, '65).
- The Markov process  $(\mathbf{L}_t^{\text{KLD}}; \mathbf{V}_t^{\text{KLD}})$  is positive recurrent and its invariant distribution is absolutely continuous wrt the Lebesgue measure on  $\mathbb{R}^{2p}$ . The corresponding invariant density is

$$p_*(\boldsymbol{\theta}, \mathbf{v}) \propto \exp(-f(\boldsymbol{\theta}) - \|\mathbf{v}\|_2^2/2). \quad (8)$$

- Mixing-time in Wasserstein distance is of order  $\exp(-t/\kappa)$ . (Eberle et al., 2017)



# Kinetic Langevin Monte-Carlo

**KLMC** is a discretization of Kinetic LD:

$$\mathbf{v}_{k+1} = \psi_0(h)\mathbf{v}_k - \psi_1(h)\nabla f(\boldsymbol{\theta}_k) + \sqrt{2\eta}\boldsymbol{\xi}_{k+1};$$

$$\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k + \psi_1(h)\mathbf{v}_k - \psi_2(h)\nabla f(\boldsymbol{\theta}_k) + \sqrt{2\eta}\boldsymbol{\xi}'_{k+1}.$$

- where  $\psi_0(x) = \exp(-\eta x)$  and  $\psi_{k+1}(x) = \int_0^x \psi_k(x) dx$ .
- $(\boldsymbol{\xi}_k, \boldsymbol{\xi}'_k)_{k \in \mathbb{N}}$  is a sequence of standard Gaussians independent from the initial conditions.

## Theorem (Dalalyan and Riou-Durand, '18)

Assume that  $f$  satisfies  $A(m, M)$ , with  $m > 0$ . Then for  $\eta > \sqrt{m + M}$  and  $h < m/4\eta M$ , the following bound is true:

$$W_2(\nu_k, \pi) \leq \sqrt{2} \left(1 - \frac{0.75mh}{\eta}\right)^k W_2(\nu_0, \pi) + \frac{Mh\sqrt{2p}}{m}.$$

## Remarks on KLMC

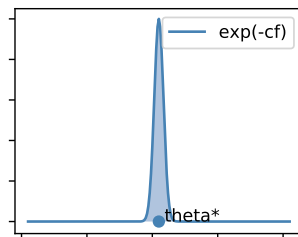
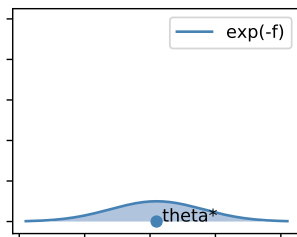
- From the theorem we deduce that in order to have  $\varepsilon$  Wasserstein sampling error it suffices to do

$$K_{\text{KLMC}} = \tilde{O} \left( \kappa^{3/2} \left\{ \kappa \sqrt{\frac{p}{m\varepsilon^2}} \right\}^{1/2} \right) \quad (9)$$

iterations of KLMC.

- Compared to LMC, KLMC has better convergence rate in terms of  $p$  and  $\varepsilon$ . This is due to smoother trajectories of KLD. (Cheng et al, '17).
- Nevertheless, LMC and KLMC are not comparable. In the case when  $\nabla f$  increases rapidly, LMC performs better.

## Relation with optimization



- Suppose we can have exact samples from  $\pi \propto \exp(-f)$ . Then we can also sample from  $\pi^c \propto \exp(-cf)$ . When  $c$  is large, the samples are concentrated around the minimum point  $\theta^*$ .

## Relation with optimization

- In the case of non-convex  $f$ , one can use LMC for the minimization problem. The Gaussian noise term allows to escape local minima.
- The gradient flow of the function  $f$  is defined as the following differential equation:

$$d\mathbf{L}_t^{\text{GF}} = -\nabla f(\mathbf{L}_t^{\text{GF}})dt. \quad (10)$$

GF can be seen as a limit version of LD scaled over time axis, when the scaling parameter tends to infinity.  
(discussed later)

## Relation with optimization

We define the following functionals for every  $\mu \in \mathcal{P}_2(\mathbb{R}^p)$ :

- **Energy functional:**  $\mathcal{E}(\mu) := \mathbb{E}_\mu[f]$ .
- **Gibbs-Boltzmann functional:**  $\mathcal{H}(\mu) := \mathbb{E}_\mu[\log(d\mu/d\lambda)]$ , where  $\lambda$  is the Lebesgue measure.

It is known that  $\pi$  minimizes  $\mathcal{F} := \mathcal{H} + \mathcal{E}$ , and that the Langevin diffusion is its gradient flux.

**What happens when  $f$  is convex but not strongly?**

## Discrete case

- (Durmus et al, 2018) have shown that LMC with averaging converges:

$$\text{KL}(\nu'_K|\pi) \leq \varepsilon, \quad \text{when} \quad K = O\left(\frac{\kappa p}{\varepsilon^2}\right), \quad (11)$$

where  $\nu'_K$  is the mixture of the first  $K$  iterates of the LMC.

**The idea:** leverage the connection between Langevin dynamics and optimization of KL divergence (Jordan et al. '98)(not Michael Jordan).

- (Dalalyan, K. Riou-Durand, 2019) proposed to use constant linear penalty in LMC and KLMC. **The idea:** sample from the adjusted density  $\exp(-f(\theta) - \alpha\|\theta\|_2^2/2)$  and optimize over  $\alpha$ . Convergence in Wasserstein-1 and Wasserstein-2 distance are obtained.

## LD + Poincaré inequality

**Poincaré inequality:** We say  $\pi$  satisfies the Poincaré inequality, if for  $\forall g \in L^2(\pi)$  locally-Lipschitz, we have

$$\text{var}_\pi[g] \leq C_P \mathbb{E}_\pi \left[ \|\nabla g\|^2 \right]. \quad (\text{P})$$

Chewi et al (2020) have shown that if  $\pi$  satisfies (P), then

$$W_2^2(\mu_t, \pi) \leq 2C_P e^{-\frac{2t}{C_P}} \chi^2(\mu_0 \| \pi). \quad (12)$$

The convergence is exponential but the constants can be large.



## KLS conjecture

$C_p$  of any distribution  $\pi$  is bounded by some universal constant  $C_{\text{KLS}}$  times the largest eigenvalue of the covariance matrix of  $\pi$ :

$$C_p \leq C_{\text{KLS}} \times \mathbb{E}_\pi [\|\mathbf{X}^T \mathbf{X}\|_{\text{op}}]. \quad (13)$$

## Penalized Langevin Dynamics

We propose to modify the Langevin equation by adding a vanishing linear penalty:

$$d\mathbf{L}_t^{\text{PLD}} = -(\nabla f(\mathbf{L}_t^{\text{PLD}}) + \alpha(t)\mathbf{L}_t^{\text{PLD}}) dt + \sqrt{2} d\mathbf{W}_t, \quad (14)$$

where  $\alpha : [0, \infty) \rightarrow [0, \infty)$  is a positive time-dependent penalty factor converging to zero as  $t \rightarrow \infty$ .

- For every fixed time instant  $t$  the potential is equal to  $f(\cdot) + \alpha(t)\|\cdot\|^2/2$ , which is  $\alpha(t)$ -strongly convex.

## Convergence of PLD

Define  $\mu_2 := \mathbb{E}_\pi[\|\boldsymbol{\theta}\|^2]$ .

### Proposition

Suppose  $f$  satisfies A(0, M). If  $\alpha(t) = 1/(2\mu_2 + 2t)$ , then

$$W_2(\nu_t^{\text{PLD}}, \pi) \leq \frac{10\mu_2[1 + \log(1 + (t/\mu_2))]}{\sqrt{t + \mu_2}}, \quad (15)$$

### Remark:

- We obtain polynomial convergence.
- The result only depends on the second order moment.
- In the  $m$ -strongly convex case  $\mu_2 \leq p/m$ .

## Main theorem

### Theorem

Suppose  $f$  satisfies  $A(0, M)$ . Then, for every positive number  $t$  and for  $\beta(t) = \int_0^t \alpha(u) du$ , we have

$$W_2(\nu_t^{PLD}, \pi) \leq \sqrt{\mu_2} e^{-\beta(t)} + 11\mu_2 e^{-\beta(t)} \int_0^t \frac{|\alpha'(s)| e^{\beta(s)}}{\sqrt{\alpha(s)}} ds + \sqrt{\alpha(t)} \mu_2.$$

- The error bound stated in the previous proposition is obtained by “optimizing” this upper bound w.r.t.  $\alpha$ .

## Sketch of the proof

We apply the triangle inequality to  $W_2(\nu_t^{\text{PLD}}, \pi)$ :

$$W_2\left(\nu_t^{\text{PLD}}, \pi\right) \leq W_2\left(\nu_t^{\text{PLD}}, \pi_{\alpha(t)}\right) + W_2\left(\pi_{\alpha(t)}, \pi\right), \quad (16)$$

where  $\pi_\gamma(\boldsymbol{\theta}) \propto \exp(-f(\boldsymbol{\theta}) - \gamma\|\boldsymbol{\theta}\|^2/2)$ , for  $\forall \gamma > 0$ . We then define  $\phi(t) := W_2\left(\nu_t^{\text{PLD}}, \pi_{\alpha(t)}\right)$  and prove the following bound:

$$\phi'(t) \leq -\alpha(t)\phi(t) - \frac{11\alpha'(t)}{\sqrt{\alpha(t)}} \cdot \mu_2\left(\pi_{\alpha(t)}\right). \quad (17)$$

Applying Gronwall inequality to  $\phi$  and transportation inequality to the term  $W_2\left(\pi_{\alpha(t)}, \pi\right)$  we obtain the proof.

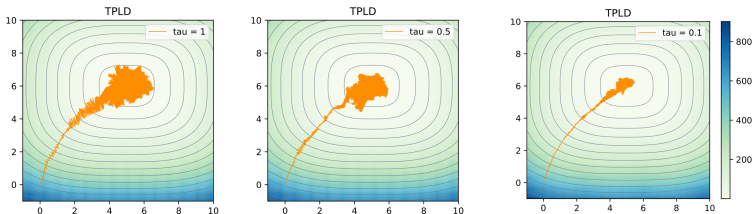
## Tempered PLD

- Define  $f_\tau(\cdot) := f(\cdot)/\tau$ ,  $\pi^\tau \propto \exp(-f_\tau)$  and

$$\mathbf{L}_t^\tau = \mathbf{L}_0^\tau - \frac{1}{\tau} \int_0^t (\nabla f(\mathbf{L}_s^\tau) + \alpha(s/\tau)\mathbf{L}_s^\tau) ds + \sqrt{2}\mathbf{W}_t.$$

- We know that  $\mathbf{L}_t^\tau \rightarrow \pi^\tau$ . Tempered PLD is defined as  $\mathbf{X}^{\text{TPLD}} := \mathbf{L}_{\tau t}^\tau$  and it satisfies

$$d\mathbf{X}_t^{\text{TPLD}} = - \left( \nabla f(\mathbf{X}_t^{\text{TPLD}}) + \alpha(t)\mathbf{X}_t^{\text{TPLD}} \right) dt + \sqrt{2\tau}d\mathbf{W}_t. \quad (\text{TPLD})$$



**Figure:** Plots of TPLD for different temperature parameter.

- Sampling  $\xrightarrow{\tau \rightarrow 0}$  optimization and (TPLD)  $\xrightarrow{\tau \rightarrow 0}$  (PGF):

$$d\mathbf{X}_u^{\text{PGF}} = - \left( \nabla f \left( \mathbf{X}_u^{\text{PGF}} \right) + \alpha(u) \mathbf{X}_u^{\text{PGF}} \right) du. \quad (\text{PGF})$$

## The convergence of PGF

**Technical assumption B(D,q):**

$$\|\mathbf{x}_\gamma - \mathbf{x}_{\tilde{\gamma}}\|_2 \leq \frac{D}{\tilde{\gamma}^q} (\tilde{\gamma} - \gamma) \|\mathbf{x}_*\|_2, \quad \forall \gamma < \tilde{\gamma} \quad (18)$$

where  $\mathbf{x}_\gamma := \arg \min \{f(\mathbf{x}) + \gamma \|\mathbf{x}\|^2/2\}$ .

### Theorem

*If  $f$  satisfies A(0, M) and B(D, q), then taking  $\alpha(t) = \frac{(1-q)}{(t+A)}$  we have*

$$\left\| \mathbf{X}_t^{\text{PGF}} - \boldsymbol{\theta}_* \right\|_2 \leq \frac{A^{1-q} + D + D \log(1 + (t/A))}{(t + A)^{1-q}} \|\boldsymbol{\theta}_*\|_2. \quad (19)$$



# Conclusion and outlook

## Take home message

- Penalized Langevin Dynamics gives explicit non-asymptotic bounds on the Wasserstein error of sampling.

## Future work

- Assessing the error of the penalized kinetic Langevin diffusion in continuous-time.
- Discretization of the time-continuous dynamics.
- Gain understanding on the condition required in the theorem for the PGF.

**This is the last slide.  
Thank you!**