

# **CPSC 661: Sampling Algorithms in ML**

---

Andre Wibisono

April 26, 2021

Yale University

## Last time

- Langevin dynamics
- Unadjusted Langevin Algorithm
- Variants of Langevin dynamics (first-order):  
  {Weighted, Mirror, Newton, Interacting} Langevin Dynamics

**Today:** Underdamped Langevin Dynamics

# First vs. Second-Order Dynamics

$$\dot{x}_t = -\nabla f(x_t)$$

- First order (in time) dynamics

- Gradient flow

- Dissipative: Minimizes  $f$

$$\frac{d}{dt} f(x_t) = -\|\nabla f(x_t)\|^2 \leq 0$$

- Greedy descent flow



$$\frac{d^2}{dt^2} x_t = \ddot{x}_t = -\nabla f(x_t)$$

- Second-order (in time) dynamics

- Newton's Law:

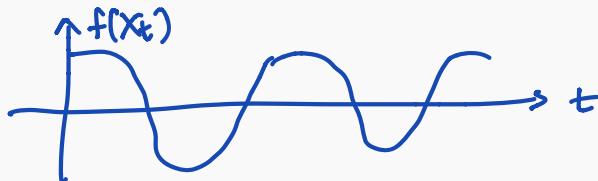
$$\underbrace{\text{mass}}_{m=1} \times \underbrace{\text{Acceleration}}_{\ddot{x}_t} = \underbrace{\text{Force}}_{-\nabla f(x_t)}$$

- Conservative: Preserves Hamiltonian

$$H(t) = \frac{1}{2}\|\dot{x}_t\|^2 + f(x_t)$$

$$\frac{d}{dt} H(t) = 0$$

- Function value  $f(x_t)$  is oscillatory



$$\text{Eg. } f(x) = \frac{1}{2} x^T A x$$

$$\nabla f(x) = Ax, \quad A = A^T > 0$$

Gradient flow:

$$\dot{x}_t = -\nabla f(x_t) = -Ax_t$$

$$\Rightarrow x_t = e^{-At} x_0 \rightarrow 0$$

Hamiltonian flow:

$$\ddot{x}_t = -\nabla f(x_t) = -Ax_t$$

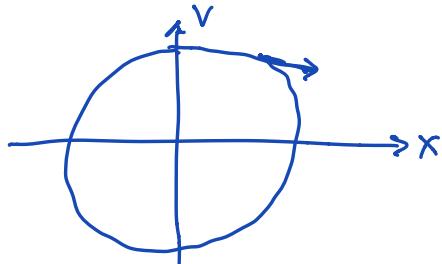
$$x_t = (\cos At) \beta_1 + (\sin At) \beta_2$$

where  $\beta_1, \beta_2 \in \mathbb{R}^n$  determined by

$$x_0, \dot{x}_0 \in \mathbb{R}^n$$

in phase space:

$$\begin{aligned} \dot{x}_t &= v_t \\ \dot{v}_t &= -Ax_t \end{aligned} \quad \left. \right\} \Leftrightarrow \ddot{x}_t = -Ax_t$$



# Optimization: Acceleration

Newton's Law with friction / damping

(HB)  $\ddot{\tilde{x}}_t + \gamma \dot{\tilde{x}}_t + \nabla f(\tilde{x}_t) = 0 , \quad \gamma > 0 \text{ damping}$

(also called \* Heavy ball [Polyak, Nesterov]

\* accelerated gradient flow )

\* no damping :  $\gamma = 0 \Rightarrow$  Newton's Law  $\tilde{\tilde{x}}_t + \nabla F(\tilde{x}_t) = 0$

\* Overdamped :  $\gamma \rightarrow \infty$  :  $\tilde{\tilde{x}}_t = \tilde{x}_{gt}$  regular gradient flow

$$\dot{\tilde{x}}_t = -\nabla f(\tilde{x}_t)$$

because  $\tilde{\tilde{x}}_t$  satisfies

$$\frac{1}{\gamma} \ddot{\tilde{x}}_t + \dot{\tilde{x}}_t + \nabla f(\tilde{x}_t) = 0$$

Eg.  $f(x) = \frac{\alpha}{2}x^2$  on  $x \in \mathbb{R} \rightarrow$  strongly convex w. param  $\alpha$

$$\text{then (HB): } \ddot{x}_t + \gamma \dot{x}_t + \alpha x_t = 0$$

$$\text{Characteristic eq: } r^2 + \gamma r + \alpha = 0$$

$$\Leftrightarrow r_{12} = \frac{-\gamma \pm \sqrt{\gamma^2 - 4\alpha}}{2}$$

$$\text{Solution to (HB) is } x_t = c_1 \cdot e^{r_1 t} + c_2 \cdot e^{r_2 t}, \quad c_1, c_2 \in \mathbb{R}$$

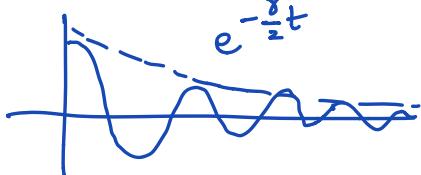
Case 1: Underdamped case:

$$\text{if } \gamma < 2\sqrt{\alpha}, \text{ then } r_{12} = \frac{-\gamma \pm i\omega}{2} \quad \omega = \sqrt{4\alpha - \gamma^2} \in \mathbb{R}$$

$$e^{r_{12} t} = e^{-\frac{\gamma \pm i\omega}{2} t}$$

$$(e^{i\theta} = \cos \theta + i \sin \theta) \quad = e^{-\frac{\gamma}{2} t} \left( \cos \left( \frac{\omega t}{2} \right) \pm i \sin \left( \frac{\omega t}{2} \right) \right)$$

$$\text{solution is } x_t = e^{-\frac{\gamma}{2} t} \left( c_1 \cdot \cos \left( \frac{\omega t}{2} \right) + c_2 \cdot \sin \left( \frac{\omega t}{2} \right) \right)$$



\* best rate is  $\gamma = 2\sqrt{\alpha}$  (critical damping)

$$\Rightarrow O(e^{-\sqrt{\alpha} t}) \text{ convergence}$$

$$(\text{vs. Gradient flow } \dot{x}_t = -\alpha x_t)$$

$$\Rightarrow x_t = e^{-\alpha t} x_0 \Rightarrow O(e^{-\alpha t}) \text{ convergence}$$

when  $0 < \alpha \ll 1, \quad \sqrt{\alpha} \gg \alpha$

Heavy ball:  $\ddot{x}_t + \gamma \dot{x}_t + \nabla f(x_t) = 0$

\* Hamiltonian is decreasing:

$$\frac{d}{dt} \left( \frac{1}{2} \|\dot{x}_t\|^2 + f(x_t) \right) = -\gamma \|\dot{x}_t\|^2 \leq 0$$

\* Theorem:

Assume  $f$  is  $\alpha$ -strongly convex ( $\nabla^2 f(x) \succeq \alpha I$ )

$$\text{and } 0 < \gamma \leq 2\sqrt{\alpha}$$

then along heavy ball,

$$\mathcal{E}_t = f(x_t) - f(x^*) + \frac{\alpha}{2} \|x_t - x^* + \frac{2}{\gamma} \ddot{x}_t\|^2$$

is a Lyapunov function, and it converges to 0 exp. fast:

$$\begin{aligned} \frac{d}{dt} \mathcal{E}_t &\leq -\frac{\gamma}{2} \mathcal{E}_t \\ \Rightarrow \mathcal{E}_t &\leq e^{-\frac{\gamma}{2} t} \mathcal{E}_0 \\ \Rightarrow f(x_t) - f(x^*) &\leq \mathcal{E}_t \leq e^{-\frac{\gamma}{2} t} \mathcal{E}_0 \end{aligned}$$

Then best rate is  $r = 2\sqrt{\alpha}$ :

$$\Rightarrow f(x_t) - f(x^*) \leq O(e^{-\sqrt{\alpha} t})$$

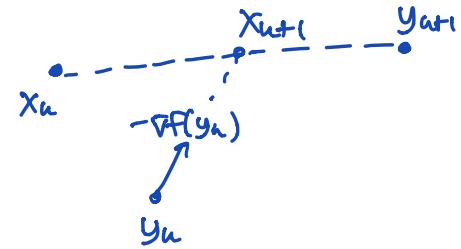
(c.f. Gradient flow  
 $f(x_t) - f(x^*) \leq e^{-\alpha t}$ )

↓  
discretize:

Accelerated gradient descent (momentum method)

$$(AGD) \quad x_{n+1} = y_n - \eta \nabla f(y_n)$$

$$y_{n+1} = x_n + \beta \cdot (x_{n+1} - x_n)$$



Thm: If  $f$  is  $\alpha$ -strongly convex,  $L$ -smooth

$$\alpha I \leq \nabla^2 f(x) \leq L I, \quad \kappa = \frac{L}{\alpha}$$

then with  $\eta = \frac{1}{L}$ ,  $\beta = \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1}$ :

AGD has  $f(x_n) - f(x^*) \leq \left(1 - \frac{1}{\sqrt{\kappa}}\right)^n \cdot C_0$

⇒ complexity is  $O(\sqrt{\kappa})$

(c.f. gradient descent: complexity  $O(\kappa)$ )

& AGD  $O(\sqrt{\kappa})$  is optimal,  $\exists$  lower bound  $\Omega(\sqrt{\kappa})$

Note: • HB:  $\ddot{x}_t + \gamma \dot{x}_t + \nabla^2 f(x_t) = 0$  for  $f$  strongly convex

• can use  $\ddot{x}_t + \frac{3}{t} \dot{x}_t + \nabla^2 f(x_t) = 0$  for  $f$  convex

(gets  $O(\frac{1}{t^2})$  convergence rate vs.  $O(\frac{1}{t})$  for GF)

# Underdamped Langevin Dynamics

$$(ULD) \quad \ddot{x}_t + \gamma \dot{x}_t + \nabla f(x_t) + \sqrt{2\gamma} \xi_t = 0$$

where  $\gamma > 0$  is damping

$$\xi_t = \frac{dw_t}{dt} \text{ is "white noise"}$$

(ULD)  $\Leftrightarrow$

$$\begin{aligned} dx_t &= v_t dt \\ dv_t &= (-\gamma v_t - \nabla f(x_t)) dt + \sqrt{2\gamma} dw_t \end{aligned}$$

Facts: \* In  $(x, v)$  space, stationary distribution is

$$\rho_\infty(x, v) = \nu(x) \cdot \mathcal{N}(v)$$

where  $\nu(x) \propto e^{-f(x)}$  for  $x \in \mathbb{R}^n$

$$\mathcal{N}(v) = \mathcal{N}(v; 0, I) \propto e^{-\frac{1}{2}\|v\|^2} \text{ for } v \in \mathbb{R}^n$$

$$\text{Let } H(x, v) = f(x) + \frac{1}{2} \|v\|^2$$

$$\text{then } p_\infty(x, v) \propto e^{-H(x, v)}$$

\* As  $\gamma \rightarrow \infty$ ,  $\tilde{x}_t = X_{\gamma t}$  recovers (overdamped) Langevin dynamics:

$$d\tilde{x}_t = -\nabla f(\tilde{x}_t) dt + \sqrt{2} dW_t$$

\* Hypoelliptic diffusion, not uniformly elliptic  
(diffusion only on  $v$  variable)

\* Lemma: If  $(x_t, v_t) \sim p_t(x, v)$  follows ULD

then  $p_t$  follows the "kinetic Fokker-Planck equation":

$$\frac{\partial p_t}{\partial t} = -\langle v, \nabla_x p_t(x, v) \rangle + \langle \nabla f(x), \nabla_v p_t(x, v) \rangle \quad \left. \begin{array}{l} \text{comes from} \\ \text{Hamiltonian} \\ \text{flow} \end{array} \right\}$$

$$+ \gamma \left( \nabla_v \circ (v p_t(x, v)) + \Delta_v p_t(x, v) \right) \quad \left. \begin{array}{l} \text{comes from} \\ \text{Langevin} \\ \text{dynamics} \end{array} \right\}$$

because ULD is a mixture:

$$\begin{pmatrix} dx_t \\ dv_t \end{pmatrix} = \underbrace{\begin{pmatrix} v & dt \\ -\nabla f(x) & dt \end{pmatrix}}_{\substack{\text{Hamiltonian} \\ \text{flow}}} + \underbrace{\begin{pmatrix} 0 \\ -\gamma v dt + \sqrt{2\gamma} dW_t \end{pmatrix}}_{\substack{\text{Langevin dynamics} \\ \text{for } v \text{ for target } \mathcal{N}(0, I)}}$$

$$\begin{cases} \dot{x} = v \\ \dot{v} = -\nabla f(x) \end{cases} \Rightarrow \ddot{x}_t = -\nabla f(x_t)$$

\* So for example,

$$\frac{d}{dt} \chi_{S_\infty}^2(S_t) = -2\gamma \mathbb{E}_{S_\infty} \left[ \left\| \nabla_v \frac{s_t(x, v)}{s_\infty(x, v)} \right\|^2 \right] \leq 0$$

How to get convergence rate?

\* Villani, Hypocoercivity, 2009

\* Eberle, Guillin, Zimmer, Synchronous + Reflection coupling, 2019

Will see:

\* Cao, Lu, Wang, Poincaré inequality

\* Dalalyan & Rio-Durand, Synchronous coupling

→ discrete-time algorithm

## ① Convergence Rate via Poincaré:

Theorem: [Cao, Lu, Wang, 2020]

Assume: •  $v(x) \leq e^{-f(x)}$  satisfies  $\alpha$ -Poincaré inequality

- $\|\nabla^2 f(x)\| \leq M(1 + \|\nabla f(x)\|)$  for all  $x \in \mathbb{R}^n$   
some  $M \geq 1$

- $f$  is super linear:

$$\exists \beta > 1 : \lim_{x \rightarrow \infty} \frac{f(x)}{\|x\|^\beta} = \infty$$

( $\Rightarrow$  embedding  $H^1(\Omega) \hookrightarrow L^2(\Omega)$  is compact)

Then ULD:

$$dx = V dt$$

$$dV = (-\gamma V - \nabla f(x)) dt + \sqrt{2\gamma} dW$$

has:  $\chi_{S_0}^2(s_t) \leq c_0 \cdot e^{-\lambda t} \chi_{S_0}^2(s_0)$

where  $\lambda = \sqrt{\alpha} \cdot \log \left( 1 + \frac{\gamma \sqrt{\alpha}}{c_1 (\sqrt{\alpha} + R + \gamma)^2} \right)$

where: 1) if  $f$  convex, then  $R=0$

$$\text{then with } \gamma = \sqrt{\alpha} : \lambda = \sqrt{\alpha} \cdot \log \left( 1 + \frac{1}{c_1} \right) = \Theta(\sqrt{\alpha})$$

2) if  $\nabla^2 f(x) \succeq -K \cdot I$ ,  $K > 0$ , then  $R = \sqrt{K}$  for  $\alpha \ll 1$ :

$$\text{then with } \gamma = \sqrt{\alpha} + K : \lambda = \sqrt{\alpha} \cdot \log \left( 1 + \frac{\sqrt{\alpha}}{\sqrt{\alpha} + K} \right) \approx \frac{\alpha}{\sqrt{\alpha} + K}$$

## ② Convergence rate by Synchronous Coupling

[Cheng et al. 2018, Dalalyan & Riau-Djoudi 2020]

Theorem: Assume  $f$  is  $\alpha$ -strongly convex,  $L$ -smooth

Then ULD:  $dx_t = v_t dt$

$$dv_t = (-\gamma v_t - \underline{\nabla f(x_t)}) dt + \sqrt{2\gamma} dW_t$$

$$\text{with } \gamma \geq \sqrt{L+\alpha},$$

ULD has exponential contraction, in particular to  $s_\infty$ :

$$W_2(s_t, s_\infty) \leq \sqrt{2} \cdot e^{-\frac{\alpha}{\gamma} t} W_2(s_0, s_\infty)$$

in discrete time: approximate ULD by

$$dx_t = v_t dt$$

$$dv_t = (-\gamma v_t - \underline{\nabla f(x_0)}) dt + \sqrt{2\gamma} dW_t$$

Can integrate, get algorithm:

Kinetic Langevin Monte Carlo (KMLC):

$$x_{k+1} = x_k + \Psi_1(\eta) v_k - \Psi_2(\eta) \nabla f(x_k) + \sqrt{2\gamma} \xi_k$$

$$v_{k+1} = \Psi_0(\eta) v_k - \Psi_1(\eta) \nabla f(x_k) + \sqrt{2\gamma} \xi'_k$$

$$\text{where } \Psi_0(\eta) = e^{-\gamma\eta} (\approx 1-\gamma\eta)$$

$$\Psi_1(\eta) = \int_0^\eta \Psi_0(t) dt = \frac{1-e^{-\gamma\eta}}{\gamma} (\approx \eta)$$

$$\Psi_2(\eta) = \int_0^\eta \Psi_1(t) dt = \frac{e^{-\gamma\eta}-1+\gamma\eta}{\gamma^2} (\approx \frac{\eta^2}{2})$$

$$\text{and } \tilde{s}_k = (\tilde{s}_{k,1}, \dots, \tilde{s}_{k,n})$$

$$\tilde{s}'_k = (\underbrace{\tilde{s}'_{k,1}, \dots, \tilde{s}'_{k,n}}_{\sim})$$

$$\forall i: (\tilde{s}_{k,i}, \tilde{s}'_{k,i}) \sim \mathcal{N}(0, \Sigma)$$

$$\text{where } \Sigma = \int_0^T \begin{pmatrix} \Psi_0(t) \\ \Psi_1(t) \end{pmatrix} \begin{pmatrix} \Psi_0(t) & \Psi_1(t) \end{pmatrix} dt$$

Theorem: [Dalalyan & Riou-Durand]

Assume  $f$  is  $\alpha$ -strongly convex,  $L$ -smooth

$$(\alpha I \leq \nabla^2 f(x) \leq L I) \quad K = \frac{L}{\alpha}$$

then KL-MC from any  $x_0 \sim \mathcal{S}_0$

$$v_0 \sim \mathcal{N}(0, I)$$

$$\text{with } \gamma \geq \sqrt{L+\alpha}$$

$$\text{satisfies } W_2(\mathcal{S}_k, \mathcal{S}_0) \leq \sqrt{2} \left(1 - \frac{\alpha n}{\gamma}\right)^k W_2(\mathcal{S}_0, \mathcal{S}_0) + K \cdot \eta \sqrt{n}$$

$$\Rightarrow \text{to reach } W_2(\mathcal{S}_k, \mathcal{S}_0) \leq \epsilon$$

$$\text{mixing time of KL-MC is } \tilde{O}\left(\frac{K^{3/2} \cdot \sqrt{n}}{\sqrt{\alpha} \cdot \epsilon}\right)$$

$$\text{compare with ULA: } \tilde{O}\left(\frac{K^2 \cdot n}{\alpha \cdot \epsilon^2}\right) \text{ from coupling}$$

$$\Rightarrow \tilde{O}\left(\frac{K \cdot n}{\alpha \cdot \epsilon^2}\right) \text{ from convex opt.}$$

# References

- Villani, *Hypocoercivity*, Memoirs of the American Mathematical Society, 2009
- Eberle, Guillin, & Zimmer, *Couplings and quantitative contraction rates of Langevin dynamics*, The Annals of Probability, 2019

# References

- Cheng, Chatterji, Bartlett, & Jordan, *Underdamped Langevin MCMC: A non-asymptotic analysis*, COLT 2018
- Shen & Lee, *The Randomized Midpoint Method for Log-Concave Sampling*, NeurIPS 2019
- Dalalyan & Riou-Durand, *On sampling from a log-concave density using kinetic Langevin diffusions*, Bernoulli, 2020
- Cao, Lu, & Wang, *On explicit  $\ell_2$ -convergence rate estimate for underdamped Langevin dynamics*, arXiv 2019
- Ma, Chatterji, Cheng, Flammarion, Bartlett, & Jordan, *Is there an analog of Nesterov acceleration for MCMC?*, arXiv 2019
- Li, Zha, & Tao, *Hessian-Free High-Resolution Nesterov Acceleration for Sampling*, arXiv 2020