

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

---

Andre Wibisono

April 21, 2021

Yale University

## Last time

- Langevin dynamics for sampling from  $\nu \propto e^{-f}$

$$dX_t = -\nabla f(X_t) dt + \sqrt{2} dW_t$$

- Unadjusted Langevin Algorithm

$$X_{u+t} = X_u - \eta \nabla f(X_u) + \sqrt{2\eta} Z_u$$

# Plan

## Variations of Langevin dynamics

$$dX_t = -\nabla f(X_t) dt + \sqrt{2} dW_t$$

### 1. First-order dynamics

Weighted Langevin dynamics

Mirror Langevin dynamics

Newton Langevin dynamics

Interacting Langevin dynamics

### 2. Second-order dynamics

Underdamped Langevin dynamics

Hamiltonian Monte Carlo

# 1. First-order dynamics

---

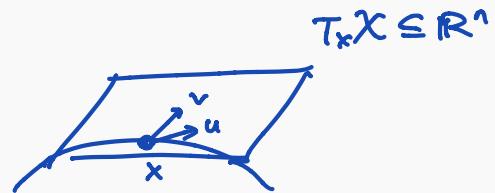
# Optimization 1: Natural gradient flow

On manifold  $X$ , want  $\min_{x \in X} F(x)$

Natural gradient flow  $\dot{x}_t = -\text{grad } F(x_t)$

- if  $X \subseteq \mathbb{R}^n$ , metric (as a matrix)  $g(x) > 0$

$$\underline{\text{NGF:}} \quad \dot{x}_t = -\underbrace{g(x_t)^{-1}}_{M(x_t) > 0} \nabla F(x_t)$$



$$\langle u, v \rangle_x = u^\top g(x) v$$

- dissipative:

$$\begin{aligned} \frac{d}{dt} F(x_t) &= \left\langle \text{grad } F(x_t), \dot{x}_t \right\rangle_{x_t} \\ &= \left\langle g(x_t)^{-1} \nabla F(x_t), -g(x_t)^{-1} \nabla F(x_t) \right\rangle_{x_t} \\ &= -\|g(x_t)^{-1} \nabla F(x_t)\|_{x_t}^2 \\ &= -\|\nabla F(x_t)\|_{g(x_t)^{-1}}^2 \leq 0 \end{aligned}$$

- exp. convergence rate if Gradient Domination:  $F(x) - F(x^*) \leq \frac{1}{2\alpha} \|\nabla F(x)\|_{g(x)}^2$

# Sampling 1: Weighted Langevin Dynamics

For  $M(x) \succ 0$ , Weighted Langevin Dynamics:

$$(WLD) \quad dX_t = (\nabla \cdot M(x_t) - M(x_t) \nabla f(x_t)) dt + \sqrt{2M(x_t)} dW_t$$

• e.g.,  $M(x) = I$ , this is Langevin dynamics

$$dX = -\nabla f(X) dt + \sqrt{2} dW_t$$

- Lemma: (WLD) has  $\nu \propto e^{-f}$  as stationary distribution
- Fokker-Planck eq: for density  $\pi_t$  of  $X_t$

$$\frac{\partial \pi_t}{\partial t} = \nabla \cdot \left( \pi_t M \nabla \log \frac{\pi_t}{\nu} \right)$$

$$\|v\|_M^2 = v^\top M(x)v$$

• de Bruijn's identity:  $\frac{d}{dt} H_\nu(\pi_t) = -\mathbb{E}_{\pi_t} \left[ \left\| \nabla \log \frac{\pi_t}{\nu} \right\|_M^2 \right]$

$$\frac{d}{dt} \mathcal{K}_\nu(\pi_t) = -2 \mathbb{E}_\nu \left[ \left\| \nabla \frac{\pi_t}{\nu} \right\|_M^2 \right]$$

• So if we have M-weighted LSI / Poincaré :

$$(M\text{-LSI}) \quad H_\nu(g) \leq \frac{1}{2\alpha} \mathbb{E}_g \left[ \|\nabla \log \frac{g}{\nu}\|_M^2 \right]$$

$$(M\text{-PI}) \quad \text{Var}_\nu(g) \leq \frac{1}{\alpha} \mathbb{E}_\nu [\|\nabla g\|_M^2]$$

then we have exponential convergence rate along (WLD) :

$$M\text{-LSI} \Rightarrow \quad H_\nu(g_t) \leq e^{-2\alpha t} H_\nu(g_0)$$

$$M\text{-PI} \Rightarrow \quad \chi^2_\nu(g_t) \leq e^{-2\alpha t} \chi^2_\nu(g_0)$$

## Optimization 2: Mirror flow

Hessian manifold:  $g(x) = \nabla^2 \phi(x)$  for some  $\phi: X \rightarrow \mathbb{R}$  convex

Natural gradient flow:  $\dot{x}_t = -g(x_t)^{-1} \nabla f(x_t)$

$$\hookrightarrow \dot{x}_t = -(\nabla^2 \phi(x_t))^{-1} \nabla f(x_t)$$

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

Mirror flow:

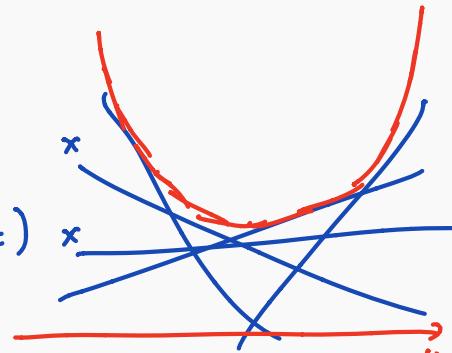
$$\boxed{\frac{d}{dt} \nabla \phi(x_t) = -\nabla f(x_t)}$$

(discretize:

$$\nabla \phi(x_{u+1}) = \nabla \phi(x_u) - \eta \nabla f(x_u)$$

$$\Leftrightarrow x_{u+1} = \nabla \phi^*(\nabla \phi(x_u) - \eta \nabla f(x_u))$$

$$(\text{Mirror Descent}) = \arg \min_x \left\{ \langle \nabla f(x_u), x - x_u \rangle + \frac{1}{\eta} D_\phi(x, x_u) \right\} - \langle \nabla \phi(y), x - y \rangle$$



Convex

$$\text{Dual } \phi^*(y) = \sup_{x \in X} \langle x, y \rangle - \phi(x)$$

$$\nabla \phi^*(y) = \arg \max_{x \in X} \langle x, y \rangle - \phi(x)$$

$$\nabla \phi^* = (\nabla \phi)^{-1}$$

Bregman divergence

$$\hookleftarrow D_\phi(x, y) = \phi(x) - \phi(y)$$

## Sampling 2: Mirror Langevin Dynamics

$$M(x) = \nabla^2 \phi(x)^{-1} \text{ for some } \phi: X \rightarrow \mathbb{R} \text{ convex}$$

Weighted Langevin Dynamics:

$$dx = (\nabla \circ (\nabla^2 \phi(x)^{-1}) - \nabla^2 \phi(x)^{-1} \nabla f(x)) dt + \sqrt{2 \nabla^2 \phi(x)^{-1}} dW$$

↳  $\Rightarrow$  Stationary distribution is  $\nu \propto e^{-f}$

Ito's Lemma (change of variable for SDE):

$$d \nabla \phi(x_t) = - \nabla f(x_t) dt + \sqrt{2 \nabla^2 \phi(x_t)} dW_t$$

Mirror  
Langevin  
Dynamics

$$\gamma = \nabla \phi(x) \sim \tilde{\nu} \text{ if } x \sim \nu$$

$$\text{Stationary distribution is } \tilde{\nu} = (\nabla \phi)_\# \nu$$

Can get convergence rates if have LSI / PI with  $M = (\nabla^2 \phi)^{-1}$

$\Rightarrow$  Mirror LSI

Mirror PI

In discrete time, can discretize Mirror Langevin Dynamics:

$$1) \quad \nabla\phi(x_{u+1}) = \nabla\phi(x_u) - \eta \nabla f(x_u) + \sqrt{2\eta \nabla^2\phi(x_u)} z_u, \quad z_u \sim \mathcal{N}(0, I)$$

$$\Leftrightarrow x_{u+1} = \nabla\phi^*(\nabla\phi(x_u) - \eta \nabla f(x_u) + \sqrt{2\eta \nabla^2\phi(x_u)} z_u)$$

$$= \underset{x \in X}{\operatorname{argmin}} \left\{ \langle \nabla\phi(x_u) - \sqrt{2\eta \nabla^2\phi(x_u)} z_u, x - x_u \rangle + D_\phi(x, x_u) \right\}$$

[Zhang et al. 2020]

2) Mirror Langevin Algorithm [Ahn & Chewi 2020]

$$\nabla\phi(x_{u+\frac{1}{2}}) = \nabla\phi(x_u) - \eta \nabla f(x_u)$$

$$\Leftrightarrow x_{u+\frac{1}{2}} = \underset{x \in X}{\operatorname{argmin}} \left\{ \langle \nabla f(x_u), x - x_u \rangle + \frac{1}{\eta} D_\phi(x, x_u) \right\}$$

$$x_{u+1} = \nabla\phi^*(Y_\eta) \quad \text{where} \quad dY_t = \sqrt{2 \nabla^2\phi^*(Y_t)^{-1}} dw_t$$

$$Y_0 = \nabla\phi(x_{u+\frac{1}{2}})$$

$\Rightarrow$  Convergence analysis using self-concordance of  $\phi$

& relative smoothness + strong convexity of  $f$  wrt  $\phi$

## Optimization 3: Newton flow

Choose  $g(x) = \nabla^2 f(x)$  for  $\min_{x \in \mathbb{R}^n} f(x)$  ( $\phi = f$ )

Natural gradient flow:  $\dot{x}_t = -(\nabla^2 f(x_t))^{-1} \nabla f(x_t)$



Newton flow :  $\frac{d}{dt} \nabla f(x_t) = -\nabla f(x_t)$

$$\Leftrightarrow \nabla f(x_t) = e^{-t} \nabla f(x_0)$$

Ex:  $f(x) = \frac{1}{2} x^T A x$

Convergence  $\rightarrow 0$  exp. fast  
with uniform rate

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

$$\Rightarrow x_t = e^{-At} x_0 \rightarrow 0 \text{ with rate } \lambda_{\min}(A)$$

NF:  $\dot{x}_t = -(\nabla^2 f(x_t))^{-1} \nabla f(x_t)$

$$= -A^{-1} \cdot Ax_t = -x_t$$

$$\Rightarrow x_t = e^{-t} x_0 \rightarrow 0 \text{ with rate } 1.$$

in discrete time:

$$\text{Newton's method: } x_{k+1} = x_k - \eta (\nabla^2 f(x_k))^{-1} \nabla f(x_k)$$

Analysis via self-concordance + barrier properties

# Sampling 3: Newton Langevin Dynamics

Mirror Langevin Dynamics with  $\phi = f$ :

$$(NLD): d\nabla f(x_t) = -\nabla f(x_t) dt + \sqrt{2 \nabla^2 f(x_t)} dW_t$$

- Affine invariance
- Exponential convergence rate for any log-concave  $\nu \propto e^{-f}$

Theorem [Brascamp-Lieb inequality '76]

If strictly convex  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $\nu \propto e^{-f}$  satisfies:

$$\text{Var}_\nu(h) \leq 1 \cdot \mathbb{E}_\nu [\|\nabla h\|_{(\nabla^2 f)^{-1}}^2]$$

$\Leftrightarrow \nu$  satisfies Poincaré ineq. with constant 1  
wrt. metric  $\nabla^2 f$

$\Rightarrow$  Along (NLD),  $X_\nu^2(s_t) \leq e^{-2t} X_\nu^2(s)$ .

# References 1

- Brascamp & Lieb, *On extensions of the Brunn-Minkowski and Prékopa-Leindler theorems, including inequalities for log concave functions, and with an application to the diffusion equation*, *Journal of Functional Analysis*, 1976
- Hsieh, Kavis, Rolland, & Cevher, *Mirrored Langevin dynamics*, NeurIPS 2018
- Chewi, Le Gouic, Lu, Maunu, Rigollet, & Stromme, *Exponential ergodicity of mirror-Langevin diffusions*, NeurIPS 2020
- Zhang, Peyré, Fadili, & Pereyra, *Wasserstein control of mirror Langevin Monte Carlo*, COLT 2020
- Ahn & Chewi, *Efficient constrained sampling via the mirror-Langevin algorithm*, arXiv 2020

## Sampling 4: Interacting Langevin Dynamics

Another way to get affine invacione:

Precondition by covariance

$$dX_t = -C_t \nabla f(X_t) + \sqrt{2 C_t} dW_t$$

where  $C_t = \text{Cov}(S_t) = \mathbb{E}_{S_t} [(X_t - \mu_t)(X_t - \mu_t)^\top]$ ,  $X_t \sim S_t$

$$\mu_t = \mathbb{E}_{S_t} [X_t]$$

Fokker-Planck eq: now nonlinear

$$\frac{\partial \pi_t}{\partial t} = \nabla \cdot (\pi_t C(\pi_t) \nabla \log \frac{\pi_t}{\pi})$$

This is "Kalman-Wasserstein gradient flow"

(metric  $\phi = -\nabla \cdot (\pi C(\pi) \nabla u)$ )

$$\|\phi\|_S^2 = \mathbb{E}_S [\|\nabla u\|_{C(S)}^2]$$



this is mean-field ( $N \rightarrow \infty$ )

### Affine - Invariant Lagrangian Dynamics (ALDI)

have particles  $X_t = (x_t^{(1)}, \dots, x_t^{(N)})$ ,  $x_t^{(i)} \in \mathbb{R}^n$

$$\mathbb{R}^{n \times N}$$

empirical mean  $\mu(X_t) = \frac{1}{N} \sum_{i=1}^N x_t^{(i)}$

covariance  $C(X_t) = \frac{1}{N} \sum_{i=1}^N (x_t^{(i)} - \mu(X_t))(x_t^{(i)} - \mu(X_t))^T$

$$\text{ALDI: } dX_t^{(i)} = -C(X_t) \left( \nabla f(x_t^{(i)}) - \frac{n+1}{2} \nabla \log \det C(X_t) \right) dt$$

for  $i = 1, \dots, N$

$$+ \sqrt{2 C(X_t)} dw_t^{(i)}$$

$$\begin{array}{c} \rightarrow \\ \rightarrow \\ \rightarrow \\ \downarrow \end{array} \quad \Leftrightarrow \quad dX_t^{(i)} = \left( -C(X_t) \nabla f(x_t) + \frac{n+1}{N} (x_t^{(i)} - \mu(X_t)) \right) dt$$

$$+ \sqrt{2 C(X_t)} dw_t^{(i)}$$

## References 2

- Garbuno-Inigo, Hoffmann, Li, & Stuart, *Interacting Langevin Diffusions: Gradient Structure And Ensemble Kalman Sampler*, SIAM Journal on Applied Dynamical Systems, 2020
- Garbuno-Inigo, Nüsken, & Reich, *Affine invariant interacting Langevin dynamics for Bayesian inference*, SIAM Journal on Applied Dynamical Systems, 2020

Example: Gaussian  $v = \mathcal{N}(0, \Sigma)$   $\rightarrow \nabla F(x) = \Sigma^{-1}x$   
 $\nabla^2 F(x) = \Sigma^{-1}$

1) OU (Langevin Dynamics)

$$dx = -\Sigma^{-1}x dt + \sqrt{2} dW$$

$$\Rightarrow x_t \stackrel{d}{=} e^{-\Sigma^{-1}t} x_0 + \sqrt{\Sigma(I - e^{-2\Sigma^{-1}t})} z$$

$$s_t \rightarrow v \text{ with rate } \lambda_{\min}(\Sigma^{-1}) = \frac{1}{\lambda_{\max}(\Sigma)}$$

2) Newton Langevin Dynamics:

$$d\nabla F(x) = -\nabla F(x)dt + \sqrt{2\nabla^2 F(x)} dW$$

$$d(\Sigma^{-1}x) = -\Sigma^{-1}x dt + \sqrt{2\Sigma^{-1}} dW$$

$$dx = -x dt + \sqrt{2\Sigma} dW$$

$$d(e^t x) = e^t (dx + x dt) = \sqrt{2\Sigma} e^t dW$$

$$\Rightarrow x_t \stackrel{d}{=} e^{-t} x_0 + \sqrt{(1-e^{-2t})\Sigma} z$$

$$s_t \rightarrow v \text{ with rate 1.}$$

3) Interacting Langevin Dynamics

If  $x_0 \sim s_0 = \mathcal{N}(0, \Sigma_0)$

Then  $x_t \sim s_t = \mathcal{N}(0, \Sigma_t)$

$$\Sigma_t^{-1} = (\Sigma_0^{-1} - \Sigma^{-1}) e^{-2t} + \Sigma^{-1}$$

$$s_t \rightarrow v \text{ with rate 2.}$$