Recent Developments in Algorithm Design: Sampling from High-Dimensional Distributions

Prof. Christopher Musco, New York University

COMMUNICATION-FREE COUPLING

Draft: NYU is a private research university in the city of New York .



Desired Output: NYU is a private research university in New York City.



NYU is a private research **university**

NYU is a private research university in

NYU is a private research <u>university</u> in <u>New</u>



NYU is a private research university in the city

••••

Issue: Even if then next token distribution for the drafter model, \mathcal{P} , and the product model, \mathcal{Q} are very similar, it could be unlikely for the draft to be correct.



If $a \sim P$ and $b \sim Q$, $\Pr[a = b] \approx \forall 3$

COUPLING

Solution: Coordinate the sampling!

Definition (Coupling)

Let \mathcal{P} and \mathcal{Q} be distributions over $\{1, \ldots, n\}$. A <u>coupling</u> between \mathcal{P} and \mathcal{Q} is any distribution over pairs $(a,b) \in \{1, \ldots, n\} \times \{1, \ldots, n\}$ such that *a*'s marginal distribution is \mathcal{P} and *b*'s marginal distribution is \mathcal{Q} .

Goal: Efficiently sample from a coupling \mathcal{C} between the small and large model distributions which maximizes

$$\left(\mathsf{Pr}[a=b]. \right)$$

Always possible to find a coupling which ensures that $\Pr[a = b] = 1 - D_{TV}(\mathcal{P}, \mathcal{Q}).$

Reminder: For discrete distributions \mathcal{P} and \mathcal{Q} over $\{\underline{1, \dots, n}\}$ represented by length *n* probability vectors $\mathbf{p}, \mathbf{q} \in [0, 1]^n$,

$$D_{TV}(\mathcal{P},\mathcal{Q}) = \left(1 - \sum_{i=1}^{n} \min(p_i, q_i)\right)$$

The following procedure achieves the optimal bound of $P_{\underline{r}[\underline{a}=\underline{b}]} = 1 - D_{TV}(\mathcal{P}, \mathcal{Q})$.

Drafter:

 $\mathbf{f} \cdot \text{Sample} (\mathbf{a} \sim \mathcal{P})$. Sends both a and \mathbf{p} to FullModel. (Full Model:) $\mathbf{h} \sim \mathbf{Q}$

• Await (*a*, **p**) from Drafter.

With probability min(1,
$$q_a/p_a$$
) return $\underline{b} = a$.

• Otherwise, sample *b* from $Q' = \{q'_1, \dots, q'_n\}$, where:

$$q'_i = \frac{\max(0, q_i - p_i)}{\sum_{i=1}^n \max(0, q_j - p_j)}$$

Important that the Drafter could sample without knowing the Full Model's Distribution, q! There is only "one-way communication".

Pla=57

Is it posssible to do anything with <u>(no communication</u>) between the samplers?

Why would we care? The output of the Full Model is always sampled from Q, but the exact value sampled <u>depends on the</u> Drafter distribution P.

(Cannot immediately verify that adding speculative) decoding did not change the model distribution.

' If drafter changes, model output is not deterministic from the user's point of view given a fixed random seed.

("Coupling without Communication and Drafter-Invariant Speculative Decoding" [Daliri, Musco, Suresh) ISIT 2025]. Basically the same idea appeared in:

- Anari, Gao, Rubinstein, STOC 2024 🕯
- Liu, Yin, STOC 2022
- · Bavarian, Ghazi, Haramaty, Kamath, Rivest, Sudan, 2020.

WEIGHTED MINHASH COUPLING



WEIGHTED MINHASH COUPLING

£

Claim:
$$\Pr[a = b] \ge \frac{\sum_{i=1}^{n} \min(p_i, q_i)}{\sum_{i=1}^{n} \max(p_i, q_i)} = \frac{1 - \mathcal{D}_{1V}(P, \Delta)}{1 + \mathcal{D}_{1V}(P, \Delta)} \le 1 - 2 \mathcal{D}_{1V}(P, \Delta)$$

10

COMMUNICATION FREE COUPLING

Optimal Coupling:

$$\frac{\overline{f}}{p_{i-1}} \operatorname{wirw}(p_{i}, g_{i}) + \overline{f} \operatorname{wirw}(p_{i}, g_{i}) = 2$$

$$\operatorname{Pr}[a = b] = 1 - D_{TV}(\mathcal{P}, \mathcal{Q}) = 2 - \overline{f} \operatorname{wirw}(p_{i}, g_{i})$$

$$\operatorname{Pr}[a = b] \geq \frac{\sum_{i=1}^{n} \min(p_{i}, q_{i})}{\sum_{i=1}^{n} \max(p_{i}, q_{i})} = \frac{1 - D_{TV}(\mathcal{P}, \mathcal{Q})}{1 + D_{TV}(\mathcal{P}, \mathcal{Q})} \quad \frac{1}{2}$$

Takeaway: Pay very little for no communication!

Possible to show that this is optimal. No communication-free protocol can achieve for all distributions:

$$\Pr[a = b] > \frac{1 - D_{TV}(\mathcal{P}, \mathcal{Q})}{1 + D_{TV}(\mathcal{P}, \mathcal{Q})}.$$

′ [Bavarian, Ghazi, Haramaty, Kamath, Rivest, Sudan, 2020]. 🖊

Fix public random variables $\underline{u_1}, u_2, \ldots \sim \text{Unif}[0, 1]$.

Drafter:

• Return
$$\underline{a} = \arg\min_{i \in \{1,...,n\}} \frac{-\ln(u_i)}{p_i}$$
.

Full Model:

• Return
$$b = \arg\min_{i \in \{1,...,n\}} \frac{-\ln(u_i)}{q_i} = \alpha/\beta m - \frac{\ln(1/u_i)}{q_i}$$

This is already how samples are typically obtained! In جمع کر کرد کرد particular, standard to use the "Gumbel Max Trick": - مرج مرامی (هن)

$$\int_{i \in \{1,\dots,n\}} b = \arg\max_{i \in \{1,\dots,n\}} \left[\ln(q_i) - \ln(\ln(1/u_i)) \right].$$

Not too hard to check that $a \sim \mathcal{P}$ and $b \sim \mathcal{Q}$.

Gumbel sampling gives a <u>pareto improvement</u> over weighted MinHash.

Theorem (Daliri, Musco, Suresh, ISIT 2025)

For any two distributions \mathcal{P}, \mathcal{Q} ,

$$\Pr_{(a,b)\sim Gumbel}[a=b] \ge \Pr_{(a,b)\sim MinHash}[a=b],$$

and there exist distributions where inequality is strict.

Question one group is studying for the project: Is Gumbel pareto optimal?

EFFICIENT SAMPLING IN HIGH-DIMENSIONS

Increasingly common goal in machine learning: Sample from a distribution over \mathbb{R}^d with density $p(\mathbf{x})$ η((μ, Ξ) p(x) - e^{-(x-μ)⁵Ξ(x-μ)} f(x)= (x-~) [x-~) $p(\mathbf{x})$

Assume $p(\underline{x}) \propto \exp(-f(\underline{x}))$ for some function $f : \mathbb{R}^d \to \mathbb{R}$ and that we are given gradient oracle access to $\nabla f(\underline{x})$.

What I hope to cover:

- Where/why does this problem arise in machine learning?
- What is the <u>(stochastic) gradient Langevin dynamics</u> algorithm and why does is work?
- Where is the area headed / where are opportunities for algorithms research?

GRADIENT REMINDER

Recall that $\nabla f : \mathbb{R}^d \to \mathbb{R}^d$ returns the vector of partial derivatives at a point **x**:

$$\nabla f(\mathbf{x}) = \begin{bmatrix} \frac{\partial}{\partial X_1} f(\mathbf{x}) \\ \vdots \\ \frac{\partial}{\partial X_d} f(\mathbf{x}) \end{bmatrix} ,$$

The gradient determines the instantanious change in *f*'s value with respect to changes in the input variables:

$$\lim_{h\to 0}\frac{f(\mathbf{x}+h\mathbf{v})-f(\mathbf{x})}{h}=\langle \nabla f(\mathbf{x}),\mathbf{v}\rangle.$$

Where do gradients show up in machine learning?

Let $M_{\mathbf{x}} : \mathbb{R}^d \to \mathbb{R}$ be a model parameterized by **x**. Given a labeled dataset $(\mathbf{a}_1, b_1), \ldots, (\mathbf{a}_n, b_n)$, goal in supervised learning is to find parameters such that:

 $M_{\mathbf{x}}(\mathbf{a}_i) \approx b_i.$

Typically accomplished by writing down some loss function $f(\mathbf{X})$ and minimizing. For example, least squares loss:

$$h_{0}$$
, h_{in} , h_{in} , $f(\underline{x}) = \sum_{i=1}^{n} (M_{\underline{x}}(\underline{a}_i) - \underline{b}_i)^2$.

Goal: Find $\underline{x}^* = \arg \min_{x} f(x)$.

GRADIENT DESCENT

Most common algorithm to do so: gradient descent.

- Choose starting point $\underline{x}_0 \in \mathbb{R}^d$, step size $\underline{\eta}$.
- For $t = 0, \ldots, T$

•
$$\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t - \eta \cdot \nabla f(\mathbf{x}_t).$$



Gradient descent:

• Choose starting point $\mathbf{x}_0 \in \mathbb{R}^d$, step size η .

• For
$$t = 0, ..., T$$

• $\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t - \eta \cdot \nabla f(\mathbf{x}_t)$.
 $\mathbf{X} + \mathcal{M} \cdot \mathbf{V}$

Justification: We want to make a small change, $\eta \cdot \mathbf{v}$ to \mathbf{x}_t that decreases the value of f. we know the value of f with the value of f and $\mathbf{x}_t = \mathbf{v} \cdot \mathbf{v}$.

Choosing $\mathbf{v} = -\nabla f(\mathbf{x})$ ensures that, if we take $\eta \to 0$,

WHY GRADIENT DESCENT

- Simple and general. We only need to implement a gradient oracle for computing ∇f(x). For almost all models with <u>d</u> parameters, can be done in O(<u>nd</u>) time.
- Stochastic approximation of gradient is even faster. Typically O(d) time: $f(x) = \frac{1}{2} \mathcal{L}(x, \alpha; \beta; \beta; \beta)$

$$\nabla f(\mathbf{x}) = \sum_{i=1}^{n} \nabla \ell(\mathbf{x}, \mathbf{a}_i, b_i).$$

- Guaranteed to converge to a stationary point (e.g., local min) of *f* for sufficiently small step size.
- Dimension independent convergence rates can be obtained under mild assumptions.



¹Other methods (e.g., Center-of-Gravity Method) can achieve a better dependence on ϵ , but at the cost of a dependence on *d*.

Huge amount of algorithmic research centered around gradient descent and its variants.

- Acceleration/momentum to speed up convergence.
- Generalized steppest descent, mirror descent, etc.
- Stochastic gradient methods, variance reduction.
- Preconditioning, quasi-second order methods, adaptive step size methods.
- Lower bounds (e.g, in first order oracle model).

Where does least squares loss comes from?

Assume fixed dataset $\underline{a}_1, \ldots, \underline{a}_n$ with targets generated from ground truth model, $\underline{M}_{\mathbf{x}^*}$ plus Gaussian noise:

Would like to choose params. <u>most likely</u> to have generated the targets we observed. Likelihood of data given parameters:

$$\underline{L(\mathbf{x})} = \underline{p(b_1, \ldots, b_n \mid \underline{\mathbf{x}})} \propto \prod_{i=1}^n \exp\left(-\frac{(b_i - M_{\mathbf{x}}(\mathbf{a}_i))^2}{2\sigma^2}\right).$$

E

Goal: Compute the maximium likelihood estimator (MLE):

$$\underline{\mathbf{x}^{*}} = \arg \max_{\mathbf{x}} L(\mathbf{x}). \quad \text{is or } \mathbf{y} \text{ or } [\mathbf{v}_{i}(\mathbf{L}(\mathbf{x}))]$$
quivalent to minimizing the negative log-likelihood:
$$f(\mathbf{x}) = -\log L(\mathbf{x}) = \prod_{i=1}^{n} (b_{i} - M_{\mathbf{x}}(\mathbf{a}_{i}))^{2} + \operatorname{const.}$$

Most standard ML loss functions are negative log-likelihoods for some other data generation process, including logistic/cross-entropy loss, ℓ_1 loss, etc.

BAYESIAN MACHINE LEARNING

One step further: Assume prior distribution over parameters **x**.
E.g.
$$x_i \sim \mathcal{N}(0, \gamma^2)$$
 for all *i*.
Lets us define a posterior probability of **x** given the data:
 $p(\mathbf{x}, b_1, \dots, b_n) = \frac{p(b_1, \dots, b_n \mid \mathbf{x}) \cdot p(\mathbf{x})}{p(b_1, \dots, b_n)} = \frac{\text{likelihood} \cdot \text{prior}}{\text{evidence}}.$

Goal: Compute the <u>maximum a posteriori</u> (<u>MAP)</u> estimator:

$$\mathbf{x}^{*} = \arg \max_{\mathbf{x}} \underbrace{p(\mathbf{x} \mid b_{1}, \dots, b_{n})}_{\mathbf{x}}$$

$$= \arg \max_{\mathbf{x}} \underbrace{p(b_{1}, \dots, b_{n} \mid \mathbf{x}) \cdot \underline{p(\mathbf{x})}}_{\mathbf{x}}$$

$$\operatorname{arg}_{\mathbf{x}}^{\mathsf{arg}} \underbrace{f_{\bullet}}_{\mathbf{x}} \left(\underbrace{p(b_{1}, \dots, b_{n} \mid \mathbf{x})}_{\mathbf{x}} \right) * \operatorname{leg}_{\mathbf{x}}^{\mathsf{arg}} \right)$$

$$\operatorname{arg}_{\mathbf{x}}^{\mathsf{arg}} \underbrace{- (11)}_{\mathbf{x}}$$

25

BAYESIAN MACHINE LEARNING

Again, can equivalently minimize the negative log-posterior:

$$f(\mathbf{x}) = -\log(p(b_1,\ldots,b_n \mid \mathbf{x})) - \log(p(\mathbf{x})).$$

Example: Least squares loss with Gaussian prior.

$$p(b_1, \dots, b_n \mid \mathbf{x}) \propto \prod_{i=1}^n \exp\left(-\frac{(b_i - M_{\mathbf{x}}(\mathbf{a}_i))^2}{2\sigma^2}\right) \mathbf{x}$$

$$\mathbf{x}_i \sim \mathbf{N}(\mathbf{o}, \mathbf{y})$$

$$f(\mathbf{x}) = \frac{1}{2\sigma^2} \sum_{i=1}^n (b_i - M_{\mathbf{x}}(\mathbf{a}_i))^2 + \left(\frac{1}{2\gamma^2} \sum_{i=1}^d x_i^2\right)$$

$$= \left(\frac{1}{2\sigma^2} \sum_{i=1}^n (b_i - M_{\mathbf{x}}(\mathbf{a}_i))^2\right) + \frac{1}{2\gamma^2} ||\mathbf{x}||_2^2.$$

Optimization (usually solved with gradient descent) computes the <u>mode</u> of the posterior distribution $p(\mathbf{x} | b_1, ..., b_n)$.



Another important goal: Sample parameter vector **x** from the posterior distribution. I.e., sample $\mathbf{x} \sim c \cdot e^{-f(\mathbf{x})}$ given a gradient oracle for f

$$C = \int_{1 \in \mathbb{R}^{d}} e^{-f(x)}$$

BAYESIAN MACHINE LEARNING

Original Bayesian motivation: Confidence intervals and uncertainty quantification. For new data point \mathbf{a}_{n+1} with unknown label y_{n+1} , can sample from $p(y_i | b_1, ..., b_n)$ by sampling **x** from posterior and computing $M_{\mathbf{x}}(\mathbf{a}_{n+1})$.



For simple models (linear models, GLMS, kernel or Gaussian process regression, etc.) we have model-specific methods to sample from the posterior or compute confidence intervals.

Goal: Extend posterior sampling to any model that we can efficiently compute the gradient of (e.g., neural <u>networks</u>).

Why work with negative log posterior instead of directly working with posterior?

$$\prod_{i=1}^{n} \exp\left(-\frac{(b_i - M_{\mathbf{x}}(\mathbf{a}_i))^2}{2\sigma^2}\right) \quad \text{vs.} \quad \sum_{i=1}^{n} (b_i - M_{\mathbf{x}}(\mathbf{a}_i))^2$$

UNADJUSTED LANGEVIN ALGORITHM

? P, Pr ... (Pr) Unadjusted Langevin algorithm to sample from $e^{-f(x)}$:

• Choose starting point $\underline{\mathbf{x}_0} \in \mathbb{R}^d$, step size $\underline{\eta}$.



 Widely used throughout computational science, statistics, and other fields since at least the 1990s.

 Bayesian Learning for Neural Networks Radford M. Neal
 A thesis submitted in conformity with the requirements for the degree of Doctor of Philosophy, Graduate Department of Computer Science, in the University of Toronto Convocation of March (1995)

Bayesian Learning via Stochastic Gradient Langevin Dynamics 261

 (Max Welling)
 WELLING@ICS.UCI.EDU

 D. Bren School of Information and Computer Science, University of California, Irvine, CA 92697-3425, USA
 Yee Whye Teh

 Yee Whye Teh
 YWTEH@GATSBY.UCL.AC.UK

 Gatsby Computational Neuroscience Unit, UCL, 17 Queen Square, London WC1N 3AR, UK

Unadjusted Langevin algorithm:

- Choose starting point $\mathbf{x}_0 \in \mathbb{R}^d$, step size η ,
- For t = 0, ..., T
 - Sample $\mathbf{g}_t \sim \mathcal{N}(0, I)$.
 - $\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t \eta \cdot \nabla f(\mathbf{x}_t) + \sqrt{2\eta} \cdot \mathbf{g}_t$.

Informal claim: For $\eta \rightarrow 0$, the distribution of \underline{x}_t converges to $c \cdot e^{-f(\mathbf{x})}$ for many natural distributions.

Suffices for *f* to be is strongly convex mixture of distributions) with this property, anything with Poincaré inequality, etc.

Non-asymptotic convergence rates have only been proven relatively recently, starting with [Durmus, Moulines, 2017].

Distribution 14 log concove

MORE RECENT MOTIVATION: GENERATIVE AI

We have seen insane progress in (conditional) image generation.



Just a few years ago, we got excited about images like:



Older methods: Variational Auto-Encoders, Generative Adversarial Networks) normalizing flows, energy-based models, etc.

(Leading modern methods: Denoising Diffusion Models, Score-based Generative Models . الالجان سالجس لموا"

Denoising Diffusion Probabilistic Models



ENERGY-BASED MODELS

View image generation as a sampling problem $p(\mathbf{x})$ is the distribution over "natural images". Want to sample from p, or p conditioned on some prompt.

Energy-based models: Train model \underline{M}_{θ} that takes in an image **x** and returns a negative log probability. Given a training set of images $\mathbf{x}_1, \ldots, \mathbf{x}_n$, goal is to minimize: $f(\mathbf{x})$

$$\sum_{i=1} \underline{\mathcal{M}_{\theta}(\mathbf{x}_i)} = \prod_{i=1} e^{-\underline{\mathcal{M}_{\theta}(\mathbf{x}_i)}}.$$

Can sample new images using Langevin dynamics, where $f(\mathbf{x}) = M_{\theta}(\mathbf{x}_i)$.

Lots of issues with normalization... how do you ensure M_{θ} models a normalized probability density?

SCORE-BASED MODELS

Train a model that <u>directly predicts</u> $\nabla(-\log(p(\mathbf{x})))$. This is $\nabla f(\mathbf{x})$ called the <u>(score function</u>) but it is no different from the gradient ∇f we needed to implement Langevin dynamics. (How to train the model without input/output pairs?) $(\mathbf{x}_i, \nabla(-\log(p(\mathbf{x}_i))))$





Intuitivtely,
$$\left(\nabla f(\mathbf{x} + \mathbf{n}) \right) = -\mathbf{n}$$
.

Unadjusted Langevin algorithm:



- Choose starting point $\mathbf{x}_0 \in \mathbb{R}^d$, step size η .
- For t = 0, ..., T
 - Sample $\mathbf{g}_t \sim \mathcal{N}(0, I)$.

•
$$\mathbf{X}_{t+1} \leftarrow \mathbf{X}_t - \eta \cdot \nabla f(\mathbf{X}_t) + \sqrt{2\eta} \cdot \mathbf{g}_t.$$

Informal claim: For $\eta \to 0$, the distribution of \mathbf{x}_t converges to $c \cdot e^{-f(\mathbf{x})}$ for <u>many</u> natural distributions.



Continuous-time Langevin Dynamics: Typical analysis begins by considering the continuous-time limit of the unadjusted (Langevin algorithm as $\eta \rightarrow 0$.)



To do so, we need to define a Brownian motion which is the continuous limit of a Gaussian random walk.

(One-dimensional) Brownian motion: A Brownian motion $\underline{B_t}$ is a continuous function of time $t \ge 0$ with the properties:



STEADY-STATE OF LANGEVIN EQUATION



 $\begin{pmatrix} \text{Let } p_t \text{ be the distribution of } X_t. \text{ The goal is to show that} \\ \frac{d}{dt} p_t(X) = 0 \text{ for all } X \in \mathbb{R} \text{ when } p_t(X) = c \cdot e^{-f(X)}. \end{cases}$

Fokker-Planck Equation: If
$$dX_t = z(X_t)dt + \sqrt{2}dB_t$$
, then

$$\frac{d}{dt}p_t(X) = -\frac{d}{dX}[z(X)p_t(X)] + p_t''(X)$$

$$= -z'(X)p_t(X) - z(X)p_t'(X) + p_t''(X)$$

 $(\underline{p}'_t \text{ and } \underline{p}''_t \text{ denote the derivatives of } p_t \text{ with respect to X.})$



Fokker-Planck Equation: If $dX_t = z(X_t)dt + \sqrt{2}dB_t$, then $\frac{d}{dt}p_t(X) = -z'(X)p_t(X) - z(X)p'_t(X) + p''_t(X)$ $= f''(x)p_t(x) + f'(x)p_t'(x) + p_t''(x)$ **Claim:** If $\underline{z(X)} = -\underline{f'(X)}$ for some *f*, then $\frac{d}{dt}\overline{p_t(X)} = \underline{0}$ for all $X \in \mathbb{R}$ when $p_t(X) = \underline{c \cdot e^{-f(X)}}$. Proof: $p_{*}'(x) = -f'(x) \cdot ce^{-f(x)} = -f'(x) p_{*}(x)$ $p_{+}^{*}(x) = -f'(x) \cdot -f'(x) \cdot (e^{-f(x)} + -f''(x)) \cdot (e^{-f(x)})$ $= (f'(x))^2 \cdot p_{+}(x) - f''(x) \cdot p_{+}(x)$ $\frac{1}{4t} p_{+}(x) = f''(x)p_{+}(x) + f'(x)(-f'(x))p_{+}(x) + (f'(x))^{2}p_{+}(x) - f'''(x)p_{+}(x)$ = O

/ - f'(x)

$$\frac{d}{dt}p_t(X) = p_t''(X)$$

Adding Gaussian noise at each time step smooths the distribution. Can be thought of as a moving average.





$$\frac{d}{dt}p_t(X) = p_t''(X)$$

$$p_{1+L}(x) \stackrel{1}{\to} \int_{x-q}^{x+a} p_{+}(x) dx$$

How does moving avergage change $p_t(X)$?



$$\frac{d}{dt}p_t(X) = p_t''(X)$$

How does moving avergage change $p_t(X)$?



FOKKER-PLANCK INTUITION

Diffusion-only Fokker-Planck Equation: If $dX_t = \sqrt{2}dB_t$, then

$$\frac{d}{dt}p_t(X) = p_t''(X)$$

$$\frac{d}{dt}p_t(X) = \lim_{h \to 0} \frac{p_{t+h}(X) - p_t(X)}{h}$$

Very informal argument:

$$p_{t+h}(x) \approx \frac{2}{\sqrt{h}} \int_{x-\sqrt{h}}^{x+\sqrt{h}} p_t(y) dy$$

$$\approx \frac{2}{\sqrt{h}} \int_{x-\sqrt{h}}^{x+\sqrt{h}} p_t(x) + p'_t(x)(y-x) + \frac{1}{2}p''_t(x)(y-x)^2 + o(h) dy$$

$$p_{t+h}(x) - p_t(x) \approx \frac{2}{\sqrt{h}} \int_{x-\sqrt{h}}^{x+\sqrt{h}} \frac{1}{2} p_t''(x) (y-x)^2 \, dy$$

= $\frac{1}{\sqrt{h}} p_t''(x) \frac{1}{3} (y-x)^3 \Big|_{x-\sqrt{h}}^{x+\sqrt{h}} \propto p_t''(x) h.$

45



Drift-only Fokker-Planck Equation: If $dX_t = z(X_t)dt$, then



Drift-only Fokker-Planck Equation: If $dX_t = z(X_t)dt$, then



$$dX_t = -f'(X_t)dt + \sqrt{2}dB_t$$

Claim: The distribution $c \cdot e^{-f(X)}$ is an invariant distribution of the Langevin SDE. If $X_0 \sim c \cdot e^{-f(X)}$, then $X_t \sim c \cdot e^{-f(X)} \forall t > 0$.

To get meaningful algorithmic results, need to show:

- 1. (Fast) convergence to this invariant distribution.
- 2. Discretization argument to show that the discrete-time ULA also converges close ot the invariant distribution.

Flavor of result people are interested in proving:

Theorem (See e.g., Chewi 2024)

Suppose f is an α -smooth, β -strongly convex function with condition number $\kappa = \alpha/\beta$. Then after:

 $T = \tilde{O}(\kappa d/\epsilon^2)$ iterations,

the unadjusted Langevin algorithm returns a sample from a distribution \mathcal{P} satisfying:

$$(W_2(\mathcal{P}, c \cdot e^{-f(\mathbf{x})})) \succeq \epsilon,$$

where W₂ is the Wasserstein-2 distance.

Suppose we want to sample from
$$\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$
. Let $\mathbf{H} = \boldsymbol{\Sigma}^{-1}$.
 $p(\mathbf{x}) \propto \exp(-f(\mathbf{x}))$ where $f(\mathbf{x}) = \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{H}(\mathbf{x} - \boldsymbol{\mu})$

$$\nabla f(\mathbf{x}) = \mathsf{H}(\mathbf{x} - \boldsymbol{\mu}).$$

Unadjusted Langevin algorithm:

•
$$\mathbf{x}_t \leftarrow \mathbf{x}_{t-1} - \eta \cdot \nabla f(\mathbf{x}_{t-1}) + \sqrt{2\eta} \cdot \mathbf{g}_{t-1}$$
.

If we initialized \mathbf{x}_0 as a Gaussain, then every iterate is Gaussian distributed. I.e. $\mathbf{x}_t \sim \mathcal{N}(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t)$. Want to show:

$$oldsymbol{\mu}_t
ightarrow oldsymbol{\mu}$$
 $oldsymbol{\Sigma}_t
ightarrow oldsymbol{\Sigma}$

Suppose we want to sample from $\mathcal{N}(\mu, \mathbf{\Sigma})$. Let $\mathbf{H} = \mathbf{\Sigma}^{-1}$. $p(\mathbf{x}) \propto \exp(-f(\mathbf{x}))$ where $f(\mathbf{x}) = \frac{1}{2}(\mathbf{x} - \mu)^{\mathsf{T}}\mathbf{H}(\mathbf{x} - \mu)$

$$\nabla f(\mathbf{x}) = \mathsf{H}(\mathbf{x} - \boldsymbol{\mu}).$$

Unadjusted Langevin algorithm:

•
$$\mathbf{x}_t \leftarrow \mathbf{x}_{t-1} - \eta \cdot \nabla f(\mathbf{x}_{t-1}) + \sqrt{2\eta} \cdot \mathbf{g}_{t-1}$$
.

Plugging in definition of gradient:

$$\begin{aligned} \mathbf{x}_t &= \mathbf{x}_{t-1} - \eta \mathsf{H}(\mathbf{x}_{t-1} - \boldsymbol{\mu}) + \sqrt{2\eta} \cdot \mathbf{g}_{t-1} \\ (\mathbf{x}_t - \boldsymbol{\mu}) &= (\mathbf{x}_{t-1} - \boldsymbol{\mu}) - \eta \mathsf{H}(\mathbf{x}_{t-1} - \boldsymbol{\mu}) + \sqrt{2\eta} \cdot \mathbf{g}_{t-1} \\ &= (\mathsf{I} - \eta \mathsf{H})(\mathbf{x}_{t-1} - \boldsymbol{\mu}) + \sqrt{2\eta} \cdot \mathbf{g}_{t-1} \end{aligned}$$

Unrolling the iteration:

$$\begin{aligned} (\mathbf{x}_t - \boldsymbol{\mu}) &= (\mathbf{I} - \eta \mathbf{H})(\mathbf{x}_{t-1} - \boldsymbol{\mu}) + \sqrt{2\eta} \cdot \mathbf{g}_{t-1} \\ &= (\mathbf{I} - \eta \mathbf{H})((\mathbf{I} - \eta \mathbf{H})(\mathbf{x}_{t-2} - \boldsymbol{\mu}) + \sqrt{2\eta} \cdot \mathbf{g}_{t-2}) + \sqrt{2\eta} \cdot \mathbf{g}_{t-1} \\ &\vdots \\ &= (\mathbf{I} - \eta \mathbf{H})^t (\mathbf{x}_0 - \boldsymbol{\mu}) + \sqrt{2\eta} \sum_{i=0}^{t-1} (\mathbf{I} - \eta \mathbf{H})^{t-1-i} \mathbf{g}_i \end{aligned}$$

First observation: If we choose $\eta = 1/\lambda_{\max}(\mathbf{H})$, then $\|\mathbb{E}[\mathbf{x}_t] - \boldsymbol{\mu}\|_2 \le \epsilon \|\mathbb{E}[\mathbf{x}_0] - \boldsymbol{\mu}\|_2$ after $t = O(\kappa \log(1/\epsilon))$ iterations.

In other words, we quickly converge to a Gaussian distribution with the correct mean!

$$(\mathbf{x}_t - \boldsymbol{\mu}) = (\mathbf{I} - \eta \mathbf{H})^t (\mathbf{x}_0 - \boldsymbol{\mu}) + \sqrt{2\eta} \sum_{i=0}^{t-1} (\mathbf{I} - \eta \mathbf{H})^{t-1-i} \mathbf{g}_i$$

First observation: If we choose $\eta = 1/\lambda_{\max}(\mathsf{H})$, then $\|\mathbb{E}[\mathsf{x}_t] - \boldsymbol{\mu}\|_2 \le \epsilon \|\mathbb{E}[\mathsf{x}_0] - \boldsymbol{\mu}\|_2$ after $t = O(\kappa \log(1/\epsilon))$ iterations. What about the covariance matrix?

$$\begin{split} \boldsymbol{\Sigma}_{t} &= \mathbb{E}\left[(\mathbf{X}_{t} - \boldsymbol{\mu}_{t})(\mathbf{X}_{t} - \boldsymbol{\mu}_{t})^{\mathsf{T}} \right] \\ &= \left((\mathbf{I} - \eta \mathbf{H})^{t}(\mathbf{X}_{0} - \boldsymbol{\mu}) + \sqrt{2\eta} \sum_{i=0}^{t-1} (\mathbf{I} - \eta \mathbf{H})^{t-1-i} \mathbf{g}_{i} - (\boldsymbol{\mu}_{t} - \boldsymbol{\mu}_{t}) \right) \left(\cdots \right)^{\mathsf{T}} \\ &= \left((\mathbf{I} - \eta \mathbf{H})^{t}(\mathbf{X}_{0} - \boldsymbol{\mu}_{0}) + \sqrt{2\eta} \sum_{i=0}^{t-1} (\mathbf{I} - \eta \mathbf{H})^{t-1-i} \mathbf{g}_{i} \right) \left(\cdots \right)^{\mathsf{T}} \end{split}$$

Basically all cross-terms cancel. If we assume $x_0 \sim \mathcal{N}(\mu_0, l)$, we get:

$$\Sigma_{t} = (I - \eta H)^{2t} + 2\eta \sum_{i=0}^{t-1} (I - \eta H)^{2t-2-2i}$$
$$= (I - \eta H)^{2t} + 2\eta \sum_{i=0}^{t-1} (I - \eta H)^{2i}$$

CONVERGENCE OF COVARIANCE

$$\boldsymbol{\Sigma}_t = (\mathbf{I} - \eta \mathbf{H})^{2t} + 2\eta \sum_{i=0}^{t-1} (\mathbf{I} - \eta \mathbf{H})^{2i}$$

We have that $\sum_{i=0}^{\infty} \mathbf{A}^i = (\mathbf{I} - \mathbf{A})^{-1}$ and thus:

$$\sum_{i=0}^{t-1} \mathsf{A}^i = (\mathsf{I} - \mathsf{A})^{-1} - \mathsf{A}^t (\mathsf{I} - \mathsf{A})^{-1}.$$

Apply to
$$\mathbf{A} = (\mathbf{I} - \eta \mathbf{H})^2 = \mathbf{I} - 2\eta \mathbf{H} + \eta^2 \mathbf{H}^2$$

 $\mathbf{\Sigma}_t = 2\eta \left(2\eta \mathbf{H} - \eta^2 \mathbf{H}\right)^{-1} - (\mathbf{I} - \eta \mathbf{H})^{2t} \left(2\eta \mathbf{H} - \eta^2 \mathbf{H}\right)^{-1} + (\mathbf{I} - \eta \mathbf{H})^{2t}$
 $= (\mathbf{H} + .5\eta \mathbf{H}^2)^{-1} - (\mathbf{I} - \eta \mathbf{H})^{2t} \left(2\eta \mathbf{H} - \eta^2 \mathbf{H}\right)^{-1} + (\mathbf{I} - \eta \mathbf{H})^{2t}$
 $\approx \mathbf{H}^{-1}$

for small enough η .

Can we accelerate convergence using the existing toolkit of optimization tricks?

• Acceleration/momentum, preconditioning, variance reduction, etc.

Can we take advantage of additional oracles, e.g. that can draw samples $\mathbf{x} \sim e^{-f(\mathbf{x})}$?

• See e.g. [Koehler, Vuong, 2023]

Lower bounds on gradient oracle complexity?

• See e.g. [Chewi, de Dios Pont, Li, Lu, Narayanan, 2024]