

CPSC 661: Sampling Algorithms in ML

Andre Wibisono

February 17, 2021

Yale University

Last time

- Reversible Markov chain $\nu(x) \cdot P_x(y) = \nu(y) \cdot P_y(x)$
- Spectral gap \Leftrightarrow Conductance $\frac{\phi^2}{2} \leq \gamma \leq 2\phi$
 $\gamma = \lambda_2(L)$ ϕ
- s -Conductance \Rightarrow Mixing time in TV distance

$$0 \leq s \leq \frac{1}{2}$$

$$\tau(\varepsilon) = \tilde{O}\left(\frac{1}{\phi_s^2}\right), \quad s = \frac{\varepsilon}{2M}$$

$$M = M_{\nu}^{\infty}(s_0) = \sup_x \left| \frac{s_0(x)}{\nu(x)} - 1 \right|$$

Questions:

1. How to construct reversible Markov chain P ?
 \Rightarrow Today: Metropolis-Hastings
2. How to bound s -conductance? (next time)

References

- Dwivedi, Chen, Wainwright, and Yu, *Log-Concave Sampling: Metropolis-Hastings Algorithms are Fast*, Journal of Machine Learning Research, 2019
- Billera and Diaconis, *A Geometric Interpretation of the Metropolis-Hastings Algorithm*, Statistical Science, 2001

Metropolis Algorithm

THE JOURNAL OF CHEMICAL PHYSICS

VOLUME 21, NUMBER 6

JUNE, 1953

Equation of State Calculations by Fast Computing Machines

NICHOLAS METROPOLIS, ARIANNA W. ROSENBLUTH, MARSHALL N. ROSENBLUTH, AND AUGUSTA H. TELLER,
Los Alamos Scientific Laboratory, Los Alamos, New Mexico

AND

EDWARD TELLER,* *Department of Physics, University of Chicago, Chicago, Illinois*

(Received March 6, 1953)

THOSE WE'VE LOST

Arianna Rosenbluth Dies at 93; Pioneering Figure in Data Science

Dr. Rosenbluth, who received her physics Ph.D. at 21, helped create an algorithm that has become a foundation of understanding huge quantities of data. She died of complications of the coronavirus.

By Katie Hafner

Published Feb. 9, 2021 Updated Feb. 15, 2021

The Metropolis algorithm, a technique for generating random samplings, started out as a way to understand a fundamental problem: how atoms rearrange themselves as solids melt.

Over the decades, the Metropolis algorithm and its subsequent variations have been put to a vast number of uses and now serve as an underpinning to understanding critical challenges of our age, including making sense of huge volumes of data, predicting election outcomes and understanding Covid-19's spread.



Dr. Arianna Wright Rosenbluth in 2013. She helped create what has become one of the most important algorithms of all time. via Rosenbluth family

Metropolis Algorithm

To sample from ν on X :

- Start with a symmetric Markov chain $P = (P_x : x \in X)$

$$P_x(y) = P_y(x) \quad \leftarrow \begin{array}{l} \text{stationary dist for } P \\ \text{is uniform} \end{array}$$

- Do random walk with accept/reject step:

- From x , draw $y \sim \underline{P_x}$ proposal

- Accept y with probability $\min \left\{ 1, \frac{\nu(y)}{\nu(x)} \right\}$

Metropolis filter else reject y

new point $x' = \begin{cases} y \sim P_x \text{ w.p. } \alpha_x(y) = \min \left\{ 1, \frac{\nu(y)}{\nu(x)} \right\} \\ x \quad \text{w.p. } 1 - \alpha_x(y) \end{cases}$

notes:

1) intuitive: y is "better" ($v(y) \geq v(x)$)

\Rightarrow always move

else stay with some probability

2) Zero-order algorithm: only depends on $v(x)$ (not gradient)

& only up to normalizing constant $\left(\frac{v(y)}{v(x)} \right)$

good for: Bayesian computation

$$p(\theta | \text{data}) \propto p(\theta) \cdot p(\text{data} | \theta)$$

$$\underbrace{p(\text{data})}_{\substack{\text{can't compute} \\ \text{don't need}}} = \int_{\Theta} p(\theta) \cdot p(\text{data} | \theta) d\theta$$

large/exp-size

also in physics, ML, ...

Gibbs distribution

$$v(x) \propto e^{-f(x)}$$

from potential function / energy

$$f: X \rightarrow \mathbb{R}$$

$$\begin{pmatrix} q_i \\ p_i \end{pmatrix} \in \mathbb{R}^2 \times \mathbb{R}^2 = \mathbb{R}^4$$

e.g. in physics: $f = \text{Hamiltonian of } N \text{ particles on } \mathbb{R}^2$

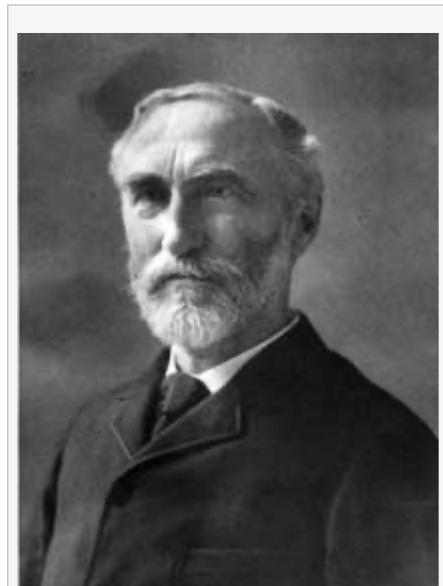
$$X = \mathbb{R}^{4N}$$

• convex geometry / analysis

$$f \text{ convex} \iff v \propto e^{-f} \text{ log-concave}$$

Gibbs

- Josiah Willard Gibbs (1839–1903)
- New Haven, CT
- Yale College:
 - ★ BS 1858,
 - ★ PhD 1863
 - ★ Professor of Mathematical Physics 1871
- Thermodynamics, Gibbs free energy,
statistical mechanics



Josiah Willard Gibbs

“the greatest mind in American history” (Einstein)



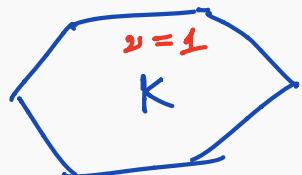
Example: Metropolis Random Walk (MRW)

Start with $P = \text{Brownian motion} / \text{Gaussian noise with step size } \eta$

From x , draw $y = x + \sqrt{2\eta} z$, $z \sim \mathcal{N}(0, I)$ on \mathbb{R}^n

Accept y w.p. $\min \left\{ 1, \frac{\nu(y)}{\nu(x)} \right\}$, else reject

- e.g. $\nu = \text{uniform on } K \subset \mathbb{R}^n$



$$\nu = 1$$

$$\nu(x) \propto e^{-f(x)}$$

$$-\log \nu = f = \mathbb{1}_K : \mathbb{R}^n \rightarrow \mathbb{R}$$

$$\mathbb{1}_K(x) = \begin{cases} 0 & \text{if } x \in K \\ \infty & \text{else} \end{cases}$$



convex if K convex set function

- e.g. $\nu = \text{Gaussian} \Leftrightarrow f \text{ quadratic}$

Generalization: Metropolis-Hastings (MH) Algorithm

To sample from ν :

- Start with any Markov chain P
- Do Random walk with accept/reject:

1. From x , draw $y \sim P_x$

*Metropolis
-Hastings
alg./
filter*

2. Accept y w.p. $\min \left\{ 1, \frac{\nu(y)}{\nu(x)} \frac{P_y(x)}{P_x(y)} \right\} = \alpha_x(y)$
else reject & stay at x

notes: * Hastings (1970)

* Zero-order, only needs ν up to constant

* Universal algorithm : construct from any $P \leftarrow$ has wrong stationary dist μ
MH : force $\mu \rightarrow \nu$

Why does MH work?

Let P be any Markov chain

Let \hat{P} be MH version wrt. ν

Lemma: \hat{P} is reversible wrt. ν ($\Rightarrow \nu$ stationary for \hat{P})

Proof:

$$\hat{P}_x(y) = P_x(y) \cdot \alpha_x(y) + \delta_x(y) \cdot A(x)$$

$$\text{where } A(x) = P_x(\{x\}) + \int_{X \setminus \{x\}} (1 - \alpha_x(y)) \cdot P_x(y) dy$$

$$= 1 - \int_{X \setminus \{x\}} \alpha_x(y) P_x(y) dy$$

$$\text{wts: } \nu(x) \cdot \hat{P}_x(y) = \nu(y) \cdot \hat{P}_y(x) \quad \forall x, y \in X$$

- If $x=y$, true

- If $x \neq y$:
$$\begin{aligned} v(x) \cdot \hat{P}_x(y) &= v(x) \cdot P_x(y) \cdot \alpha_x(y) \\ &= v(x) \cdot P_x(y) \cdot \min \left\{ 1, \frac{v(y) \cdot P_y(x)}{v(x)} \right\} \\ &= \min \left\{ v(x) \cdot P_x(y), v(y) \cdot P_y(x) \right\} \\ &\quad \uparrow \\ &\quad \text{symmetric in } x, y \\ &= v(y) \cdot \hat{P}_y(x). \end{aligned}$$

□

More general family

Given Markov chain $P = (P_x : x \in X)$

Want to make $P \mapsto \hat{P}$ reversible wrt. ν

Accept/reject step:

1. From x , draw $y \sim P_x$
 2. Accept y with some probability $\alpha_{x(y)}$
- } New MC
 \hat{P}

Let $R_x(y) = \frac{\nu(y)}{\nu(x)} \cdot \frac{P_y(x)}{P_x(y)} = \frac{1}{R_y(x)}$

Lemma: Choose any $0 \leq \alpha_{x(y)} \leq \min \{ 1, R_x(y) \}$

[Biller & Diaconis

Prop. 3.1]

set $\alpha_y(x) = \frac{\alpha_{x(y)}}{R_x(y)}$

Then \hat{P} is reversible wrt. ν

Ex: 1) Metropolis (1950) : $\alpha_{x(y)} = \min \{ 1, R_x(y) \}$

2) Barker (1965) :

$$\alpha_{x(y)} = \frac{\nu(y) \cdot p_y(x)}{\nu(x) \cdot p_x(y) + \nu(y) \cdot p_y(x)} = \frac{R_x(y)}{1 + R_x(y)}$$

e.g., P is rev wrt. $\nu \Leftrightarrow R = 1$

\Rightarrow Metropolis: accept w.p. 1

\Rightarrow Barker : accept w.p. $\frac{1}{2}$

Optimization interpretation of MH

[Billera & Diaconis 2001]

Let $\mathcal{M} = \{ \text{Markov chains on } \mathcal{X} \}$

Space of prob. distributions
on \mathcal{X}

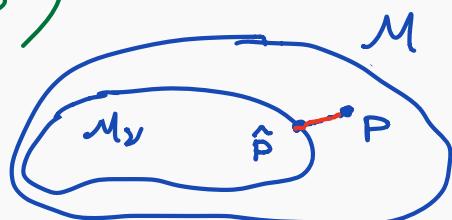
$$= \{ P = (P_x : x \in \mathcal{X}), \quad P_x \in \mathcal{P}(x) \}$$

$$\left(= \mathcal{P}(\mathcal{X})^{\mathcal{X}} \quad = \quad \{ P : \mathcal{X} \rightarrow \mathcal{P}(x) \} \right)$$

Given $\nu \in \mathcal{P}(\mathcal{X})$

Let $\mathcal{M}_\nu = \{ P \in \mathcal{M} : P \text{ reversible wrt. } \nu \}$

then MH filter: $\mathcal{M} \ni P \mapsto \hat{P} \in \mathcal{M}_\nu$



Define metric on \mathcal{M} :

$$d(P, P') = \int_{X \times X} |P_x(y) - P'_x(y)| \nu(x) dx dy$$

Δ

$\Delta = \{(x, x) : x \in \mathcal{X}\}$

Theorem: [Biller & Diaconis Thm 1]

MH filter is projection in d distance

$$\hat{P} = \arg \min_{\tilde{P} \in \mathcal{M}_\nu} d(\tilde{P}, P)$$

after MH

starting MC

(and \hat{P} is unique such MC with $\hat{P}_x(y) \leq P_x(y)$)

- Why metric d ?

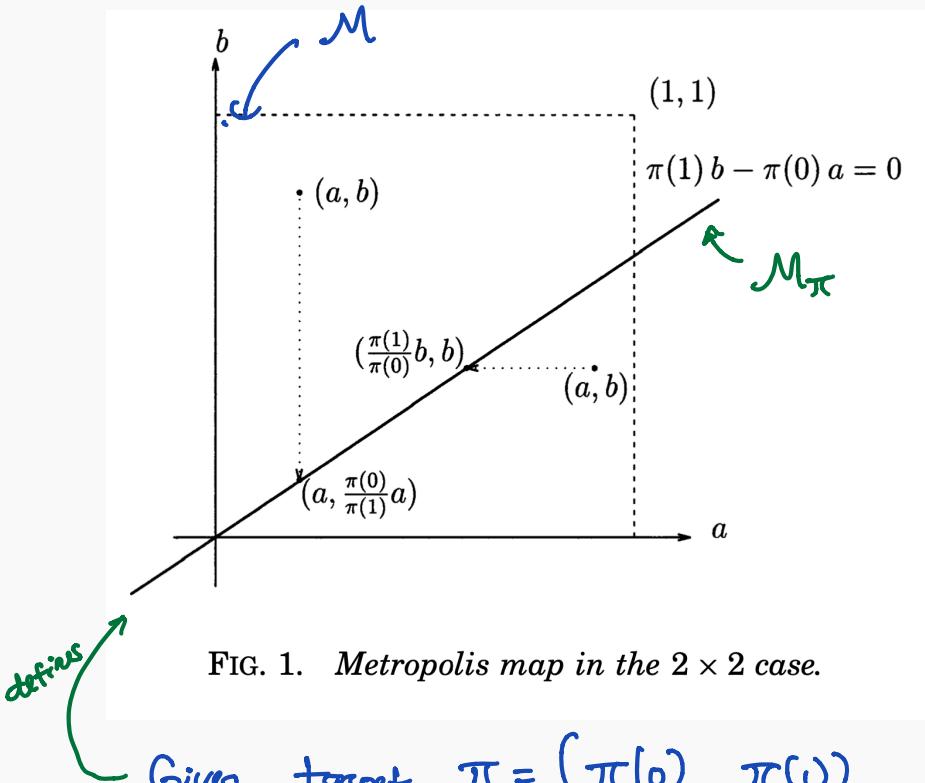
Given P , given ν , define:

$$\begin{aligned} (\nu P)(x, y) &= \nu(x) \cdot P_x(y) \\ (\overline{\nu P})(x, y) &= \nu(y) \cdot P_y(x) \end{aligned} \quad \left. \begin{array}{l} \in \mathcal{P}(X \times X) \\ P \text{ rev. wrt. } \nu \Leftrightarrow \nu P = \overline{\nu P} \end{array} \right\}$$

then

$$d(P, \mathcal{M}_\nu) = TV(\nu P, \overline{\nu P})$$

Example

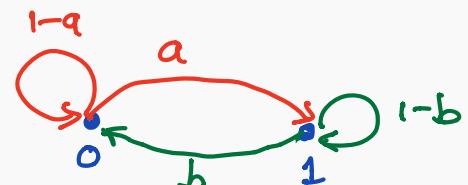


Given target $\pi = (\pi(0), \pi(1))$

P stationarity wrt. $\pi \Leftrightarrow \pi(0) \cdot a = \pi(1) \cdot b$

2-point space

$$\mathcal{X} = \{0, 1\}$$



$$P = \begin{pmatrix} 1-a & a \\ b & 1-b \end{pmatrix}$$

↑
stationary: $\left(\frac{b}{a+b}, \frac{a}{a+b} \right)$

Extensions

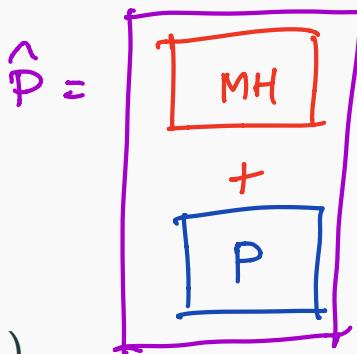
Projects:

Q: What if projection in other distances? KL divergence, χ^2 -div, Φ -div

- • Diaconis and Miclo, *On characterizations of Metropolis-type algorithms in continuous time*, Latin American Journal of Probability and Mathematical Statistics, 2009
- Choi, *Metropolis-Hastings reversibilizations of non-reversible Markov chains*, Stochastic Processes and their Applications, 2020
 $\min \rightarrow \max$
- Chewi, Lu, Ahn, Cheng, Gouic, Rigollet, *Optimal dimension dependence of the Metropolis-Adjusted Langevin Algorithm*, arXiv:2012.12810, 2020

Basic algorithms for sampling

1. Metropolis Random Walk (**MRW**)
2. Unadjusted Langevin Algorithm (**ULA**)
3. Metropolis-Adjusted Langevin Algorithm (**MALA**)



(also underdamped Langevin algorithm

Hamiltonian Monte Carlo

...

)

1. Metropolis Random Walk (MRW)

$P = \text{Brownian motion}$ $\downarrow z_k \sim N(0, I)$

$$x_{k+1} = \begin{cases} y_k = x_k + \sqrt{2\eta} z_k & \text{w.p. } \min \left\{ 1, \frac{\nu(y_k)}{\nu(x_k)} \right\} \\ x_k & \text{else} \end{cases}$$



2. Unadjusted Langevin Algorithm (ULA)

To sample from $\nu \propto e^{-f}$:

$$x_{n+1} = x_n - \eta \nabla f(x_n) + \sqrt{2\eta} z_n, \quad z_n \sim \mathcal{N}(0, I)$$

- why? From continuous time (later)
- has asymptotic bias
(stationary $\neq \nu$)

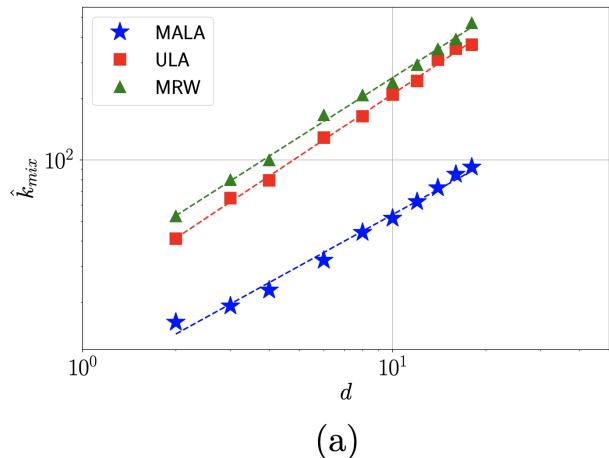
3. Metropolis-Adjusted Langevin Algorithm (MALA)

choose $P = \text{ULA}$
+ MH filter
 +
= MALA

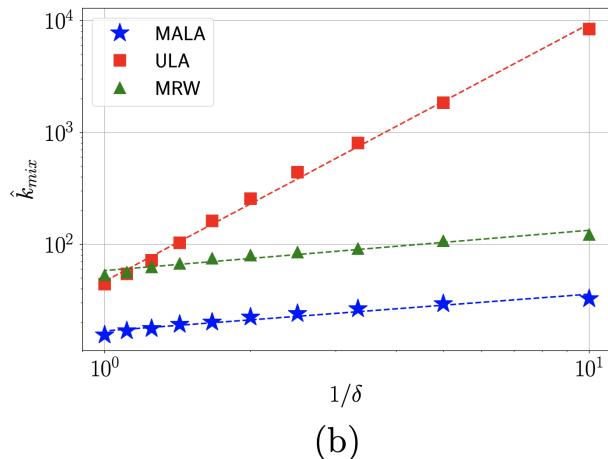
Comparison: Gaussian

$$\nu = \mathcal{N}(0, \Sigma)$$

DWIVEDI, CHEN, WAINWRIGHT AND YU



(a)

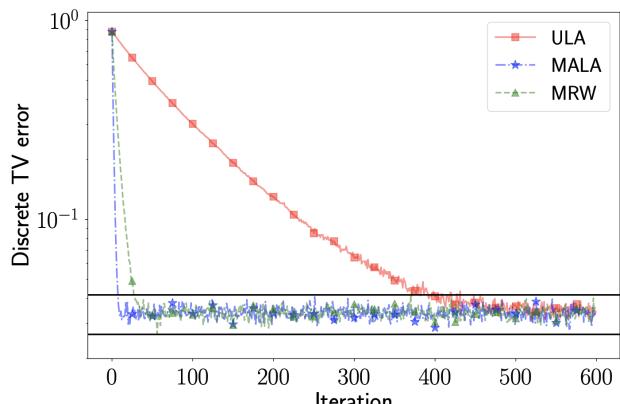


(b)

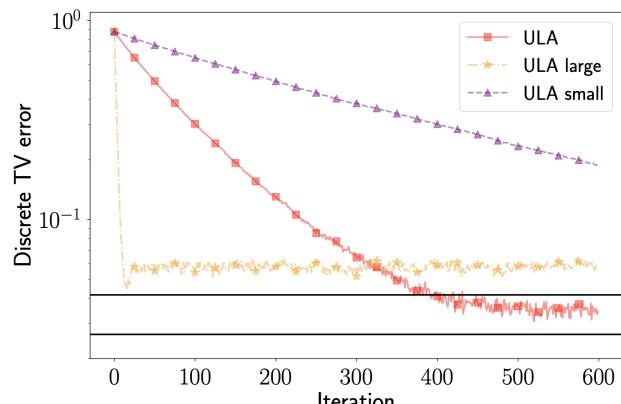
Figure 1. Scaling of the approximate mixing time \hat{k}_{mix} (refer to the discussion after equation (19) for the definition) on multivariate Gaussian density (19) where the covariance has condition number $\kappa = 4$. (a) Dimension dependency. (b) Error-tolerance dependency.

$$\kappa = \frac{\lambda_{\max}(\Sigma)}{\lambda_{\min}(\Sigma)}$$

Comparison: Mixture of Gaussians



(a)



(b)

Figure 4. Discrete TV error on a two component Gaussian mixture. (a) Behavior of three different random walks. (b) Behavior of ULA with different choices of step sizes.

Comparison: Mixing time

Random walk	Distribution μ_*	$t_{\text{mix}}(\delta; \mu_*)$
ULA (Cheng and Bartlett, 2018)	$\mathcal{N}(x^*, m^{-1}\mathbb{I}_d)$	$\mathcal{O}\left(\frac{d\kappa^2 \log(d\kappa/\delta)}{\delta^2}\right)$
ULA (Dalalyan, 2016)	$\mathcal{N}(x^*, L^{-1}\mathbb{I}_d)$	$\mathcal{O}\left(\frac{(d^3 + d \log^2(1/\delta))\kappa^2}{\delta^2}\right)$
MRW (this work)	$\mathcal{N}(x^*, L^{-1}\mathbb{I}_d)$	$\mathcal{O}\left(d^2\kappa^2 \log^{1.5}\left(\frac{\kappa}{\delta}\right)\right)$
MALA (this work)	$\mathcal{N}(x^*, L^{-1}\mathbb{I}_d)$	$\mathcal{O}\left(d^2\kappa \log\left(\frac{\kappa}{\delta}\right)\right)$

Table 2. Scalings of upper bounds on δ -mixing time, from the starting distribution μ_* given in column two, for different random walks in \mathbb{R}^d with target $\pi \propto e^{-f}$ such that $m\mathbb{I}_d \preceq \nabla^2 f(x) \preceq L\mathbb{I}_d$ for any $x \in \mathbb{R}^d$ and $\kappa := L/m$. Here x^* denotes the unique mode of the target density π .

Recap

Metropolis-Hastings: Universal filter to construct reversible MC

Next time: How to bound s -conductance?

Please fill out survey: <https://forms.gle/mCTWsWPdsdoTEkvR9>