

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

---

Andre Wibisono

February 24, 2021

Yale University

## Last time

- Reversible Markov chain  $v(x) \cdot P_x(y) = v(y) \cdot P_y(x)$

- Spectral gap  $\Leftrightarrow$  Conductance  $\gamma = \lambda_2(L)$   $\phi$  *Cheeger:*  $\frac{\phi^2}{2} \leq \gamma \leq 2\phi$
- $s$ -Conductance  $\Rightarrow$  Mixing time in TV distance

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

$M = M_{\infty}(S_0)$  warmth of  $S_0$

# Last time

Questions:

1. How to construct reversible Markov chain?

⇒ Metropolis-Hastings : accept/reject step

$$P + MH_{\nu} = \tilde{P} \quad \text{reversible wrt } \nu$$

from  $x$ , draw  $y \sim P_x$

$$\text{set } x' = \begin{cases} y & \text{w.p. } \min \left\{ 1, \frac{\nu(y) \cdot P_y(x)}{\nu(x) \cdot P_x(y)} \right\} \\ x & \text{w.p. } 1 - \alpha_x(y) \end{cases} = \alpha_x(y)$$

e.g., •  $P$  = Brownian motion  $\Rightarrow \tilde{P}$  = Metropolis Random Walk (MRW)

•  $P$  = Unadjusted Langevin Alg. (ULA)  $\Rightarrow \tilde{P}$  = Metropolis-Adjusted Langevin Alg. (MALA)

2. How to bound conductance?

Langevin Alg. (MALA)

⇒ Today: Isoperimetry

# References

- Vempala, *Geometric Random Walk: A Survey*, Combinatorial and Computational Geometry, 2005
- Dwivedi, Chen, Wainwright, and Yu, *Log-Concave Sampling: Metropolis-Hastings Algorithms are Fast*, Journal of Machine Learning Research, 2019

# From conductance to isoperimetry

Let  $P$  be Markov chain reversible w.r.t  $\nu$

Suppose  $P$  has large conductance

$$\phi = \inf_{A \subset \mathcal{X}} \frac{\Phi(A)}{\min\{\nu(A), 1 - \nu(A)\}}$$

where  $\Phi(A)$  is ergodic flow

$$\sim \Pr(X_0 \in A, X_t \notin A), \quad \begin{matrix} X_0 \sim \nu \\ X_t | X_0 \sim P_{X_0} \end{matrix}$$

What does this mean for  $\nu$ ?

\* In general, nothing: there are always M.C. with large conductance

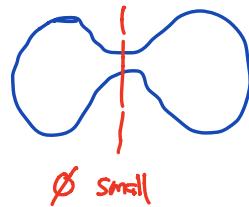
e.g. the ideal MC :  $P_x = \nu$

(random walk: from  $x$ , jump to  $y \sim \nu$  independent)

Converges in one step (but not implementable)

Exercise:  $\phi = \frac{1}{2}$  (for any  $\nu$ )

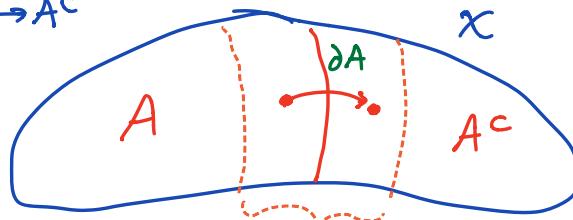
But there are bad targets, e.g.  $\mathcal{V}$  = uniform on



\* Suppose  $P$  is local:  $x$  and  $y \sim P_x$  are close in some distance  $d$

$$d(x, y) \ll 1 \quad (\text{w.h.p.})$$

then points that cross over  $A \rightarrow A^c$   
are near boundary  $\partial A$



\* Large conductance

$\Rightarrow$  region near boundary has large volume  
(large sets have large boundaries)

### isoperimetry

\* Conversely, will show isoperimetry + one-step overlap

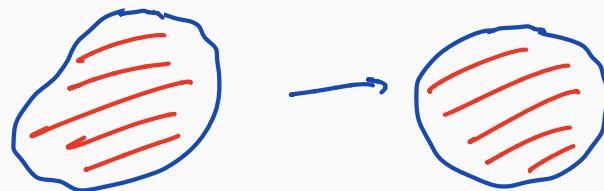
$\Rightarrow$  large conductance

# Isoperimetry

Problem: Among all regions of equal perimeter / area, which has largest volume?

Dido [c. 800 BC]

Queen of Carthage

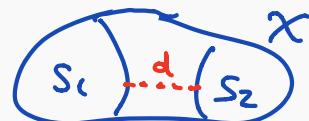


Formally:

Let  $X$  have a distance metric  $d(x, y)$

(e.g.  $X = \mathbb{R}^n$ ,  $d(x, y) = \|x - y\|_2$  Euclidean distance)

Given 2 sets  $S_1, S_2 \subset X$

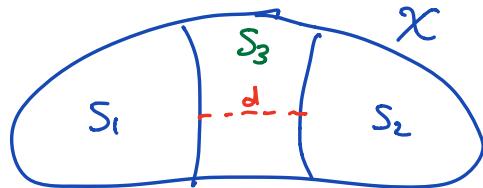


define distance  $d(S_1, S_2) = \inf \{ d(x, y) : x \in S_1, y \in S_2 \}$

### Definition:

Given a prob dist  $\nu$  on  $X$ , we say  $\nu$  is isoperimetric (or satisfies isoperimetric inequality) with constant  $\Psi > 0$  if  $\forall$  partition  $X = S_1 \cup S_2 \cup S_3$ :

$$\nu(S_3) \geq \Psi \cdot d(S_1, S_2) \cdot \min\{\nu(S_1), \nu(S_2)\}$$



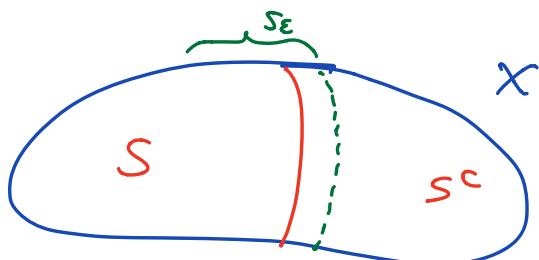
### Notes:

1. Equivalent to (differential formulation):

$$\forall S \subset X, \quad \nu(\partial S) \geq \Psi \cdot \min\{\nu(S), \nu(S^c)\}$$

where  $\nu(\partial S) = \lim_{\varepsilon \rightarrow 0} \frac{\nu(S_\varepsilon) - \nu(S)}{\varepsilon}$

$$S_\varepsilon = \{x \in X : d(x, S) \leq \varepsilon\}$$



2. Equivalent to Poincaré inequality (functional form):

$\forall g : \mathcal{X} \rightarrow \mathbb{R}$

$$\mathbb{E}_\nu [\|\nabla g\|^2] \geq \frac{\psi^2}{4} \text{Var}_\nu(g)$$

(this is the spectral gap of Dirichlet form of  $\nu$ )

Equivalently,

$$\psi = \Psi_\nu = \inf_{S \subset X} \frac{\nu(\partial S)}{\min\{\nu(S), \nu(S^c)\}}$$

# Isoperimetric Theorems

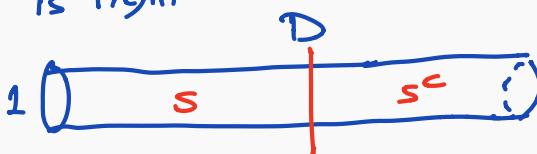
1. Theorem [Dyer & Frieze '91]:

$\nu = \text{uniform on convex } K \subset \mathbb{R}^n$  with diameter  $D$ .

$$\Rightarrow \Psi \geq \frac{2}{D}$$

Note: dependence on  $D$  is tight

e.g.,  $K = \text{cylinder}$



$$\Rightarrow \Psi \sim \frac{1}{D}$$

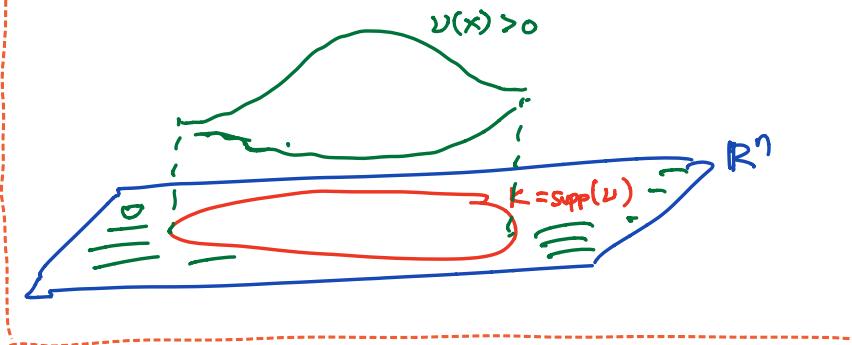
2. Can be generalized to:

$\nu = \text{log-concave distribution on } \mathbb{R}^n$  with support of diameter  $D$

$$v(x) \propto e^{-f(x)}, \quad f: \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$$

$$\text{support}(v) = \{x \in \mathbb{R}^n : v(x) > 0\}$$

$$= \{x \in \mathbb{R}^n : f(x) < \infty\} = \text{dom}(f)$$



$$\Rightarrow \Psi \geq \frac{2}{D}$$

3. Theorem [KLS '95]

$v$  = log-concave distribution on  $\mathbb{R}^n$

$$\Rightarrow \Psi \geq \frac{\log 2}{M_1(v)} \stackrel{\log 2 = 0.3 > \frac{1}{4}}{\geq} \frac{1}{4 M_1(v)} \geq \frac{1}{4 \sqrt{\text{Var}_v(x)}}$$

where  $\mu = \mathbb{E}_v[X]$

$$M_1(v) = \mathbb{E}_v[\|X - \mu\|_2]$$

$$(\text{note: } M_1(v) \leq \mathbb{E}_v[\|X - \mu\|_2^2]^{1/2} = \sqrt{\text{Var}_v(x)})$$

e.g. for  $v$  isotropic ( $\mu = \mathbb{E}_v[X] = 0$   
 $\Sigma = \text{Cov}_v(x) = \mathbb{E}_v[(x - \mu)(x - \mu)^T] = I$ )

$$\Rightarrow \text{Var}_v(x) = n \Rightarrow \text{for isotropic log-concave } v, \quad \Psi \geq \Omega\left(\frac{1}{\sqrt{n}}\right).$$

# KLS Conjecture

Let  $\nu$  be log-concave on  $\mathbb{R}^n$

Let  $\Sigma = \text{Cov}_\nu(X) \in \mathbb{R}^{n \times n}$  with eigenvalues  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n > 0$ .

Theorem [KLS '95] :  $\Psi = \Omega\left(\frac{1}{\sqrt{\sum_{i=1}^n \lambda_i}}\right) = \Omega\left(\frac{1}{\sqrt{\text{Tr}(\Sigma)}}\right)$

$$\Sigma = I$$

$$\frac{1}{\sqrt{n}}$$

Conjecture [KLS '95]:  $\Psi \stackrel{?}{=} \Omega\left(\frac{1}{\sqrt{\lambda_1}}\right) = \Omega\left(\frac{1}{\|\Sigma\|_{\text{op}}^{1/2}}\right)$

$$1 ?$$

Theorem [LV '17] :  $\Psi = \Omega\left(\frac{1}{\left(\sum_{i=1}^n \lambda_i^2\right)^{1/4}}\right) = \Omega\left(\frac{1}{\|\Sigma\|_{\text{HS}}^{1/2}}\right)$

$$\frac{1}{n^{1/4}}$$

Theorem [C '20] :  $\Psi = \Omega\left(\frac{1}{n^{o(1)} \cdot \sqrt{\lambda_1}}\right)$

$$\frac{1}{n^{o(1)}}$$

little-o

# References for KLS Conjecture

- Kannan, Lovász, & Simonovits, *Isoperimetric problems for convex bodies and a localization lemma*, Discrete & Computational Geometry, 1995
- Lee & Vempala, *The Kannan-Lovász-Simonovits Conjecture*, arXiv:1807.03465, 2018
- Lee & Vempala, *Eldan's Stochastic Localization and the KLS Conjecture: Isoperimetry, Concentration and Mixing*, arXiv:1612.01507v3, 2019
- Chen, *An Almost Constant Lower Bound of the Isoperimetric Coefficient in the KLS Conjecture*, arXiv:2011.13661, 2020

# Strongly Log-Concave Distribution

Def:  $v$  is  $\alpha$ -strongly log-concave ( $\alpha$ -SLC) for some  $\alpha > 0$

If  $v(x) \propto e^{-f(x)}$ ,  $f: \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$

such that  $f$  is  $\alpha$ -strongly convex



(Recall: this means:

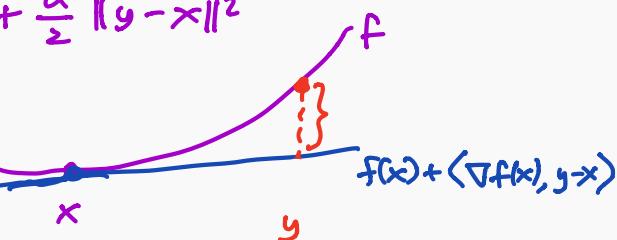
$$\frac{f(x) + f(y)}{2} - f\left(\frac{x+y}{2}\right) \geq \frac{\alpha}{8} \|y-x\|^2$$

- if  $f$  is differentiable:

$$\Leftrightarrow f(y) \geq f(x) + \langle \nabla f(x), y-x \rangle + \frac{\alpha}{2} \|y-x\|^2$$

$\nabla f(x)$  = gradient:

$$\left( \frac{\partial f}{\partial x_1}(x), \dots, \frac{\partial f}{\partial x_n}(x) \right)$$



- if  $f$  is twice-differentiable:

$$\Leftrightarrow \nabla^2 f(x) \succeq \alpha I$$

$\nabla^2 f(x)$  = Hessian matrix

$$\left( \frac{\partial^2 f}{\partial x_i \partial x_j}(x) \right)_{i,j=1}^n$$

- log-concave is when  $\alpha = 0$
- eg.  $v = \text{Gaussian } \mathcal{N}(0, \Sigma)$ ,  $\Sigma$  has eigenvalues  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n > 0$

then  $v$  is  $\alpha$ -SLC where  $\alpha = \frac{1}{\lambda_1}$

because

$$v(x) = \frac{1}{\sqrt{\det(2\pi\Sigma)}} e^{-\frac{x^\top \Sigma^{-1} x}{2}}$$

$$\text{so } f(x) = -\log v(x) = \frac{1}{2} x^\top \Sigma^{-1} x + \frac{1}{2} \log \det(2\pi\Sigma)$$

$$\nabla f(x) = \Sigma^{-1} x$$

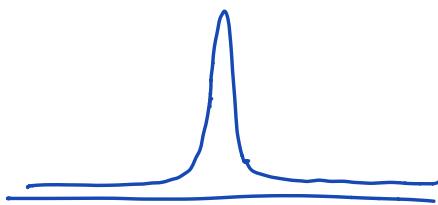
$$\nabla^2 f(x) = \Sigma^{-1} \geq \frac{1}{\lambda_1} I$$

$\uparrow$  has eigenvalues  $\frac{1}{\lambda_1} \leq \frac{1}{\lambda_2} \leq \dots \leq \frac{1}{\lambda_n}$

$$\text{eg. } \Sigma = \lambda I$$

- if  $\lambda$  is small:

$v$  very peaked



$\Rightarrow$  very SLC ( $\frac{1}{\lambda}$  large)

- if  $\lambda$  is large:

$v$  very flat



$\Rightarrow$  not very SLC ( $\frac{1}{\lambda}$  small)

Theorem: If  $\nu$  is  $\alpha$ -SLC on  $\mathbb{R}^n$

$$\text{then } \Psi \geq (\log 2) \sqrt{\alpha} \geq \frac{\sqrt{\alpha}}{4}$$

- [Cousins - Vempala '16]
- Nice: dimension-free (doesn't depend on  $n$ )
- (Naively, using earlier bound

$$\Psi \geq \frac{\log 2}{M_1(\nu)} \geq \frac{\log 2}{\sqrt{Var_{\nu}(x)}} \geq \log 2 \cdot \sqrt{\frac{\alpha}{n}} \quad )$$

↑  
bad