

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

---

Andre Wibisono

March 1, 2021

Yale University

## Last time

- Reversible Markov chain  $\nu(x) \cdot P_x(y) = \nu(y) \cdot P_y(x)$
- $s$ -Conductance  $\Rightarrow$  Mixing time in TV distance

$$\phi_s \quad T(\varepsilon) = \frac{2}{\phi_s^2} \log \frac{2M}{\varepsilon}, \quad s = \frac{\varepsilon}{2M}, \quad M = M_\nu(s_0)$$

- Metropolis-Hastings: accept - reject step

from  $x$ , go to  $x' = \begin{cases} y \sim P_x & \text{with prob } \min \left\{ 1, \frac{\nu(y) \cdot P_y(x)}{\nu(x) \cdot P_x(y)} \right\} = a_x(y) \\ x & \text{with prob } 1 - a_x(y) \end{cases}$

- Isoperimetry of  $\nu$

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

## Today:

- From isoperimetry to conductance
- Metropolis Random Walk (MRW)

# 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

# Strongly Log-Concave Distribution

Recall  $v \propto e^{-f}$  on  $\mathbb{R}^n$  is  $\alpha$ -strongly log-concave ( $\alpha$ -SLC)

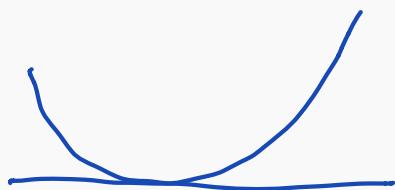
if  $f = -\log v : \mathbb{R}^n \rightarrow \mathbb{R}$  is  $\alpha$ -strongly convex

$$\Leftrightarrow \boxed{\nabla^2 f(x) \geq \alpha I} \quad \forall x \in \mathbb{R}^n$$

(Recall this means all eigenvalues  $\lambda$  of  $\nabla^2 f(x)$  are lower bounded:

$$\lambda \geq \alpha$$

)



• example:  $v = \mathcal{N}(\mu, \Sigma)$ ,  $\mu \in \mathbb{R}^n$ ,  $\Sigma > 0$

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

$$f(x) = \frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu) + \frac{1}{2} \log \det(2\pi\Sigma)$$

$$\nabla^2 f(x) = \Sigma^{-1} \geq \frac{1}{\lambda_{\max}(\Sigma)} I \leftarrow \alpha = \frac{1}{\lambda_{\max}(\Sigma)}$$

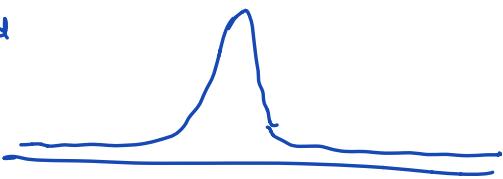
large  $\Sigma \Rightarrow \nu$  is flat

$\Rightarrow$  less SLC



Small  $\Sigma \Rightarrow \nu$  is peaked

$\Rightarrow$  more SLC



# 1. Isoperimetry

## Theorem

If  $\nu$  is  $\alpha$ -strongly log-concave on  $\mathbb{R}^n$ , then it satisfies isoperimetry with constant:

$$\psi \geq (\log 2) \sqrt{\alpha} \geq \frac{\sqrt{\alpha}}{4}$$

- Cousins & Vempala, Gaussian Cooling and  $O^*(n^3)$  Algorithms for Volume and Gaussian Volume, SIAM Journal of Computing, 2018 (Theorem 5.4)
- Proof by localization: Reduce problem from  $n$  dimension  $\rightarrow$  1 dimension
- Nice: dimension free  
(cf. KLS conjecture for weakly log-concave )

## 2. Effective Diameter

Prob distribution  $v \propto e^{-f}$  on  $\mathbb{R}^n$



- Mode:  $x^* = \arg \max_{x \in \mathbb{R}^n} v(x)$

$$= \arg \min_{x \in \mathbb{R}^n} f(x)$$

[ Convex optimization: if  $f$  is  $\alpha$ -strongly convex,  $L$ -smooth  
 $\alpha I \leq \nabla^2 f(x) \leq L I$  ]

Define condition number

$$\kappa = \frac{L}{\alpha} \geq 1$$

- Gradient descent can find  $x^*$  in time  $\tilde{O}(K)$
- Accelerated gradient descent can find  $x^*$  in time  $\tilde{O}(\sqrt{\kappa})$   
(this is optimal)

]

$$\bullet \text{ mean: } \bar{x} = \mathbb{E}_\nu[x] = \frac{\int_{\mathbb{R}^n} x e^{-f(x)} dx}{\int_{\mathbb{R}^n} e^{-f(x)} dx}$$

Lemmas: Suppose  $\nu$  is  $\alpha$ -SLC on  $\mathbb{R}^n$ . Then:

$$1) \text{Var}_\nu(x) = \mathbb{E}_\nu[\|x - \bar{x}\|_2^2] \leq \frac{n}{\alpha}$$

$$2) \mathbb{E}_\nu[\|x - x^*\|_2^2] \leq \frac{n}{\alpha}$$

$$3) \|\bar{x} - x^*\|_2^2 \leq \frac{n}{\alpha}$$

}  $\therefore$  effective diameter  
 $\sim \sqrt{\frac{n}{\alpha}}$

example:  $\nu = \mathcal{N}(\mu, \Sigma)$ ,  $\alpha$ -SLC with  $\alpha = \frac{1}{\lambda_{\max}(\Sigma)}$

$$\text{Var}_\nu(x) = \text{Tr}(\Sigma) = \sum_{i=1}^n \lambda_i(\Sigma) \leq n \cdot \lambda_{\max}(\Sigma) = \frac{n}{\alpha}$$

### 3. Concentration

For  $s > 0$ , let  $r(s) = 2 + 2 \max \left\{ \left( \frac{1}{n} \log \frac{1}{s} \right)^{\frac{1}{4}}, \left( \frac{1}{n} \log \frac{1}{s} \right)^{\frac{1}{2}} \right\}$

(note: • if  $s \geq e^{-n}$  then  $\frac{1}{n} \log \frac{1}{s} \leq 1 \Rightarrow r(s) \leq 4$

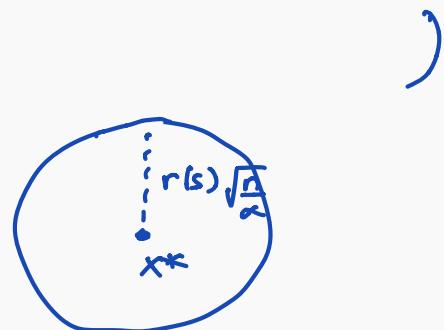
• if  $s < e^{-n}$  then  $r(s) \sim \left( \frac{1}{n} \log \frac{1}{s} \right)^{\frac{1}{2}}$

• if  $s = e^{-cn}$  then  $r(s) \sim \sqrt{c}$

Let  $R_s = \mathbb{B}(x^*, r(s)\sqrt{\frac{n}{\alpha}})$

Lemma: IF  $\nu$  is  $\alpha$ -SLC on  $\mathbb{R}^n$

then  $\nu(R_s) \geq 1 - s$



[DCWY '19, Lemma 1]

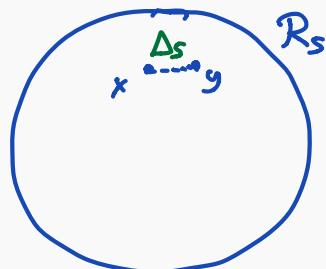
# From isoperimetry to conductance

## Lemma (DCWY'19 Lemma 2)

Suppose  $\nu$  is  $\alpha$ -strongly log-concave on  $\mathbb{R}^n$ . Fix  $0 < s \leq \frac{1}{2}$ . Let  $P$  be a Markov chain with stationary distribution  $\nu$  and satisfying the following one-step overlap property:

$$TV(P_x, P_y) \leq \frac{3}{4} \quad (\text{or } \leq 1 - c \text{ for some } c > 0)$$

for all  $x, y \in \mathcal{R}_s$  with  $\|x - y\|_2 \leq \Delta_s$ . Then  $P$  has  $s$ -conductance:



$$\phi_s \geq \min \left\{ \frac{1}{16}, \frac{\sqrt{\alpha} \Delta_s}{128} \right\} = \Omega(\sqrt{\alpha} \cdot \Delta_s)$$

Corollary: Mixing time in TV:  $\mathcal{T}(\varepsilon) = \tilde{O}\left(\frac{1}{\phi_s^2}\right) = \tilde{O}\left(\frac{1}{\alpha \cdot \Delta_s^2}\right)$

Proof: Let  $A \subset \mathbb{R}^n$  with  $s < \nu(A) < 1-s$ .

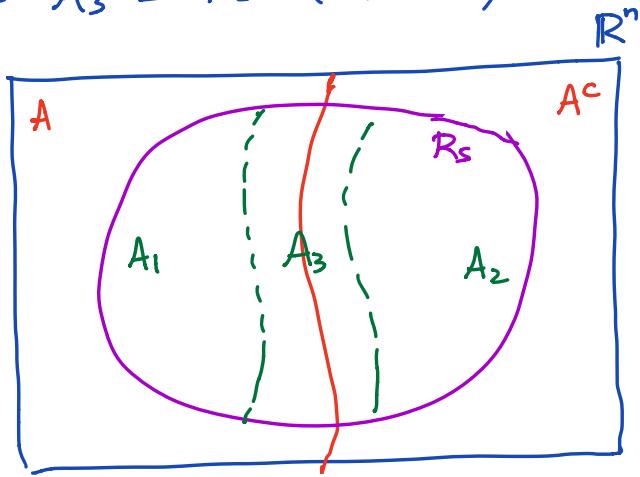
We will bound  $C(A) = \frac{\Phi(A)}{\min\{\nu(A)-s, \nu(A^c)-s\}}$

① Consider points which are "deep inside" and unlikely to cross over

$$A_1 = \{x \in A \cap R_s \mid P_x(A^c) < \frac{1}{8}\}$$

$$A_2 = \{y \in A^c \cap R_s \mid P_y(A) < \frac{1}{8}\}$$

and define  $A_3 = R_s \setminus (A_1 \cup A_2)$



Claim:  $d(A_1, A_2) \geq \Delta_s$  ... (1)

Proof of claim:  $\forall x \in A_1, y \in A_2:$

$$TV(P_x, P_y) = \sup_{B \subset \mathbb{R}^n} |P_x(B) - P_y(B)|$$

$$\geq P_x(A) - P_y(A)$$

$$= 1 - P_x(A^c) - P_y(A)$$

$$\geq 1 - \frac{1}{8} - \frac{1}{8} = \frac{3}{4}$$

So by (the contrapositive of) the one-step overlap:

$$\|x-y\|_2 \geq \Delta_s$$

Then:  $d(A_1, A_2) = \inf_{\substack{x \in A_1 \\ y \in A_2}} \|x-y\|_2 \geq \Delta_s.$

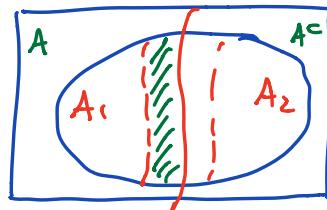
□

② Claim: if  $\nu(A_1) \leq \frac{\nu(A \cap R_s)}{2}$ , then  $C(A) \geq \frac{1}{16}$

Proof:

$$\bar{\Phi}(A) = \int_A P_x(A^c) d\nu(x)$$

$$\geq \int_{(A \cap R_s) \setminus A_1} P_x(A^c) d\nu(x) \quad \underbrace{\geq \frac{1}{8}}$$



$$\geq \frac{1}{8} \nu((A \cap R_s) \setminus A_1)$$

$$\geq \frac{1}{16} \cdot \nu(A \cap R_s) \quad \text{since } \nu(A_1) \leq \frac{\nu(A \cap R_s)}{2}$$

$$= \frac{1}{16} (\nu(A) - \nu(A \cap R_s^c))$$

$$\geq \frac{1}{16} (\nu(A) - \nu(R_s^c))$$

$$\geq \frac{1}{16} (\nu(A) - s) \quad \text{since } \nu(R_s) \geq 1-s$$

$$\geq \frac{1}{16} \min \{ \nu(A) - s, \nu(A^c) - s \}$$

$$\Rightarrow C(A) \geq \frac{1}{16}.$$

□

Similarly, if  $\nu(A_2) \leq \frac{\nu(A^c \cap R_s)}{2}$ , then  $C(A) \geq \frac{1}{16}$ .

③ Now assume:

$$\nu(A_1) \geq \frac{\nu(A \cap R_s)}{2} \geq \frac{\nu(A) - s}{2}$$

$$\nu(A_2) \geq \frac{\nu(A^c \cap R_s)}{2} \geq \frac{\nu(A^c) - s}{2}$$

... (2)

Now compute:

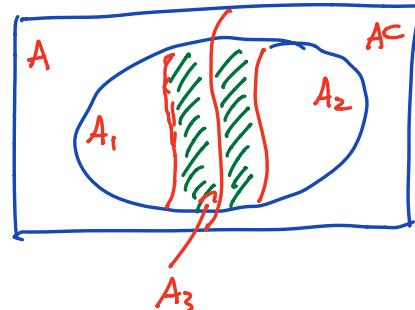
$$\overline{\Phi}(A) = \frac{1}{2} (\overline{\Phi}(A) + \overline{\Phi}(A^c))$$

$$= \frac{1}{2} \left( \int_A P_x(A^c) d\nu(x) + \int_{A^c} P_x(A) d\nu(x) \right)$$

$$\geq \frac{1}{2} \left( \int_{(A \cap R_s) \setminus A_1} P_x(A^c) d\nu(x) + \int_{(A^c \cap R_s) \setminus A_2} P_x(A) d\nu(x) \right)$$

$\geq \frac{1}{8}$

$$\geq \frac{1}{16} \nu(A_3) \dots (3)$$



④ Let  $\nu_s = \nu|_{R_s}$  be restriction of  $\nu$  on  $R_s$

$$\nu_s(B) = \frac{\nu(B \cap R_s)}{\nu(R_s)}$$

Nice:  $\nu$  is  $\alpha$ -SLC      }  
 $R_s$  is convex set      }  $\nu_s$  is also  $\alpha$ -SLC

By isoperimetry for  $\nu_s$ :

$$\text{for partition } R_s = A_1 \cup A_2 \cup A_3$$

$$\nu_s(A_3) \geq \frac{\sqrt{\alpha}}{4} \cdot d(A_1, A_2) \cdot \min \{ \nu_s(A_1), \nu_s(A_2) \}$$

$$\Leftrightarrow \nu(A_3) \geq \frac{\sqrt{\alpha}}{4} \cdot d(A_1, A_2) \cdot \min \{ \nu(A_1), \nu(A_2) \}$$

$$\stackrel{(1)}{\geq} \frac{\sqrt{\alpha}}{4} \cdot \Delta_s \cdot \min \{ \nu(A_1), \nu(A_2) \}$$

$$\stackrel{(2)}{\geq} \frac{\sqrt{\alpha} \cdot \Delta_s}{8} \cdot \min \{ \nu(A) - s, \nu(A^c) - s \}$$

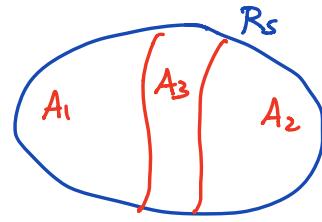
Plugging in to (3):

$$\overline{\Phi}(A) \geq \frac{1}{16} \cdot \nu(A_3)$$

$$\geq \frac{\sqrt{\alpha} \cdot \Delta_s}{128} \cdot \min \{ \nu(A) - s, \nu(A^c) - s \}$$

$$\Rightarrow C(A) \geq \frac{\sqrt{\alpha} \cdot \Delta_s}{128}.$$

□



## Recap

To sample from  $\nu$  on  $\mathbb{R}^n$ :

1. Start from any Markov chain  $P$
2. Apply Metropolis-Hastings filter to get  $\tilde{P}$  reversible wrt  $\nu$
3. Assume  $\nu$  is  $\alpha$ -SLC, so isoperimetric with  $\psi = \Omega(\sqrt{\alpha})$
4. Show  $\tilde{P}$  satisfies one-step overlap property:

$$x, y \in \mathcal{R}_s, \|x - y\|_2 \leq \Delta_s \Rightarrow \text{TV}(\tilde{P}_x, \tilde{P}_y) \leq \frac{3}{4}$$

$\Rightarrow \tilde{P}$  has  $s$ -conductance  $\phi_s = \Omega(\sqrt{\alpha} \Delta_s)$

$\Rightarrow$  mixing time in TV distance:  $\tau(\epsilon) = \tilde{O}\left(\frac{1}{\phi_s^2}\right) = \tilde{O}\left(\frac{1}{\alpha \Delta_s^2}\right)$

# What random walk?

1.  $P$  = Brownian motion (Gaussian walk)  
 $\Rightarrow \tilde{P}$  = Metropolis Random Walk (MRW)  
(Today)
  
2.  $P$  = Unadjusted Langevin Algorithm (ULA)  
 $\Rightarrow \tilde{P}$  = Metropolis-Adjusted Langevin Algorithm (MALA)  
(Next time)

# Kullback-Leibler (KL) Divergence

Let  $\rho$  and  $\nu$  be probability distributions on  $\mathcal{X}$ .

The **Kullback-Leibler (KL) divergence** of  $\rho$  with respect to  $\nu$  is

$$H_\nu(\rho) = \int_{\mathcal{X}} \rho(x) \log \frac{\rho(x)}{\nu(x)} dx$$

Notes: • Also called relative entropy

• Relative form of Shannon entropy  $H(s) = - \int_{\mathcal{X}} s(x) \log s(x) dx$

• Non-negative:  $H_\nu(\rho) \geq 0 \quad \forall \rho, \nu$

$$H_\nu(\rho) = 0 \iff \rho = \nu$$

• Not a metric: not symmetric  
does not satisfy triangle inequality

- If  $s$  has density  $h = \frac{ds}{d\nu} : X \rightarrow \mathbb{R}$  wrt  $\nu$

then  $H_\nu(s) = \mathbb{E}_\nu [h(x) \log h(x)]$

- Recall:  $TV(s, \nu) = \frac{1}{2} \mathbb{E}_\nu [|h(x) - 1|]$

- Pinsker's Inequality:  $TV(s, \nu) \leq \sqrt{\frac{1}{2} H_\nu(s)}$