

CPSC 661: Sampling Algorithms in ML

Andre Wibisono

February 10, 2021

Last time

- Reversible Markov chain P wrt ν
 $\nu(x) \cdot P(x, y) = \nu(y) \cdot P(y, x)$
- Laplacian $L = I - P \geq 0$
 $L\mathbf{1} = 0 \cdot \mathbf{1} \leftarrow \lambda_1(L) = 0$
- Spectral gap $\gamma = \lambda_2(L) = \inf_{f \in L^2(\nu)} \frac{\langle f, Lf \rangle_\nu}{\text{Var}_\nu(f)}$
- Mixing time in χ^2 -divergence $\tilde{O}\left(\frac{1}{\gamma}\right)$

Today: Conductance, Cheeger's inequality

References

- Levin, Peres, and Wilmer, *Markov Chains and Mixing Times*, AMS, 2009
- Lovász and Simonovits, *Random Walks in a Convex Body and an Improved Volume Algorithm*, Random Structures and Algorithms, 1993
- Vempala, *Geometric Random Walk: A Survey*, Combinatorial and Computational Geometry, 2005

From last time: Typo

2) What about P general? (P not sym \Rightarrow complex eigenvalues)

Assume reversibility:

$$\tilde{P}(x,y) = \sqrt{\nu(x)} \cdot \frac{P(x,y)}{\sqrt{\nu(y)}} = \tilde{P}(y,x)$$

symmetric

\Rightarrow spectral theorem

(eg. $\tilde{P} = D^{-\frac{1}{2}} A D^{\frac{1}{2}}$)

should be:

$$\begin{aligned}\tilde{P} &= D^{\frac{1}{2}} P D^{-\frac{1}{2}} \quad (P = D^{\frac{1}{2}} A D^{\frac{1}{2}}) \\ &= D^{-\frac{1}{2}} A D^{\frac{1}{2}}\end{aligned}$$

$$\tilde{L} = I - \tilde{P} = I - D^{-\frac{1}{2}} A D^{\frac{1}{2}} \quad \text{normalized Laplacian}$$

this is equivalent to working $L^2(\nu)$

From last time: Mixing time in χ^2 -divergence

Theorem

Let P be reversible with respect to ν with spectral gap γ . Along the Markov chain $X_k \sim \rho_k$:

$$\chi_\nu^2(\rho_k) \leq (1 - \gamma)^{2k} \chi_\nu^2(\rho_0) \leq e^{-2\gamma k} \chi_\nu^2(\rho_0)$$

Corollary

The mixing time in χ^2 -divergence is:

(assuming n, ϵ constant)
hides log factor

$$\tau(\epsilon) = \frac{1}{2\gamma} \log \frac{\chi_\nu^2(\rho_0)}{\epsilon} = \tilde{O}\left(\frac{1}{\gamma}\right) \leftarrow \tilde{O}\left(\frac{n}{\gamma}\right)$$

$\log \frac{1}{\epsilon}$

$$\log \chi_\nu^2(\rho_0)$$

$\mathcal{O}(n)$

in $X = \mathbb{R}^n$

Proof

Let $h_k = \frac{d\pi_k}{dy}$ be density of π_k wrt ν

P reversible $\Rightarrow h_{k+1} = Ph_k$ (exercise)

$$\Rightarrow \boxed{h_k = P^k h_0}$$

P self-adjoint \Rightarrow has real eigenvalues

$$1 = \lambda_1 > \lambda_2 \geq \dots \geq \lambda_n > -1$$

with orthonormal basis of eigenfunctions

$$1 = \phi_1, \phi_2, \dots, \phi_n$$

Spectral gap $\gamma \Leftrightarrow |\lambda_i| \leq 1-\gamma$ for $i=2,\dots,n$

Write $h_0 = \frac{d\phi}{d\nu}$ as

$$h_0 = 1 + c_2 \phi_2 + \dots + c_n \phi_n$$

where $c_i = \langle h_0, \phi_i \rangle_\nu$ ($c_i = \langle h_0, \phi_i \rangle_\nu = \langle h_0, 1 \rangle_\nu = 1$)

Apply P :

$$\begin{aligned} h_1 &= Ph_0 = P1 + c_2 P\phi_2 + \dots + c_n P\phi_n \\ &= 1 + c_2 \underbrace{\lambda_2}_{\text{P}\phi_2} \phi_2 + \dots + c_n \underbrace{\lambda_n}_{\text{P}\phi_n} \phi_n \end{aligned}$$

$$P\phi_i = \lambda_i \phi_i$$

Repeat:

$$h_k = P^k h_0 = 1 + c_2 \underbrace{\lambda_2^k}_{\text{P}^k \phi_2} \phi_2 + \dots + c_n \underbrace{\lambda_n^k}_{\text{P}^k \phi_n} \phi_n$$

Then:

$$\begin{aligned} \chi^2_\nu(S_k) &= \text{Var}_\nu(h_k) = \|h_k - 1\|_\nu^2 \\ &= \|c_2 \lambda_2^k \phi_2 + \dots + c_n \lambda_n^k \phi_n\|_\nu^2 \\ &= c_2^2 \lambda_2^{2k} + \dots + c_n^2 \lambda_n^{2k} \quad \text{since } \langle \phi_i, \phi_j \rangle_\nu = \begin{cases} 1 & i=j \\ 0 & \text{else} \end{cases} \end{aligned}$$

since $|\lambda_i| \leq 1-\gamma$:

$$\begin{aligned} &\leq (-\gamma)^{2k} (c_2^2 + \dots + c_n^2) \\ &= (-\gamma)^{2k} \|c_2 \phi_2 + \dots + c_n \phi_n\|_\nu^2 \\ &= (-\gamma)^{2k} \|h_0 - 1\|_\nu^2 \\ &= (-\gamma)^{2k} \chi^2_\nu(S_0). \end{aligned}$$

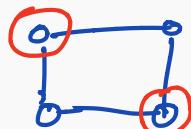
Ex: $S_0 = N(0, \alpha I)$ $\nu = N(0, I)$ $\left\{ \chi^2_\nu(S_0) = \frac{1}{(\alpha(2-\alpha))^{n_2}} \quad 0 < \alpha < 2 \right.$

Examples

Random Walk on Graph:

1) Cycle graph on n vertices

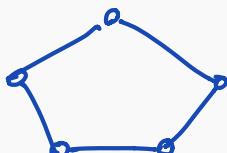
- n even:



\Rightarrow bipartite

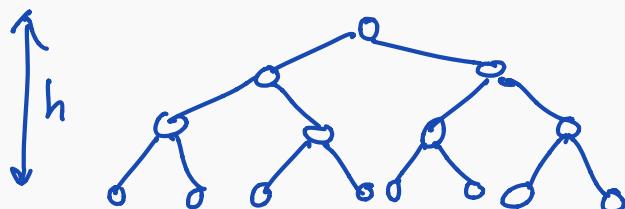
$$\Rightarrow \lambda_n = -1 \Rightarrow \gamma = 0$$

- n odd:



$$\gamma = O\left(\frac{1}{n^2}\right) \Rightarrow \tau = O(n^2)$$

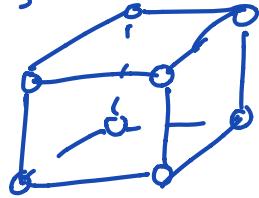
2) Complete binary tree on $n = 2^h - 1$ vertices



$$\gamma = O\left(\frac{1}{n}\right)$$

3) Hypercube on $n = 2^h$ vertices

$$\{0, 1\}^h$$



$$\gamma = O(\frac{1}{h}) = O(\frac{1}{\log n})$$

For general P :

$$\begin{aligned} \gamma &= \inf_{f \in L^2(\nu)} \frac{\langle f, Lf \rangle_\nu}{\text{Var}_\nu(f)} \\ &= \inf_{f \in L^2(\nu)} \frac{\langle f, Lf \rangle_\nu}{\|f\|_\nu^2} \\ \mathbb{E}_\nu[f] &= \langle f, \mathbf{1} \rangle_\nu = 0 \end{aligned}$$

$\mathbb{E}_\nu[(f - \mathbb{E}_\nu[f])^2]$

\tilde{f}

Exercise: Compute / bound spectral gap of OU:

$$X_{t+\eta} = e^{-\eta} X_t + \sqrt{(1-e^{-2\eta})} Z,$$

$$Z \sim N(0, I)$$

on $X = \mathbb{R}^n$, $\eta > 0$ is step size

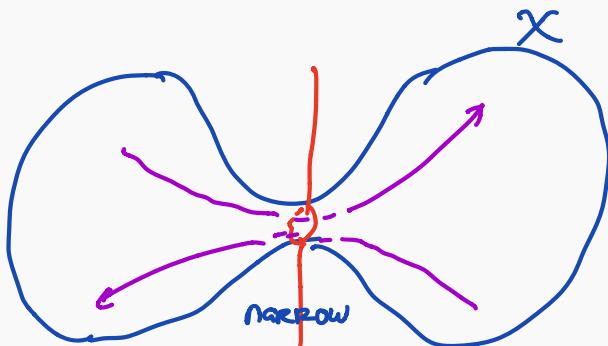
$$\nu = N(0, I)$$

now: more intuitive notion

Conductance

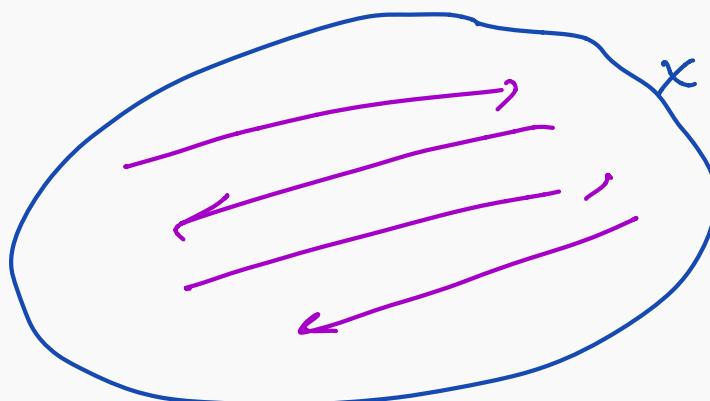
Bottleneck

What prevents a random walk from mixing quickly?



bottleneck

- ⇒ hard to pass through
- ⇒ RW slow to mix



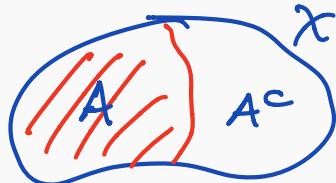
no bottleneck

- ⇒ RW fast mixing

Ergodic flow

Let P be (reversible) wrt. ν on X

A subset $A \subseteq X$ defines a partition / cut



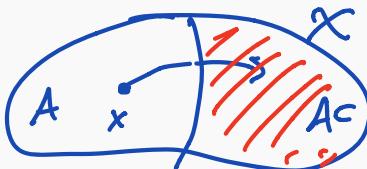
$$X = A \cup A^c$$

" $X - A$ complement

Def: The ergodic flow of A is

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

$$= \mathbb{P}_\nu(X_0 \in A, X_t \notin A)$$



$$\begin{aligned} X_0 &\sim \nu \\ X_t | X_0 &\sim P_{X_0} \end{aligned}$$

Note: Since ν is stationary for P , $\Phi(A) = \Phi(A^c)$

$$\mathbb{P}_\nu(X_0 \in A, X_t \notin A) = \mathbb{P}_\nu(X_0 \notin A, X_t \in A)$$

(does not need reversibility)

Bad case (bottleneck):

- * if have large $\Phi(A)$ with small $\Phi(A)$
- * but can make $\Phi(A)$ small by making A small

\Rightarrow should normalize:
$$\frac{\Phi(A)}{\nu(A)}$$

- * also, since $\Phi(A) = \Phi(A^c)$

$$\nu(A) = 1 - \nu(A^c)$$

\Rightarrow should consider:

cut-ratio/
expansion

$$C(A) = \frac{\Phi(A)}{\min\{\nu(A), 1 - \nu(A)\}}$$

- * Suppose $\nu(A) \leq \frac{1}{2}$:

$$\begin{aligned} C(A) &= \frac{\Phi(A)}{\nu(A)} = \frac{P_R(X_0 \in A, X \notin A)}{P_R(X_0 \in A)} \\ &= P_R(X \notin A \mid X_0 \in A) \end{aligned}$$

Conductance

The **conductance** of Markov chain P is

$$\phi = \inf_{\substack{A \subset \mathcal{X} \\ 0 < \nu(A) < 1}} \frac{\Phi(A)}{\min\{\nu(A), 1 - \nu(A)\}}$$

$$= \inf_{0 < \nu(A) \leq \frac{1}{2}} \Pr(X_i \notin A \mid X_0 \in A)$$

$X_0 \sim \nu$
 $X_i | X_0 \sim P_{X_0}$

Large conductance $\phi \Rightarrow$ every subset has good expansion
 \Rightarrow RW mixes fast.

Conductance \Leftrightarrow Spectral gap

Theorem (Cheeger's inequality)

Let P be reversible with spectral gap γ and conductance ϕ . Then:

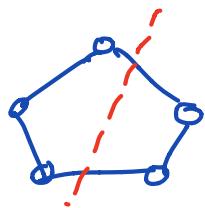
$$\frac{\phi^2}{2} \leq \gamma \leq 2\phi$$

surprising ↓ trivial ↓

- notes:
- * Cheeger [70] for compact Riemannian manifold
 - * Alon [86] for graph
 \Rightarrow used by Jerrum & Sinclair [88] for approx. permanent
 - * Lawler & Sokal [88] for general X
 - * We follow Lovász & Simonovits [93]

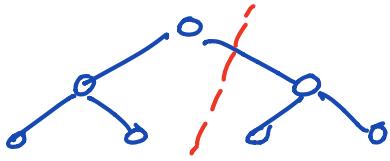
Example: RW on graph

- * Cycle on n vertices: $\phi = O(\frac{1}{n})$, $\gamma = O(\frac{1}{n^2})$



$$\gamma \sim \phi^2$$

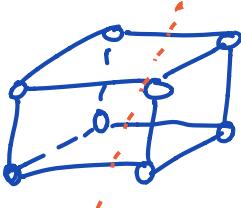
- * Complete binary tree on n vertices



$$\phi = O(\frac{1}{n}), \quad \gamma = O(\frac{1}{n})$$

$$\gamma \sim \phi$$

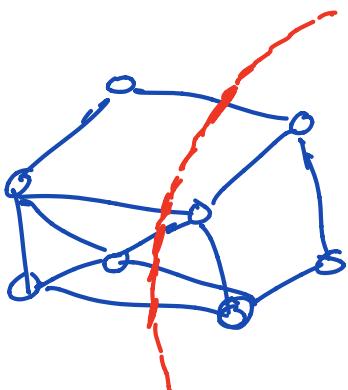
- * Hypercube $\{0,1\}^h$ on $n = 2^h$ vertices



$$\phi = O(\frac{1}{h}) = O(\frac{1}{\log n}), \quad \gamma = O(\frac{1}{\log n})$$

$$\gamma \sim \phi$$

Graph Cut / Partition



Proof of upper bound

Proof:

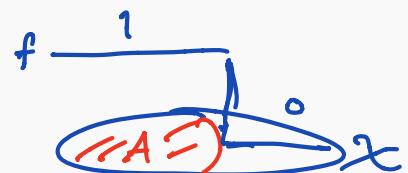
$$\gamma = \inf_{f \in L^2(\nu)} \frac{\langle f, Lf \rangle_\nu}{\text{Var}_\nu(f)}$$

$$\boxed{\gamma \leq 2\phi}$$

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

Given $A \subset X$, let $f = \mathbf{1}_A \in L^2(\nu)$

$$f(x) = \begin{cases} 1 & \text{if } x \in A \\ 0 & \text{else} \end{cases}$$



$$\langle f, Lf \rangle_\nu = \frac{1}{2} \mathbb{E}[(f(x_i) - f(x_0))^2]$$

$$f = \mathbf{1}_A := \frac{1}{2} \mathbb{E}[(\mathbf{1}_A(x_i) - \mathbf{1}_A(x_0))^2]$$

$$= \begin{cases} 1 & \text{if } (x_0 \in A, x_i \notin A) \\ & \text{or } (x_0 \notin A, x_i \in A) \\ 0 & \text{else} \end{cases}$$

$$= \frac{1}{2} (\Pr(X_0 \in A, X \notin A) + \Pr(X_0 \notin A, X \in A))$$

$$= \frac{1}{2} (\mathbb{E}(A) + \mathbb{E}(A^c))$$

$$= \mathbb{E}(A)$$

Bernoulli rv. $\begin{cases} 1 & \text{wp. } v(A) \\ 0 & \text{else} \end{cases}$

and $\text{Var}_v(f) = \text{Var}_v(\mathbf{1}_A)$

$$= v(A) (1 - v(A))$$

$$\geq \frac{1}{2} \min \{v(A), 1 - v(A)\}$$

$$\text{so } \gamma = \inf_f \frac{\langle f, Lf \rangle}{\text{Var}_v(f)} \leq \inf_{ACX} \frac{\mathbb{E}(A)}{\frac{1}{2} \min \{v(A), 1 - v(A)\}} = 2\phi.$$

$f = \mathbf{1}_A$

□

Proof of lower bound

$$\frac{\phi^2}{2} \leq \gamma$$

Proof sketch: from $f \in L^2(\nu)$, need to define a cut.

given f with $\mathbb{E}_\nu[f] = 0$, let m be median of f :

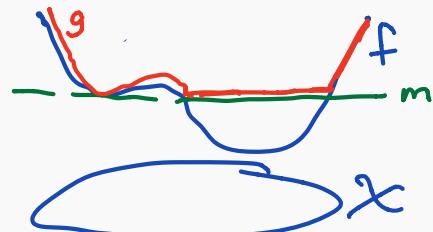
$$\nu(\{x: f(x) > m\}) \leq \frac{1}{2}$$

$$\nu(\{x: f(x) < m\}) \leq \frac{1}{2}$$

Let $g(x) = \max\{f(x) - m, 0\}$

Use g to define cut

See [Lovasz & Simonovits '93 : Lemma 1.6]



□

Mixing time bound via conductance

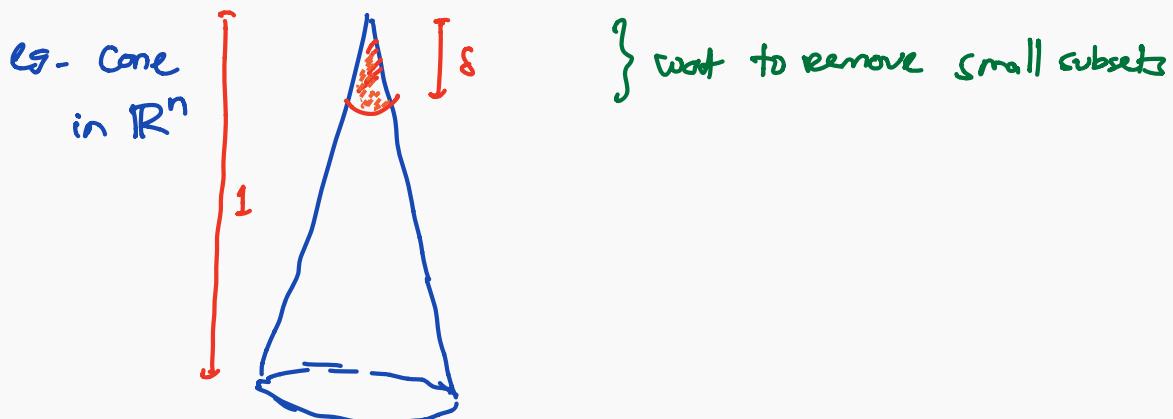
Cheeger:

$$\gamma \geq \frac{\phi^2}{z}$$

\Rightarrow mixing time in χ^2 -divergence is

$$\tau = \tilde{O}\left(\frac{1}{\gamma}\right) = \tilde{O}\left(\frac{1}{\phi^2}\right)$$

But sometimes conductance too strict



s -Conductance

The s -conductance of Markov chain P is

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

$$\phi_s = \inf_{\substack{A \subset \mathcal{X} \\ s < \nu(A) < 1-s}} \frac{\Phi(A)}{\min\{\nu(A) - s, 1 - s - \nu(A)\}}$$

↑ or too large
prevent A too small

- * Usual conductance : $s=0$
- * $s \mapsto \phi_s$ is increasing
- * can get mixing time bound in Total Variation distance
Under Warm-start [Lovász & Simonovits '93]