

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

---

Andre Wibisono

February 8, 2021

## Last time

- Markov chain  $P = (P_x : x \in X)$
- Reversibility (detailed balance)  $\Rightarrow$  Stationary  
 $v(x) \cdot P_x(y) = v(y) \cdot P_y(x) \Rightarrow v(y) = \sum_{x \in X} v(x) \cdot P_x(y)$   
 $\Leftrightarrow$  if  $x_0 \sim v$   
 $x_1 | x_0 \sim P_{x_0}$   
 $\overline{x_1 \sim v}$

Today: Spectral theory (eigenvalue analysis) of reversible MC

# Motivation from discrete space

$$X = \{1, \dots, n\}$$

$$P(X) = \Delta_n = \{ p = (p_1, \dots, p_n) : p_i \geq 0, \sum_{i=1}^n p_i = 1 \}$$

Markov chain  $P$  = stochastic matrix

e.g. RW on graph  $G = (V, E)$

adjacency matrix  $A$

diagonal degree matrix  $D$

Random walk matrix  $P = D^{-1}A$

Laplacian  $L = I - P$

Markov chain:  $X_k \sim S_k$

$X_{k+1} | X_k \sim P_{X_k}$

$$\left. \begin{array}{l} X_{k+1} \sim S_{k+1} = S_k \cdot P \\ \hline \hline = \end{array} \right\} = \boxed{P}$$

Then

$$s_k = s_0 \cdot P^k$$

← Power method

for computing eigenvectors

with max eigenvalue

1) Assume  $P$  is symmetric :  $P_k(y) = P_y(x)$

(eg.  $G$  is a regular graph)

→ linear algebra : Spectral theorem:

$P$  has real eigenvalues  $1 = \lambda_1 \geq \dots \geq \lambda_n$

with orthonormal eigenvectors  $v = v_1, \dots, v_n \in \mathbb{R}^n$

convergence of power method characterized by

$$\text{gap } \lambda_1 - \lambda_2 = 1 - \lambda_2$$

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

Assume reversibility:

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

symmetric

⇒ spectral theorem

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

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

this is equivalent to working  $L^2(v)$

$$L^2(\nu)$$

Let  $\nu$  be a probability distribution on  $X$

Define  $L^2(\nu) = \{ f: X \rightarrow \mathbb{R} \mid \int_X f(x)^2 d\nu(x) < \infty \}$

This is a Hilbert space (complete inner product space)

inner product :  $\langle f, g \rangle_\nu := \int_X f(x)g(x) d\nu(x) = \int fg d\nu = \mathbb{E}_\nu[f g]$

squared norm :  $\|f\|_\nu^2 = \langle f, f \rangle_\nu = \int_X f(x)^2 d\nu(x) = \mathbb{E}_\nu[f^2]$

note: \*  $\mathbb{E}_\nu[f] = \int_X f(x) \cdot 1 d\nu(x) = \langle f, \mathbf{1} \rangle_\nu$

\* if  $\mathbb{E}_\nu[f] = 0$ , then  $\text{Var}_\nu(f) = \mathbb{E}_\nu[f^2] = \|f\|_\nu^2$

## $\chi^2$ -divergence

How to measure distance between  $\pi$  and  $\nu$ ?

$\Rightarrow$  measure squared norm distance in  $L^2(\nu)$

Density  $h = \frac{d\pi}{d\nu} : X \rightarrow \mathbb{R}$  ( $h \in L^2(\nu)$ )

$$\mathbf{1} = \frac{d\nu}{d\nu} : X \rightarrow \mathbb{R}$$

$$(\mathbb{E}_\nu[h] = \int_X h d\nu = \int_X d\pi = \pi(X) = 1)$$

Define :  $\chi^2$ -divergence  $\chi_{\nu}^2(\pi) = \| h - \mathbf{1} \|_\nu^2$

$$= \mathbb{E}_\nu[(h(x) - 1)^2]$$

$$= \text{Var}_\nu(h)$$

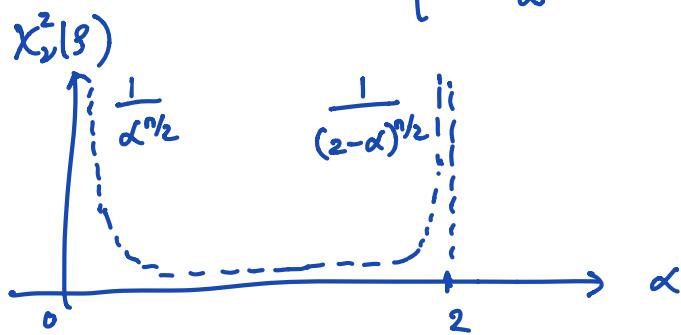
$$= \text{Var}_\nu\left(\frac{d\pi}{d\nu}\right)$$

- note:
- \* not symmetric:  $\chi^2_\nu(s) \neq \chi^2_s(\nu)$
  - \* range of  $\chi^2$  can be large (exponential in  $n$ )

e.g.  $s = \mathcal{N}(0, \alpha I)$ ,  $\alpha > 0$  on  $X = \mathbb{R}^n$

$$v = \mathcal{N}(0, I)$$

$$\Rightarrow \chi^2_\nu(s) = \begin{cases} \frac{1}{(\alpha(2-\alpha))^{n/2}} & \text{if } 0 < \alpha < 2 \\ \infty & \text{else} \end{cases}$$



# Markov chain as linear operator

Let  $P$  be reversible wrt  $\nu$  on  $X$ .

$P$  defines a linear operator :

$$\Phi : L^2(\nu) \rightarrow L^2(\nu)$$

$$f \mapsto Pf$$

given by  $(Pf)(x) = \int_X f(y) dP_x(y) = \mathbb{E}[f(X_i) | X_0 = x]$

notes: \* in discrete space, this is matrix mult. on the left:

$$\begin{array}{ccc} \| & \longmapsto & \| \\ f & & P \cdot f \end{array}$$

\* recall for distributions, matrix mult. on the right:

$$\begin{array}{ccccc} \overbrace{\hspace{1cm}} & \longmapsto & \overbrace{\hspace{1cm}} & \boxed{\phantom{000}} & = \overbrace{\hspace{1cm}} \\ s_n & & s_n \cdot P & & s_{n+1} \end{array}$$

# Reversibility $\Rightarrow$ Symmetry

**Lemma**

$$: L^2(\nu) \rightarrow L^2(\nu)$$

Assume  $P$  is reversible with respect to  $\nu$ . Then  $P$  is self-adjoint:

$$\forall f, g \in L^2(\nu) : \quad \langle f, Pg \rangle_\nu = \langle Pf, g \rangle_\nu$$

Proof:  $\langle f, Pg \rangle_\nu = \int f(x) (Pg)(x) d\nu(x)$

$$= \iint f(x) g(y) \underbrace{dP_x(y)}_{\downarrow} d\nu(x)$$
$$= dP_g(x) d\nu(y) \text{ by reversibility}$$
$$= \iint f(x) g(y) dP_g(x) d\nu(y)$$
$$= \int g(y) (Pf)(y) d\nu(y)$$
$$= \langle Pf, g \rangle_\nu$$

□

# Spectral theorem

## Theorem

A self-adjoint operator  $P$  on  $L^2(\nu)$  has real eigenvalues  $\lambda_i \in \mathbb{R}$ :

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

for  $i = 1, 2, \dots, n$  ( $n = \infty$ ) with an orthonormal basis of eigenfunctions  $\phi_i \in L^2(\nu)$ :

$$\langle \phi_i, \phi_j \rangle_\nu = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{else} \end{cases}$$

- We assume  $L^2(\nu)$  is separable (has countable basis);  
see [BGL, Appendix A.4]<sup>1</sup> for general/continuous spectrum
- e.g.  $\mathcal{X} = \mathbb{R}^n$ ,  $\nu = \mathcal{N}(0, I)$ : Eigenfunctions are Hermite polynomials

# Courant-Fischer characterization of eigenvalues

## Theorem

Assume  $P$  is self-adjoint with eigenvalues  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$  ( $n = \infty$ ) and corresponding eigenfunctions  $\phi_1, \dots, \phi_n \in L^2(\nu)$ .

Then:

$$\lambda_1 = \max_{\substack{0 \neq f \in L^2(\nu)}} \frac{\langle f, Pf \rangle_\nu}{\langle f, f \rangle_\nu}$$

$$\lambda_2 = \max_{\substack{0 \neq f \in L^2(\nu) \\ f \perp \phi_1}} \frac{\langle f, Pf \rangle_\nu}{\langle f, f \rangle_\nu}$$

$$(\langle f, \phi_1 \rangle_\nu = 0)$$

$$\lambda_n = \min_{\substack{0 \neq f \in L^2(\nu)}} \frac{\langle f, Pf \rangle_\nu}{\langle f, f \rangle_\nu}$$

# Bounds on eigenvalues

## Lemma

Assume  $P$  is reversible with respect to  $\nu$ . For all  $f \in L^2(\nu)$ :

✓ (a)  $\langle f, Pf \rangle_\nu \leq \langle f, f \rangle_\nu$

} write  $-I \leq P \leq I$

✓ (b)  $\langle f, Pf \rangle_\nu \geq -\langle f, f \rangle_\nu$

( $\Leftrightarrow \forall f: -\langle f, f \rangle_\nu \leq \langle f, Pf \rangle_\nu \leq \langle f, f \rangle_\nu$ )

Proof: Write

$$\langle f, f \rangle_\nu = \frac{1}{2} \int_X f(x)^2 d\nu(x) + \frac{1}{2} \int_X f(y)^2 d\nu(y)$$

$$= \frac{1}{2} \int_X \int_X f(x)^2 \underbrace{dP_x(y)}_{\text{equal b, rev.}} d\nu(x) + \frac{1}{2} \iint_X f(y)^2 \underbrace{dP_y(x)}_{\text{equal b, rev.}} d\nu(y)$$

$$= \frac{1}{2} \int_X \int_X (f(x)^2 + f(y)^2) dP_x(y) d\nu(x)$$

Also,  $\langle f, Pf \rangle_\nu = \int_X \int_X f(x) f(y) dP_x(y) d\nu(x)$

$$\begin{aligned}
 a) \quad & \langle f, f \rangle_{\nu} - \langle f, Pf \rangle_{\nu} = \frac{1}{2} \int_X \int_X (f(x)^2 + f(y)^2 - 2f(x)f(y)) dP_X(y) d\nu(x) \\
 & \underbrace{\qquad\qquad\qquad}_{= \langle f, (I - P)f \rangle_{\nu}} = \frac{1}{2} \int_X \int_X (f(x) - f(y))^2 dP_X(y) d\nu(x) \\
 & = \langle f, Lf \rangle_{\nu} \geq 0
 \end{aligned}$$

$$\begin{aligned}
 b) \quad & \langle f, f \rangle_{\nu} + \langle f, Pf \rangle_{\nu} = \frac{1}{2} \int_X \int_X (f(x) + f(y))^2 dP_X(y) d\nu(x) \\
 & \geq 0
 \end{aligned}$$

□

So all eigenvalues of  $P$  satisfy:

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

$\lambda_2 < 1$  if  $P$  is irreducible

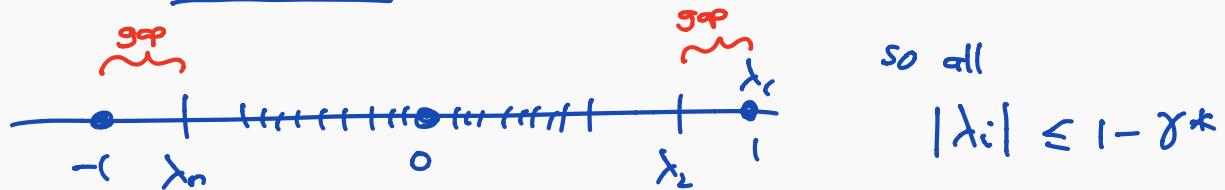
$\lambda_n = -1$  iff bipartite / periodic

# Spectral Gap

P has eigenvalues

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

Def: (Absolute) Spectral gap :  $\gamma^* = \min \{ 1 - \lambda_2, 1 + \lambda_n \}$



Notes: Can shift all  $\lambda_i \geq 0$  by Lazyfying P:

$$\hat{P} = \frac{1}{2} I + \frac{1}{2} P$$

(at each step:  
• w.p.  $\frac{1}{2}$ , stay  
• w.p.  $\frac{1}{2}$ , jump following P)

then  $\hat{P} \geq 0$

$$\begin{aligned} (\text{since } \langle f, \hat{P}f \rangle_{\nu} &= \frac{1}{2} \langle f, f \rangle_{\nu} + \frac{1}{2} \langle f, Pf \rangle_{\nu}, \\ &\geq 0) \end{aligned}$$

Def: Spectral Gap:  $\gamma = 1 - \lambda_2$

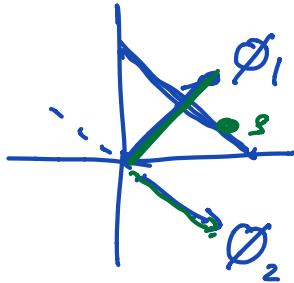
( $\gamma = \gamma^*$  if  $P$  is lazy)



$$g_1 = g_0 \cdot P$$

$$g_1(y) = \sum_x g_0(x) \cdot P_x(y) \geq 0$$

$\underbrace{\phantom{0}}_{\geq 0} \quad \underbrace{\phantom{0}}_{\geq 0}$



---

$P$  has max eigenvalue  $\lambda_1 = 1$

with eigenvector  $\phi_1 = \mathbb{1}$

# Laplacian

Def: Laplacian  $L = I - P : L^2(\nu) \rightarrow L^2(\nu)$

P reversible wrt.  $\nu \Rightarrow L$  is symmetric, positive semidefinite:

Dirichlet form  $E(f, f) := \langle f, Lf \rangle_\nu$

$$= \frac{1}{2} \mathbb{E}_{(\nu, P)} [(f(x_i) - f(x_0))^2] \geq 0$$

$$\boxed{L \geq 0}$$

$x_0 \sim \nu$   
 $x_i | x_0 \sim P_{x_0}$

eigenvalues  $0 = \lambda_1 < \lambda_2 \leq \dots \leq \lambda_n (< 2)$

eigenfunction  $\mathbf{1} = \phi_1, \phi_2, \dots, \phi_n$

Spectral Gap:  $\gamma = \lambda_2(L) = \min_{f \in L^2(\nu)} \frac{\langle f, Lf \rangle_\nu}{\langle f, f \rangle_\nu} \approx \min_{f \in L^2(\nu)} \frac{E(f, f)}{\text{Var}_\nu(f)}$

$\langle f, \mathbf{1} \rangle_\nu = 0$

# Mixing time in $\chi^2$ -divergence

## Theorem

Let  $P$  be a Markov chain reversible with respect to  $\nu$  with spectral gap  $\gamma$  ( $= \gamma^*$ ). For any  $X_0 \sim \rho_0$ , along the Markov chain  $X_k \sim \rho_k$ :

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

$\Rightarrow$  Corollary: To reach  $\chi_\nu^2(\rho_k) \leq \varepsilon$ , enough to take

mixing time  
in  $\chi_\nu^2$  :

$$T(\varepsilon, \rho_0) = \frac{1}{2\gamma} \log \frac{\chi_\nu^2(\rho_0)}{\varepsilon}$$

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

$$\Leftrightarrow k \geq \frac{1}{2\gamma} \log \frac{\chi_\nu^2(\rho_0)}{\varepsilon}$$

e.g.,  $\rho_0, \nu$  Gaussian,  $\chi_\nu^2(\rho_0) \sim \frac{1}{\alpha^n}$ , so  $\log \chi_\nu^2(\rho_0) = O(n)$

## Proof

$$X_k \sim p_k, \quad X_{k+1} \sim p_{k+1}$$

density  $h_k = \frac{dp_k}{d\nu}, \quad h_{k+1} = \frac{dp_{k+1}}{d\nu} \in L^2(\nu)$

Lemma:  $P$  reversible wrt  $\nu \Rightarrow$  
$$\boxed{h_{k+1} = Ph_k}$$

$$\Rightarrow \text{then } h_k = P^k h_0$$

$$\text{do eigenvalue analysis of } P \Rightarrow \chi_\nu^2(p_k) = \text{Var}_\nu(h_k) = \|h_k - \mathbf{1}\|_\nu^2$$