

- today:
- MC MC
 - review Markov chains
 - M-H. alg.

MCMC - Markov chain Monte-Carlo

idea: is to relax indep assumption samples

to allow adaptive proposal dist.

i.e. we'll run a chain $X_t | X_{t-1}$, s.t. $X_t \xrightarrow{t \rightarrow \infty}$ in dist. is target dist. p

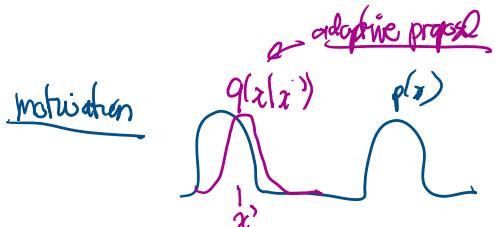
"stationary dist. of a chain"

then we can approximate

$$\mathbb{E}_p[f(x)] \text{ as } \frac{1}{T-T_0} \sum_{t=T_0+1}^T f(x_t)$$

T_0 is called "burn in period" \rightarrow depends on "mixing time" of Markov chain

- ⊕ no need to thin the samples [i.e. use Δt between samples to get more independence] as this yield higher variance \rightarrow better to use all samples after T_0 to estimate μ



[unless it is too expensive]

before: samples $x^{(t)}$ i.i.d. q

MCMC $x^{(t)} | x^{(t-1)} \sim q(\cdot | x^{(t-1)})$

P
Markov transition prob

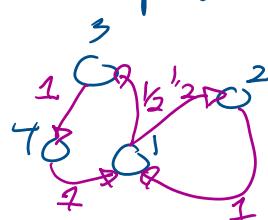
Review of (finite state space) Markov chain $[|X|=k]$

- as PGM $x^{(0)} \xrightarrow{\quad} x^{(1)} \xrightarrow{\quad} \dots \xrightarrow{\quad} x^{(t-1)} \xrightarrow{\quad} x^{(t)}$

- there is also a transition prob pt. of view: use one node per state

(probabilistic FSA)

e.g. $k=4$



[homogeneous M.C.]

$$\hookrightarrow \text{e.g. } P\{X_t = i \mid X_{t-1} = j\} = A_{ij} \quad (\text{no time dep.})$$

A is a $K \times K$ matrix s.t. $\mathbb{1}_K^T A = \mathbb{1}_K^T$

"left-stochastic matrix"

vector of ones of size K

(*) (as in HMM) suppose $P\{X_{t-1} = j\} = (\pi)_j$

$$P\{X_t = i\} = \sum_j \underbrace{P\{X_t = i \mid X_{t-1} = j\}}_{A_{ij}} \underbrace{P\{X_{t-1} = j\}}_{\pi_j}$$

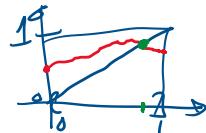
$$\begin{aligned} \pi_t &= A \pi_{t-1} \\ \Rightarrow \pi_t &= A^t \pi_0 \end{aligned}$$

stationary dist π of A is a dist. π s.t. $A\pi = \pi$

[note that this π is a right e-vector of A with e-value 1]

fact: every stochastic matrix has at least 1 stat. dist?

(by Brower's fixed pt. thm.)



def: irreducible Markov chain \Leftrightarrow there exists a positive prob "path" from any i to any j (states)

$\forall (i, j), \exists$ an integer $m_{ij}^{(positive)}$ s.t. $(A^{m_{ij}})_{ji} > 0$

(by Perron-Frobenius Thm. \Rightarrow irreducible M.C., has a unique stat. dist.)
 (multiplicity of e-value 1 is 1)

(*) In order to converge to it, we need aperiodicity as well

Irreducible and aperiodic M.C. $\Leftrightarrow \exists$ an integer m s.t. $(A^m)_{ij} > 0$

aka regular M.C. (finite state space)

$((A^m)_{ij} > 0 \forall i, j)$

or ergodic M.C.

(*) [note: a sufficient condition for an irreducible M.C. to]

be aperiodic is $\exists i$ s.t. $A_{ii} > 0$

example of a regular M.C. on K states ($K \geq 3$)

$$A = \begin{pmatrix} 0 & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & 0 & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} & 0 \end{pmatrix} = \frac{1}{K-1} (K\mathbb{I} - I)$$

$$A^2 = \frac{1}{(K-1)^2} (K\mathbb{I} - 2K\mathbb{I} + I) \\ = \frac{1}{(K-1)^2} [(K-2)\mathbb{I} + I]$$

$\boxed{\text{This is } > 0}$
 $m=2$

[but, for $K=2$ it is not aperiodic $A^2 = I$ $A = \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix}$]

Thm. If a finite M.C. is ergodic (regular)

then \exists a unique st. dist. π

and for any starting dist π_0 , $\lim_{t \rightarrow \infty} A^t \pi_0 = \pi$

the speed of convergence is related to the mixing time γ of the chain

$$\gamma \triangleq \frac{1}{1 - |\lambda_2(A)|}$$

\uparrow 2nd largest e-value of A

$$\boxed{\|A^t \pi_0 - \pi\|_1 \leq C \exp(-t/\gamma)}$$

after \sim steps, error decreases by $\frac{1}{e}$

Ex 24

② intuition (from linear algebra) [informal argument]

simpler case, suppose A is symmetric \Rightarrow psd

by spectral thm., A is diagonalizable with orthogonal U

$$A = U \Sigma U^T \text{ with } \Sigma = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_K \end{pmatrix}$$

[by Perron-Frobenius Thm]

$$\lambda_1 = 1 > |\lambda_2| \geq \dots \geq \lambda_K$$

$$\downarrow \\ U_1 = \frac{\pi}{\|\pi\|_2} \quad A\pi = \pi$$

$U \in \mathbb{R}^{n \times n}$ of e-vectors of A

$$U U^T = U \bar{U} = I \quad (\bar{U} = (U_1 \dots U_K))$$

Let α_0 be coordinates of π_0 in U basis

$$\pi_0 = \alpha_0 \bar{U} = \alpha_0 U^T \pi$$

Let α_0 be coordinates of π_0 in U basis

$$\begin{aligned} & \text{i.e. } \pi_0 = U\alpha_0 \quad (\alpha_0 = U^T \pi_0) \\ A^T \pi_0 &= (U\Sigma U^T)(U\Sigma U^T) \cdots (U\Sigma U^T)(U\alpha_0) \\ &= U\Sigma^t \alpha_0 \\ &= (\alpha_0)_1 U_1 + (\alpha_0)_2 U_2 + \dots + (\alpha_0)_k U_k \end{aligned}$$

$\Sigma^t = \begin{pmatrix} 1 & & & \\ & \lambda_2 & & \\ & & \lambda_3 & \\ & & & \ddots & \lambda_k \end{pmatrix}$

$$\|A^T \pi - \pi\|_1 = \left\| \frac{(\alpha_0)_1}{\|\alpha_0\|_2} \right\|_1 \quad \text{(if } \frac{(\alpha_0)_1}{\|\alpha_0\|_2} = 1 \text{ not true in general)} \quad \| \leq C |\lambda_2|^t \quad \text{fast e.g. gap} \\ |\lambda_2| = 1 - \varepsilon_1 \quad \varepsilon_1 \triangleq 1 - |\lambda_2|$$

$$1+x \leq \exp(x) \quad \forall x \quad |\lambda_2| \leq \exp(-\varepsilon_1) \quad |\lambda_2|^t \leq \exp(-t\varepsilon_1) \Rightarrow \uparrow = \frac{1}{1-|\lambda_2|}$$

⊗ mixing time is often (usually) exponentially long □

⊗ How do we design A s.t. $A^T \pi_0 \rightarrow \pi$?

one "easy way" reversible M.C. $\Rightarrow \exists$ dist. π s.t. $A_{ij}\pi_j = A_{ji}\pi_i \quad \forall (i,j)$
 "Detailed Balance equation"

this is sufficient condition

It means $\sum_j P_{ij}x_j = \pi_i$

to get $A\pi = \pi$ detailed balance then $P_{ij}x_j = i, x_{j-1} = j$

proof: $(A\pi)_i = \sum_j A_{ij}\pi_j = \sum_j A_{ji}\pi_i = \pi_i \left(\sum_j A_{ji} \right) = \pi_i \left(\sum_j A_{ji} \right) = \pi_i \left(\sum_j A_{ij} \right) = \pi_i$

Metropolis-Hastings alg:

goal \rightarrow construct a M.C. with flat dist. $p(x)$ [can target]

[assume $p(x) > 0 \quad \forall x$]

use some proposal $q(x'|x)$

accept new state x' with prob. $\min\{1, \frac{q(x|x')p(x')}{q(x'|x)p(x)}\}$

no dependence on normalization of p

accept new state x' with prob. $a(x'|x) \triangleq \min\{1, \frac{q(x|x')p(x')}{q(x'|x)p(x)}\}$
 if reject \rightarrow stay in same state x
 [this is still new sample]
 vs.
 rejection sampling where only "accepted states" are samples

acceptance ratio
 to satisfy detailed balance

M.H. alg.:

start at $x^{(0)}$

important design choice

for $t=1, \dots$

• propose $x^{(t)}$ $\sim q(x'|x^{(t-1)})$

flip a coin, with prob. $a(x'|x^{(t-1)})$ to be 1

• if accept (coin=1)

let $x^{(t)} = x'^{(t)}$

else

let $x^{(t)} = x^{(t-1)}$

end for

note: for symmetric proposal $q(x'|x) = q(x|x')$; always accept if $p(x') \geq p(x)$

[Metropolis alg.] \rightsquigarrow like noisy hill-climbing

[verify as exercise that M.H. satisfies detailed balance eq. with $\pi = p^{\text{target}}$]

⊗ for convergence: If M.H. chain is ergodic, then converge to unique stat. dist. ρ

sufficient conditions

for irreducibility $q(x'|x) > 0 \forall x \neq x' \in X$

for aperiodicity

\downarrow
 $A_{ii} > 0$
 for some i

either $q(x|x) > 0$ for some $x \in X$

or $a(x'|x) < 1$
 for some $x \in X$
 and $x' \in X$
 $q(x'|x) > 0$

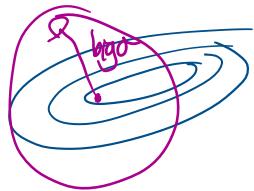
⊗ aside: it is still ok change proposal with time

$q_t(x'|x)$ as long as choice of q_t does not depend on $x^{(t-1)}$

then convergence theory above will go through

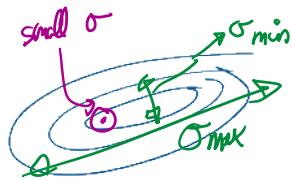
slow mixing example

suppose p is a $N(\mu, \Sigma)$



→ high prob rejection

$$\{ q(x'|x) = N(x'|x, \sigma^2 I) \}$$



"small o" → many steps needed

here best mixing time related to $\frac{\sigma_{\max}}{\sigma_{\min}}$

reference for mixing times:

Markov Chains and Mixing Times

David A. Levin, Yuval Peres, Elizabeth L. Wilmer

<https://pages.uoregon.edu/dlevin/MARKOV/markovmixing.pdf>