

today: • MCMC
• review Markov chains
• M-H alg.

MCMC - Markov chain Monte Carlo

Idea: is to relax indep assumption samples
to allow adaptive proposal dist.

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

"stationary dist. of chain"

then we can approximate

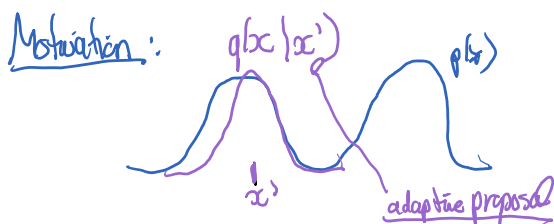
$$\mathbb{E}_p[f(x)] \text{ as } \frac{1}{T-T_0} \sum_{t=T_0+1}^T f(x_t)$$

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

(*) no need thin the samples [ie. use Δt between samples to get more independence]

as this yields higher variance

\rightarrow better to use all samples after T_0 to estimate μ
[unless it is too expensive]

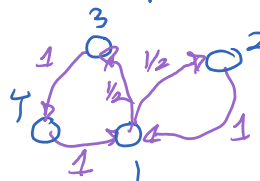


before: samples were $X^{(t)} \stackrel{iid}{\sim} q$
MCMC $X^{(t)} | X^{(t-1)} \sim q(\cdot | X^{(t-1)})$
Markov transition prob.

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

• as a DAG, $x^{(0)} \rightarrow x^{(1)} \rightarrow x^{(2)} \rightarrow \dots \rightarrow x^{(t-1)} \rightarrow x^{(t)}$

• there is also transition prob. of revs: use one per state
(probabilistic FSA) e.g. $k=4$



[homogeneous M.C.]

\hookrightarrow ie. $P\{X_t=i | X_{t-1}=j\} = A_{ij}$ (no time dep.)

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

"A is a stochastic matrix"

vector of ones $1 \times k$

A is a $k \times k$ matrix st. $\mathbb{1}_k^T A = \mathbb{1}_k^T$
 "left-stochastic matrix" vector of ones of size k

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

$$P\{X_t=i\} = \sum_j P\{X_t=i | X_{t-1}=j\} P\{X_{t-1}=j\}$$

$$= \sum_j A_{ij} (\pi)_j$$

$$\pi_{t+1} = A \pi_t$$

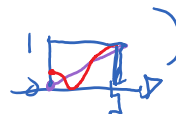
$$\Rightarrow \boxed{\pi_t = A^t \pi_0}$$

Stationary dist. π of A is a dist. π st. $A\pi = \pi$

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

Fact: every stochastic matrix has at least 1 stat. dist!

(by Brouwer's fixed pt. thm.)



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

$$\forall (i,j), \exists \text{ an } \overset{\text{positive}}{\text{integer}} m_{ij} \text{ st. } (A^{m_{ij}})_{ij} > 0$$

(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 $\overset{\text{positive}}{\text{integer}} m$ st. $(A^m)_{ij} > 0$

aka regular M.C. (finite state space)

or ergodic M.C.

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

(*) [note: a sufficient condition for an irreducible M.C. to be aperiodic is $\exists i$ st. $A_{ii} > 0$]

example of a regular M.C. on k states

$$A = \begin{pmatrix} 0 & 1/k & 1/k \\ 1/k & 0 & 1/k \\ 1/k & 1/k & 0 \end{pmatrix} = \frac{1}{k-1} (\mathbb{1}\mathbb{1}^T - I)$$

$$A^2 = \frac{1}{(k-1)^2} (\mathbb{1}\mathbb{1}^T \mathbb{1}\mathbb{1}^T - 2\mathbb{1}\mathbb{1}^T + I)$$

$$= \frac{1}{(k-1)^2} \left(\underset{k}{k} \mathbb{1}\mathbb{1}^T - 2\mathbb{1}\mathbb{1}^T + I \right) \text{ for } k \geq 3, \text{ this } > 0$$

$$= \frac{(k-1)^2}{(k-1)^2} \left(\frac{1}{k-1} \mathbf{1}\mathbf{1}^T + \mathbf{I} \right) \text{ for } k \geq 3, \text{ this } > 0$$

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

thm: if a finite M.C. is ergodic (regular)

then \exists a unique stat. dist. π

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

16h03

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

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

λ_2 is 2nd biggest e-value of A

$$\|A^t \pi_0 - \pi\|_1 \leq C \exp(-t/\tau)$$

after τ steps, error decreases factor $\frac{1}{e}$

⊛ intuition (from linear algebra) [informal argument]

Simpler case, suppose A is symmetric ^{is p.s.d.}, by spectral theorem, A is diagonalizable with orthogonal matrix U

$$A = U \Sigma U^T \text{ with } \Sigma = \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_k \end{pmatrix}$$

$U \rightarrow$ basis of e-vectors

[by Perron-Frobenius thm

$$UU^T = U^T U = \mathbf{I} \quad U = (u_1, \dots, u_k)$$

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

$$\downarrow$$

$$u_1 = \pi \leftarrow A\pi = \pi$$

let α_0 be coordinates of π_0 in U basis

$$\text{i.e. } \pi_0 = U \alpha_0 \quad (\alpha_0 = U^T \pi_0)$$

$$A^t \pi_0 = (U \Sigma U^T) (U \Sigma U^T) \dots (U \Sigma U^T) (U \alpha_0)$$

$$= U \Sigma^t \alpha_0$$

$$\Sigma^t = \begin{pmatrix} 1 & & 0 \\ & \lambda_2^t & \\ 0 & & \lambda_k^t \end{pmatrix}$$

$$A^t \pi_0 = (\alpha_0)_1 \underbrace{1^t}_{\frac{\pi}{\|\pi\|_1}} u_1 + (\alpha_0)_2 \lambda_2^t u_2 + \dots + (\alpha_0)_k \lambda_k^t u_k$$

$$\Rightarrow \text{ (if } \frac{(\alpha_0)_1}{\|\pi\|_1} = 1 \text{) [fishy]}$$

just e.g.p

$$\|A^t \pi - \pi\|_2 = \|(\alpha_0)_2 \lambda_2^t u_2 + \dots\| \leq C |\lambda_2|^t \quad |\lambda_2| = 1 - \epsilon_2 \quad \epsilon_2 = 1 - |\lambda_2|$$

$$|\lambda_2| \leq \exp(-\epsilon_1) \quad 1 - x \leq \exp(-x) \quad \forall x$$

$$|\lambda_2|^t \leq \exp(-t \epsilon_1)$$

$$\frac{1}{\gamma} \Rightarrow \gamma = \frac{1}{1 - |\lambda_2|}$$

⊛ mixing time is often (usually) exponentially big!

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

one "easy way"

reversible M.C. $\Leftrightarrow \exists$ dist. π s.t. $A_{ij}\pi_j = A_{ji}\pi_i \quad \forall (i,j)$

"detailed balance equation"

this is a sufficient condition (but not necessary) to get $A\pi = \pi$

it means when $P\{X_{t-1}=i\} = \pi_i$

then $P\{X_t=i, X_{t-1}=j\} = P\{X_t=j, X_{t-1}=i\}$

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

Metropolis Hastings alg.:

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

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

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

accept new state x' with prob. $\alpha(x'|x)$
if reject it \rightarrow stay in same state

[this is still new sample]

vs. rejection sampling where only "accepted states" are new samples.

$$\alpha(x'|x) \triangleq \min \left\{ 1, \frac{q(x|x') p(x')}{q(x'|x) p(x)} \right\}$$

no dependence on normalization of p

acceptance ratio to satisfy detailed balance

M.H. alg

start at $x^{(0)}$

for $t=1, \dots$

important design

start at $x^{(0)}$
 for $t=1, \dots$,
 • propose $x^{(t)} \sim q(x' | x^{(t-1)})$ important design choice
 • flip a coin, with prob. $q(x^{(t)} | x^{(t-1)})$ to be 1
 • if accept (coin=1)
 let $x^{(t)} = x^{(t-1)}$
 o.w.
 $x^{(t)} = x^{(t-1)}$
 end for

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

→ like a noisy hill-climbing alg.

[Metropolis alg.]

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

⊛ for convergence: if M.H. chain is ergodic, then converge to unique stationary

sufficient conditions

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

for aperiodicity

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

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

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

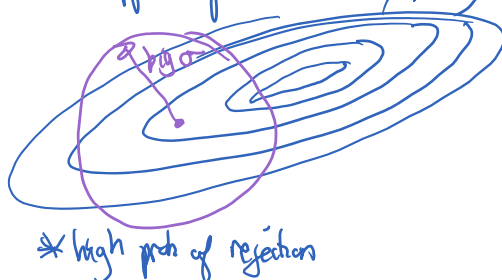
⊛ aside: it is still ok to change proposal with time (in homogeneous M.C.) $q_t(x' | x)$

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

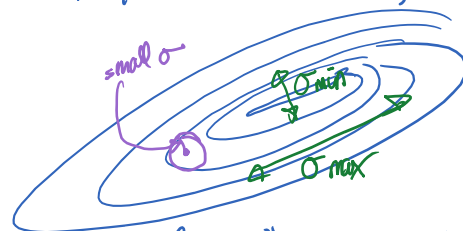
then convergence theory will go through

slow mixing example:

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



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



here the best mixing time

is related to ratio $\frac{\sigma_{\max}}{\sigma_{\min}}$

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>