

today ∴ MCMC

- review Markov chains
- M-H algo

MCMC - Markov chain Monte Carlo

idea ∴ is to relax indep. assumption between samples

to allow adaptive proposal dist.

i.e. we'll run a chain $X_t | X_{t-1}$ st. $X_t \xrightarrow{t \rightarrow \infty}$ in dist. to target dist. p
 "stationary dist. of chain"

then we can approximate

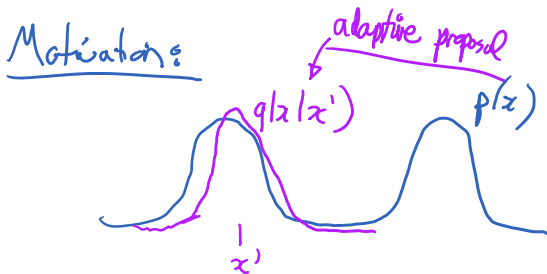
$$E_p[f(x)] \text{ as } \frac{1}{T - T_0} \sum_{t=T_0+1}^T f(x_t)$$

T_0 is called "burn-in period" \rightsquigarrow 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 were $X^{(k)} \overset{iid}{\sim} q$

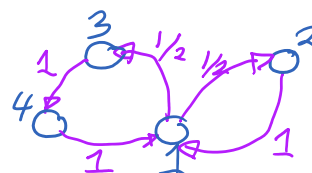
MCMC $X^{(k)} | X^{(k-1)} \sim q(\cdot | X^{(k-1)})$
 Markov transition prob.

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

• as a DGM, $X^{(0)} \rightarrow X^{(1)} \rightarrow \dots \rightarrow X^{(k-1)} \rightarrow X^{(k)}$

• there is also transition prob. pt. of view: we use one node per state (probabilistic FSA)

eg. $k=4$



[homogeneous M.C.]

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

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

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

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

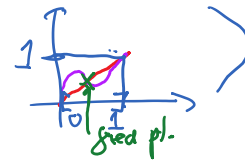
$$\begin{aligned} \pi_{t+1} &= A \pi_t \\ \Rightarrow \pi_t &= A^T \pi_0 \end{aligned}$$

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

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

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

(by Perron's fixed pt. thm.)



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

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

(by Perron-Frobenius thm.) \Rightarrow irreducible M.C. has a unique stat. dist.

[multiplicity of e-value 1 = 1]

* in order to converge to it, we need aperiodicity as well

irreducible and aperiodic M.C. $\Leftrightarrow \exists$ an integer m s.t. $A^m > 0$

aka regular M.C. (every state space)

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

or ergodic M.C.

* [notes 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

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

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

[but, for $k=2$, it is not aperiodic; $A^2 = I$

thm.: if a finite M.C. is ergodic (regular)
 then \exists a unique stationary 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 τ of the chain

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

← 2nd biggest ϵ -value of A

$$\|A^t \pi_0 - \pi\|_1 \leq C \exp(-t/\tau) \quad (\text{not } 100\% \text{ sure})$$

14h32

after τ steps, error decreases $\frac{1}{2}$

⊛ intuition (from linear algebra) [informed argument]

simpler case, suppose A is diagonalizable with orthogonal matrix U (here A symmetric)

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

$U \rightarrow$ basis of ϵ -vectors

$$U = (u_1, \dots, u_k)$$

$$U^T U = U U^T = I$$

[by Perron-Frobenius thm., $\lambda_1 = 1 > |\lambda_2| \geq \dots \geq |\lambda_k|$

$$u_1 = \frac{\pi}{\|\pi\|_2} \quad A \pi = \pi$$

let α_0 be coordinate of π_0 in U basis

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

$\rightarrow I$

$\rightarrow I$

\dots

ie. $\pi_0 = U \alpha_0$ ($\alpha_0 = U^T \pi_0$)

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

$$(\alpha_0)_1 = \frac{\langle \pi_0, \pi \rangle}{\|\pi\|_2}$$

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

$$\Sigma^t = \begin{pmatrix} \lambda_1^t & 0 \\ & \ddots \\ 0 & \lambda_k^t \end{pmatrix}$$

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

(if $\frac{(\alpha_0)_1}{\|\pi\|_2} = 1$) [fishy]

$$\|A^t \pi_0 - \pi\|_2 = \|(\alpha_0)_2 \lambda_2^t u_2 + \dots\| \leq C |\lambda_2|^t$$

first e-gap

$$|\lambda_2| = 1 - \epsilon_1 \quad \epsilon_1 \triangleq 1 - |\lambda_2|$$

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

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

⊛ mixing time is often (usually) exponentially big?

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

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

one "easy way"

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

"detailed balance equation"

this is sufficient condition to get $A\pi = \pi$

(but not necessary)

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

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

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 $q(x'|x)$

does not depend on p

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

accept new state x' with prob.
if reject it \rightarrow stay in same state

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

acceptance ratio to satisfy detailed balance

[this is still new sample]

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

M-H alg.

start at $x^{(0)}$
for $t=1, \dots,$

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

• flip a biased coin with prob. $a(x^{(t+1)}|x^{(t)})$ to be 1

• if accept (coin=1)
let $x^{(t+1)} = x^{(t)}$

o.w.

$x^{(t+1)} = x^{(t)}$

end for

important design choice

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

\rightarrow like a noisy hill-climbing alg.

[Metropolis's alg.]

[verify as exercise that it satisfies detailed balance with $\pi = p$ target]

* for convergence: if M.H. chain is ergodic, then we converge to unique stat dis p

sufficient conditions \leftarrow for irreducibility $q(x'|x) > 0 \forall x, x' \in X$

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

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

or $a(x'|x) < 1$ for some $x, x' \in X$ s.t. $q(x'|x) > 0$

⊗ aside: it is still de. to change proposal with time
(inhomogeneous M.C.) $q_t(x'|x)$

as long as choice of q_t does not depend on $x^{(t-1)}$
then convergence theory above will go through

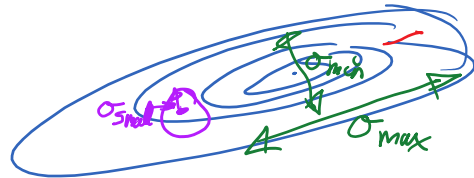
slow mixing example

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



*high prob of rejection

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



"small moves" \Rightarrow many steps needed

hence 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>