

today: MCMC
 review Markov chains
 Metropolis-Hastings

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 $X_t | X_{t-1}$ s.t. $X_t \xrightarrow{t \rightarrow \infty}$ in dist. to target dist. p

"stationary dist. of chain"

then we can approximate

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

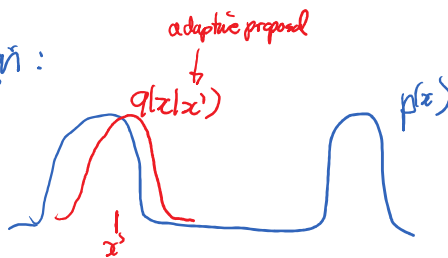
T_0 is called "burn-in" period \rightarrow depends on "mixing time" of Markov chain

⊗ no need to thin the samples [i.e. use Δt between samples to] get more independence

as this yields higher variance

\rightarrow better to use all samples after T_0 (unless it is too expensive)

Motivation:



before: samples were $X^{(i)} \text{ iid } q$

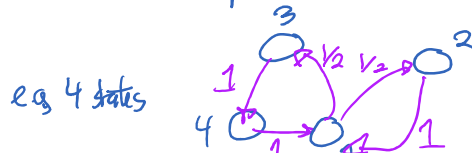
$$\text{MCMC } X^{(t)} | X^{(t-1)} \sim q(\cdot | X^{(t-1)})$$

\uparrow
 Markov transition prob

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

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

• there is also transition prob pt. of new i (probabilistic FSA) use one node per state



[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. $\mathbb{1}_k^T A = \mathbb{1}_k$

"left-stochastic matrix"

↑
vector of ones of size k

⊛ (as in HMM) suppose $P\{X_{t-1} = j\} = \pi_j$

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

A_{ij} π_j

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

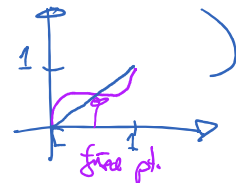
$$\Rightarrow \boxed{\pi_t = A^t \pi_0}$$

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

[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 every i to j

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

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

⊛ in order to converge fast, we need aperiodicity as well

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

aka. regular M.C. (finite state)
or ergodic M.C.

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

⊛ [note: a sufficient condition for an irreducible M.C. to be aperiodic is $\exists i$ s.t. $A_{ii} > 0$]

eg. of a regular M.C. $A = \begin{pmatrix} 0 & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & 0 & \frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} & 0 \end{pmatrix} = \frac{1}{2}(\mathbb{1}\mathbb{1}^T - I)$

$$A^2 = \frac{1}{4}(\underbrace{\mathbb{1}\mathbb{1}^T \mathbb{1}\mathbb{1}^T}_k - 2\mathbb{1}\mathbb{1}^T + I)$$

$$= \frac{1}{4}((k-2)\mathbb{1}\mathbb{1}^T + I) \quad \text{for } k \geq 3, \text{ this } > 0$$

[but for $k=2$, it is not aperiodic]

(for this for $k \neq 3$)

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

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

\leftarrow 2nd biggest e-value of A

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

* intuition (from linear algebra) [informal argument]

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

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

$U \rightarrow$ basis of e-vectors

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

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

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

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

let α_0 be coordinate of π_0 in U basis i.e. $\pi_0 = U \alpha_0$ ($\alpha_0 = U^T \pi_0$)

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

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

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

$t_1 \dots t_t \dots t_t \dots t_t \dots t_t$

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

$$= U \sum^u \alpha_0$$

$$A^t \pi_0 = (\alpha_0)_1^t u_1 + (\alpha_0)_2^t u_2 + \dots + (\alpha_0)_k^t u_k \quad \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_k^t \end{pmatrix}$$

$$A^t \pi_0 \xrightarrow{t \rightarrow \infty} (\alpha_0)_1 \frac{u_1}{\|u_1\|_2} \quad \mathbb{1}^T A^t \pi_0 = 1 \quad \forall t \Rightarrow \frac{(\alpha_0)_1}{\|u_1\|_2} = 1$$

$$\|A^t \pi_0 - \pi\|_1 = \|(\alpha_0)_2 \lambda_2^t u_2 + \dots\| \leq C |\lambda_2|^t$$

$|\lambda_2| = 1 - \epsilon_1$ $\epsilon_1 \triangleq 1 - |\lambda_2|$
 $|\lambda_2| \leq \exp(-\epsilon_1)$ $[1 - x \leq \exp(-x) \quad \forall x]$
 $|\lambda_2|^t \leq \exp(-t \epsilon_1)$
 $\frac{1}{\pi} \Rightarrow t = \frac{1}{1 - |\lambda_2|}$

* mixing time is often (usually) exponentially big!

15h16

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

one "easy way"

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

"detailed balance equation"

it means $P\{X_t = i, X_{t+1} = j\} = P\{X_t = j, X_{t+1} = i\}$

(when $P\{X_t = i\} = \pi_i$)

this is

\hookrightarrow sufficient condition (but not necessary)

$A \pi = \pi$

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

Metric-Hastings alg.

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

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

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

note $p(x)/p(x')$ does not depend on normalization of p

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

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

[this still a new sample]

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

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

does not depend
on normalization
of p

acceptance ratio to
satisfy detailed balance

M.H. alg.:

start at $x^{(0)}$

for $t=1, \dots$

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

• flip a biased coin with prob $q(x^{(t)}|x^{(t-1)})$ to be 1
acceptance ratio

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

o.w.

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

end for

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

\rightarrow like a noisy hill-climbing alg

[Metropolis alg.]

[verify as exercise that it satisfies detailed balance]

with $\pi = p$
 forget

⊗ for convergence: if M.H. chain is ergodic, then we converge to correct unique stat. dist. p

sufficient conditions \leftarrow 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 $q(x'|x) > 0$ for some x
and x'

* aside: it is still ok to change proposal with time

(inhomogeneous M.C.) $q_t(x'|x)$

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

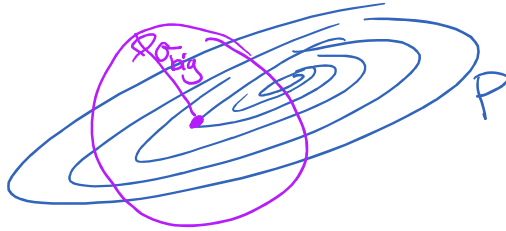
then convergence theory above will go through

[i.e. detailed balance, etc. will give
right stat. distribution]

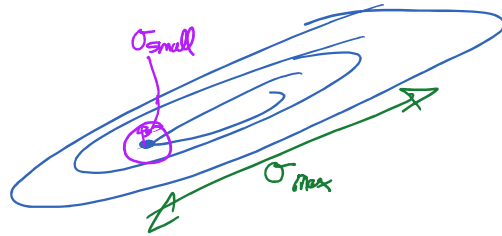
slow mixing example

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

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



* high prob of rejection



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

good book: Casella & Berger
"Monte Carlo Statistical methods"