

today: MCMC

Markov Chain Monte-Carlo:

idea: is to relax independence assumption between samples  
to allow adaptive proposal distributions

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

"stationary dist. of chain"

then, we can approximate

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

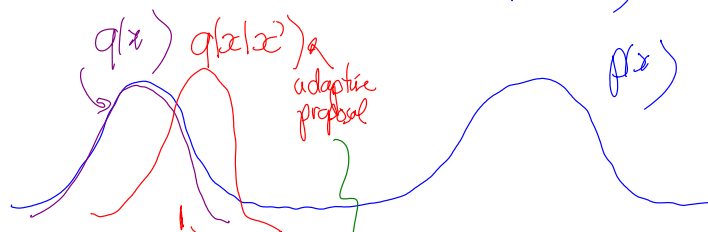
$T_0$  is called "burn-in" period  $\rightarrow$  depends on "mixing time" of Markov chain

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

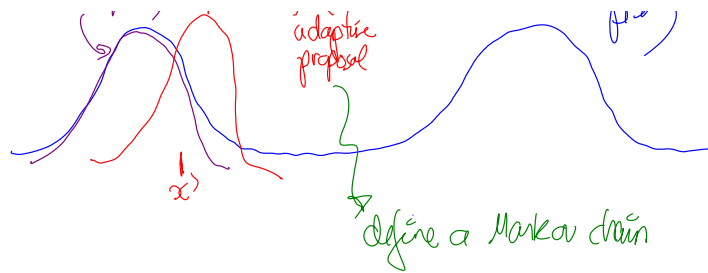
as this yields higher variance

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

Motivation:



Motivation:



before: samples were  $X^{(t)} \sim q$

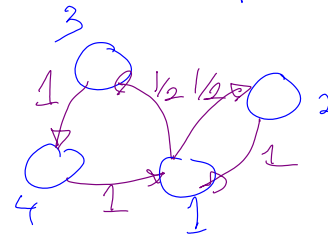
MCMC:  $X^{(t)} | X^{(t-1)} \sim q(\cdot | x^{(t-1)})$   
 $\uparrow$  Markov transition probability

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

• as a DGM  $X^{(0)} \rightarrow X^{(1)} \rightarrow X^{(2)} \rightarrow \dots$

• there is also the transition prob. point of view: use one node per state  
 (probabilistic finite state automaton) (FSA)

4 states example



[homogeneous M.C.]

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

$A$  is a  $k \times k$  matrix, s.t.  $\mathbf{1}^T A = \mathbf{1}^T$  (sum along a column is equal to 1)  
 vector of  $k$  ones

"left-stochastic matrix"

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

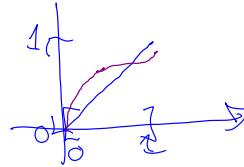
$$P\{X_t = i\} = \sum_j \underbrace{P(X_t = i | X_{t-1} = j)}_{A_{ij}} \underbrace{P(X_{t-1} = j)}_{\pi_j^0}$$

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

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

[note that  $\pi$  is a right  $e$ -vector of  $A$  with  $e$ -value of 1]

fact: every stochastic matrix has at least 1 stationary distribution  
(by Brouwer's fixed pt. thm.)



def:

irreducible MC  $\iff$  there is a positive probability "path" from every  $i \leftrightarrow j$

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

(by Perron-Frobenius thm.)  $\Rightarrow$  unique stationary dist. for irreducible MC.

$\rightarrow$  in order to converge to it, we need aperiodicity as well

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

aka. regular M.C.

(ie.  $(A^m)_{ij} > 0$ )

ergodic MC.

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

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

then  $\exists$  a unique stationary dist.  $\pi$

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

the speed of convergence is related to the mixing time  $\tau$  of the chain

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

2<sup>nd</sup> largest e-value of  $A$

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

\* intuition (from linear algebra)

suppose  $A$  is diagonalizable  $A = U \Sigma U^{-1}$  with  $\Sigma = \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_k \end{pmatrix}$

by Perron-Frobenius thm.

linear basis of e-vectors

can show that  $\lambda_1 = 1 > |\lambda_2| \geq \dots \geq |\lambda_k|$

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

take  $u_1 = \pi$  [e-value = 1]

let  $\alpha_0$  s.t.

$$\pi_0 = U \alpha_0$$

$$A^t \pi_0 = (u \leq u') \cdot (u \leq u') \cdot (u \leq u') \cdot \pi_0$$

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

$$\Sigma^t = \begin{pmatrix} 1 & 0 \\ \Lambda_2^t & 0 \\ 0 & \Lambda_k^t \end{pmatrix}$$

$$A^t \pi_0 = U \Sigma^t \alpha_0 = \left[ (\alpha_0)_1 \frac{u_1}{1} + (\alpha_0)_2 \Lambda_2^t u_2 + \dots + (\alpha_0)_k \Lambda_k^t u_k \right]$$

because  $1^T A^t \pi_0 = 1$  for all  $t \Rightarrow (\alpha_0)_k = 1$

$$\|A^t \pi_0 - \pi\|_1 \leq C |\Lambda_2|^t$$

first e-gap  
 $|\Lambda_2| = 1 - \epsilon_1 \quad \epsilon_1 \triangleq 1 - |\Lambda_2|$

$$|\Lambda_2| \leq \exp(-\epsilon_1)$$

$$|\Lambda_2|^t \leq \exp(-\epsilon_1 t)$$

$$\Rightarrow \tau = \frac{1}{1 - |\Lambda_2|}$$

mixing time can be exponentially big sometimes

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

one "easy way"

reversible M.C.

iff  $\exists$  dist.  $\pi$  s.t.  $A_{ij} \pi_j = A_{ji} \pi_i \quad \forall (i, j)$

"detailed balance equation"

it means that  $P\{X_t=i, X_{t-1}=j\} = P\{X_t=j, X_{t-1}=i\}$   
 [when  $P\{X_{t+1}=i\} = \pi_i$ ]

sufficient condition (but not necessary)

for  $A\pi = \pi$

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

Metropolis-Hastings algorithm

→ construct a MC. with stationary dist.  $p(x)$  [our target]  
 (assume  $p(x) > 0$ )

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

[i.e. if in state  $x$ , we sample  $x'|x \sim q(x'|x)$ ]

accept new state  $x'$  with prob.  
 if reject → stay in same state  $x$

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

depends only on  $p(x')/p(x)$   
 → no need to normalise

acceptance ratio to satisfy detailed balance

[this is still a new sample]  
 (vs. rejection sampling where only "accepted" states are sample)

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

if accept (coin=1)

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

O.w.

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

end

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

→ like noisy hill-climbing alg.

Let's verify detailed balance

$$A_{ij}\pi_j = A_{ji}\pi_i$$

•  $A_{ij}\pi_j = A_{ji}\pi_i$  trivially

• to have  $A_{ij}\pi_j = A_{ji}\pi_i$

need.  $q(i|j)a(i|j)p_j = q(j|i)a(j|i)p_i$

target dist.

$$\Rightarrow \text{want } \boxed{\frac{a(i|j)}{a(j|i)} = \frac{q(j|i)p_i}{q(i|j)p_j}}$$

to finish: use the max  $\{1, \frac{p_i}{p_j}\}$  to look at cases //

for convergence: if MH chain is ergodic, then we converge to correct unique stationary dist.  $p$

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

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

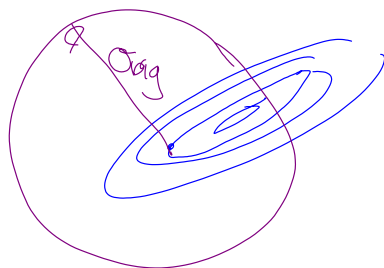
(\*) aside: it is still ok to change proposal with time  
(inhomogeneous MC.)  $q_t(x'|x)$

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

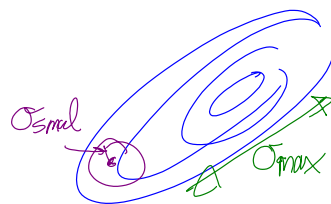
then convergence theory above will go through  
(i.e. detailed balance etc... will give right stationary dist.)

slow mixing example:

suppose  $p$  is multivariate normal  $\& q(x'|x) = N(x'|x, \sigma^2 I)$



\* high prob. of rejection



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

good book: Casella & Berger Monte Carlo Statistical Methods

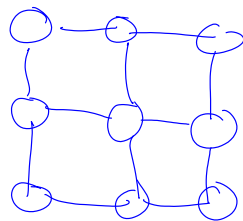
Gibbs sampling algorithm

MH with clever choice of proposal  $q_t(x'|x)$

examples of applications.  $\hookrightarrow$  UGM:  $\tilde{p}(x) = \prod_i \psi_i(x_i)$



UGM's



difficult conditional  
in DBM

$$\tilde{p}(x) = p(x, \bar{x}_E)$$

$$\propto p(x | \bar{x}_E)$$

cyclic Gibbs sampling alg.: nodes  $i=1, \dots, n$

start at some  $x^{(t)}$

for  $t=1, \dots, \infty$ ,

- pick  $i = (t \bmod n) + 1$

- sample  $x_i^{(t)} \sim p(x_i = \cdot \mid x_{-i} = x_{-i}^{(t-1)})$

- set  $x_j^{(t)} = x_j^{(t-1)}$  for  $j \neq i$

true condition as proposal

end