

Inference with Markov Chains

Outline

- I. Gibbs sampling
- II. Markov chains
- III. Metropolis-Hastings sampling

* Figures are either from the [textbook site](#).

I. Gibbs Sampling

- ◆ A Markov Chain Monte Carlo (MCMC) algorithm
 - specifies a value for every variable at the current state.
 - generates a next state by making random changes to the current state.
- ◆ *Markov chain* is a random process that generates a sequence of states.

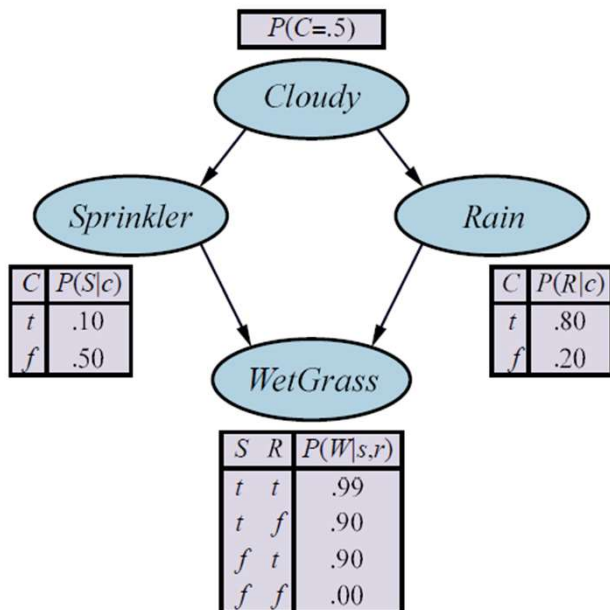
Gibbs sampling (well suited for Bayes nets) is an MCMC algorithm that

- starts with an arbitrary state,
- fix evidence variables at their observed values, and
- generates a next state by randomly sampling a value for a nonevidence variable X_i chosen according to probability distribution $\rho(i)$.

X_i is independent of all the variables outside of its *Markov blanket* (consisting its parents, children, and children's other parents).

Example of Gibbs Sampling

Gibbs sampling for X_i is conditioned on the current values of the variables in its Markov blanket.



Query $P(\text{Rain} \mid \text{Sprinkler} = \text{true}, \text{WetGrass} = \text{true})$

randomly generated values for
nonevidence variables *Cloudy* and *Rain*

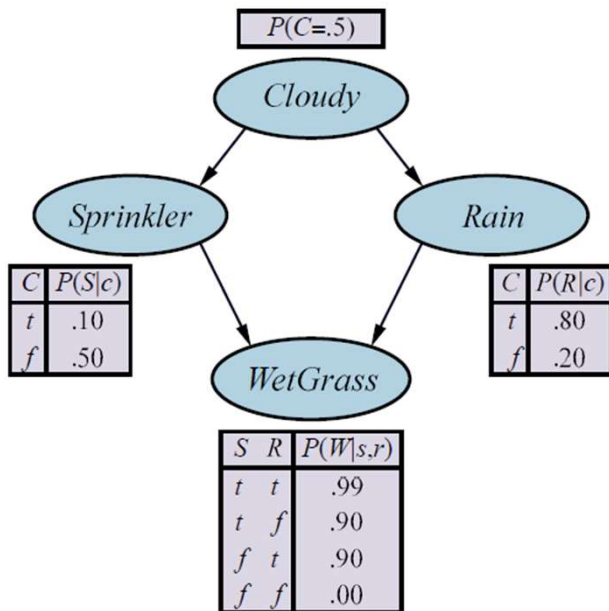
- Initial state [*true*, *true*, *false*, *true*]

evidence variables *Sprinkler* and
WetGrass fixed to their observed values

Order: *Cloudy*, *Sprinkler*, *Rain*, *WetGrass*

Example (cont'd)

Query $P(\text{Rain} \mid \text{Sprinkler} = \text{true}, \text{WetGrass} = \text{true})$



Order: *Cloudy, Sprinkler, Rain, WetGrass*

$[\text{true}, \text{true}, \text{false}, \text{true}]$
(initial state)

$\rightarrow [\text{false}, \text{true}, \text{true}, \text{true}]$
(new current state)

- Non-evidence variables are then sampled in random order following some probability distribution $\rho(i)$.

- Cloudy is chosen and sampled given the current values of its Markov blanket $\{\text{Sprinkler}, \text{Rain}\}$.

- Sampling distribution:

$$P(\text{Cloudy} \mid \text{Sprinkler} = \text{true}, \text{Rain} = \text{false})$$

- Sampling result: *Cloudy = false*.

- Rain is chosen next and sampled given the current values of its Markov blanket $\{\text{Cloudy}, \text{Sprinkler}, \text{WetGrass}\}$.

- Sampling distribution:

$$P(\text{Rain} \mid \text{Cloudy} = \text{false}, \text{Sprinkler} = \text{true}, \text{WetGrass} = \text{true})$$

- Sampling result: *Rain = true*.

How to calculate?

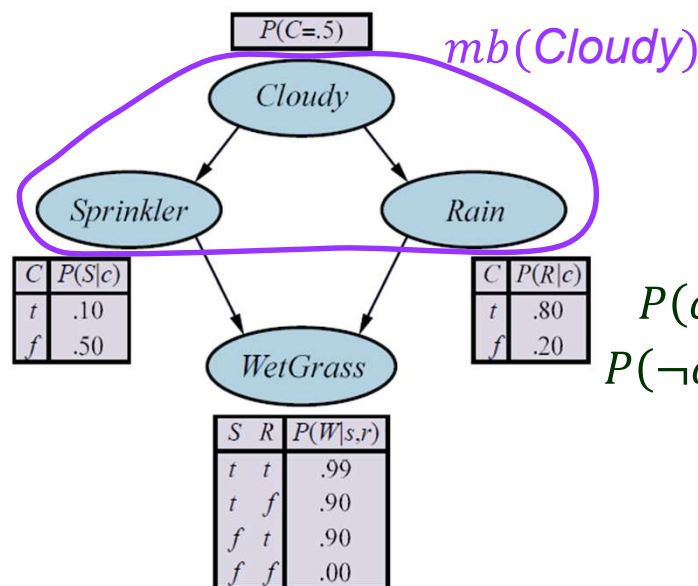
Markov Blanket Distribution

$MB(X_i)$: variables in the Markov blanket of X_i .

$mb(X_i)$: values of the variables in $MB(X_i)$.

$P(X_i \mid mb(X_i))$ is determined as follows:

$$P(x_i \mid mb(X_i)) = \alpha P(x_i \mid \text{parents}(X_i)) \prod_{Y_j \in \text{Children}(X_i)} P(y_j \mid \text{parents}(Y_j))$$



♣ Sampling distribution:

$$P(Cloudy \mid \text{Sprinkler} = \text{true}, \text{Rain} = \text{false})$$

$$P(c \mid s, \neg r) = \alpha P(c) P(s \mid c) P(\neg r \mid c) = \alpha 0.5 \cdot 0.1 \cdot 0.2$$

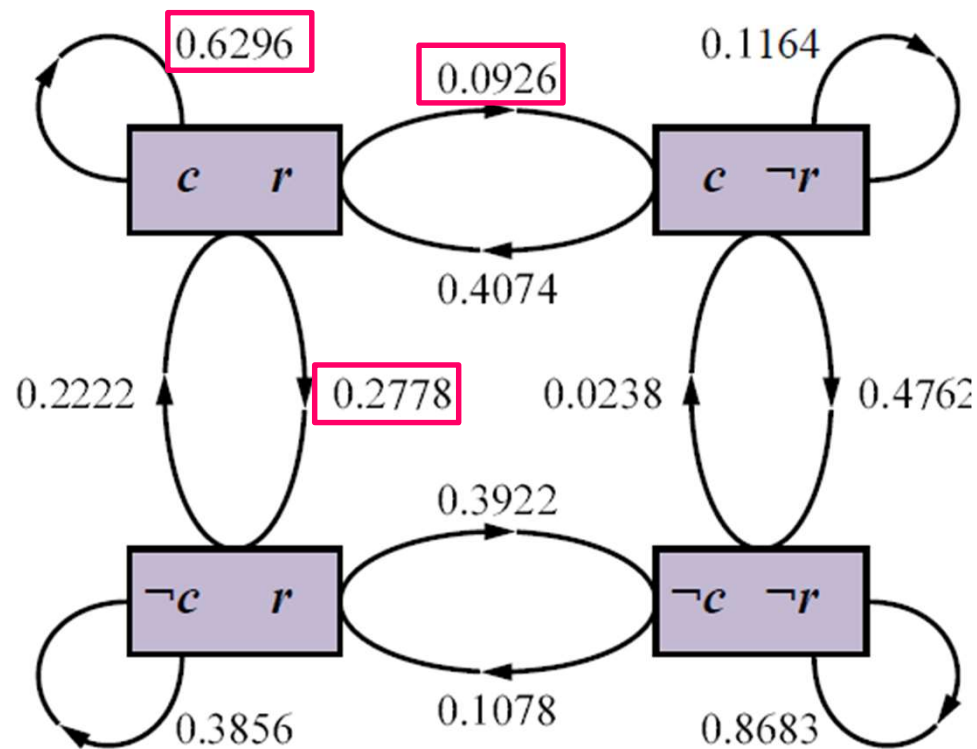
$$P(\neg c \mid s, \neg r) = \alpha P(\neg c) P(s \mid \neg c) P(\neg r \mid \neg c) = \alpha 0.5 \cdot 0.5 \cdot 0.8$$

$$P(C \mid s, \neg r) = \alpha \langle 0.001, 0.020 \rangle \approx \langle 0.048, 0.952 \rangle$$

II. Markov Chain

Query $P(\text{Rain} \mid \text{Sprinkler} = \text{true}, \text{WetGrass} = \text{true})$

A **state** need only include all nonevidence variables.



Markov chain from uniform choice of the two nonevidence variables ($p(\text{Cloudy}) = p(\text{Rain}) = 0.5$)

Probabilities with all the outgoing links of each node sum to 1, e.g., $0.6296 + 0.0926 + 0.2778 = 1$.

♦ Gibbs sampling simply wanders around in the graph, following links with probabilities.

♦ Every state visited is a sample that contributes to the estimate for the query variable *Rain*.

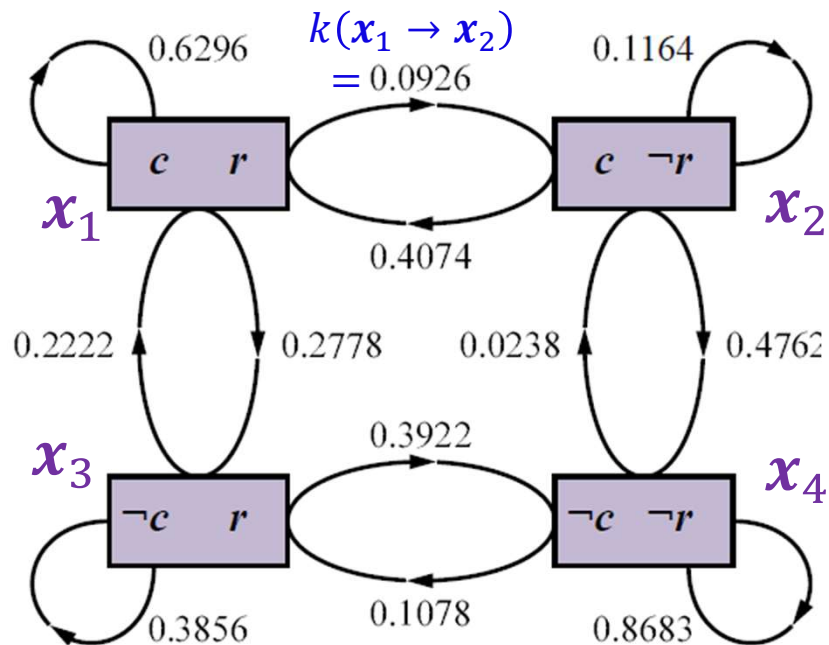
If the process visits 20 states with *Rain* = *true* and 60 states with *Rain* = *false*, then the answer to the query is $\alpha\langle 20, 60 \rangle = \langle 0.25, 0.75 \rangle$.

Analysis of Markov Chains

Why does Gibbs sampling work? Or, why does its estimates converge to correct values in the limit?

Transition kernel k assigns a probability $k(x \rightarrow x')$ to a transition from state x to state x' .

$\pi_t(x)$: probability that the system is in state x after t transitions



$$\pi_{t+1}(x') = \sum_x \pi_t(x) k(x \rightarrow x')$$

$$\pi_{t+1}(x_2) = 0.0926 \pi_t(x_1) + 0.1164 \pi_t(x_2) + 0.0238 \pi_t(x_4)$$

The chain has reached its stationary distribution if $\pi_{t+1}(x) = \pi_t(x)$ for all x .

Stationary Distribution

A distribution π of the Markov chain is *stationary* if

$$\pi(x') = \sum_x \pi(x)k(x \rightarrow x') \quad \text{for all } x, x'.$$

Such a probability distribution remains unchanged in the Markov chain as time progresses.

A kernel k is *ergodic* if every state is reachable from every other state and there exists no strictly periodic cycles.

There exists exactly one stationary distribution for every ergodic kernel of the Markov chain.

Achieving a Stationary Distribution

In a stationary distribution π , the expected “outflow” from each state is equal to the expected “inflow” from all the other states.

“population”

$$\underbrace{\pi(x')}_{\text{expected “outflow”}} = \sum_x \underbrace{\pi(x)k(x \rightarrow x')}_{\text{expected “inflow”}}$$

A *detailed balance* k with π is a distribution that satisfies

$$\pi(x)k(x \rightarrow x') = \pi(x')k(x' \rightarrow x) \quad \text{for all } x, x'.$$

The detailed balance k makes $\pi(x)$ stationary because

$$\begin{aligned} \sum_x \pi(x)k(x \rightarrow x') &= \sum_x \pi(x')k(x' \rightarrow x) \\ &= \pi(x') \sum_x k(x' \rightarrow x) = 1 \end{aligned}$$



Correctness of Gibbs Sampling

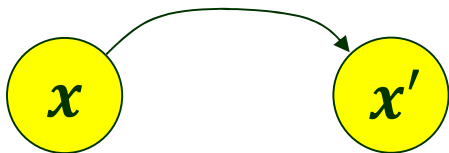
The stationary distribution of the Gibbs sampling process is exactly the posterior distribution for the nonevidence variable conditioned on the evidence.

- In Gibbs sampling, a variable X_i is chosen and sampled conditionally on
 - ♦ the current values of all the other variables,
 - ♦ equivalently, when sampling a Bayes net, the variable's Markov blanket.

Transition Kernel for Gibbs Sampling

\bar{X}_i : variables except X_i and evidence variables.

\bar{x}_i : their values.



Case 1. The states x and x' differ in ≥ 2 variables. Since Gibbs sampling changes only one variable, we set

$$k(x \rightarrow x') = 0$$

Case 2. The states x and x' differ in the value of exactly one variable X_i , which changes from x_i to x'_i . That is, $x = (x_i, \bar{x}_i)$ and $x' = (x'_i, \bar{x}_i)$.

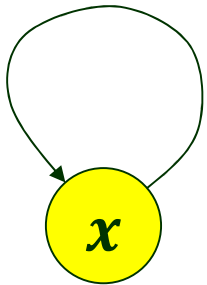
$$k(x \rightarrow x') = k((x_i, \bar{x}_i) \rightarrow (x'_i, \bar{x}_i)) = \rho(i) P(x'_i | \bar{x}_i)$$

probability of choosing X_i

The transition probability is the product of the probability of selecting the variable X_i (out of all the nonevidence variables) with the probability of selecting x'_i (out of all the values of X_i).

Completing the Definition

Case 3. The states are the same $x = x'$. Any variable could be chosen but then the sampling process reproduce the current value of the variable.



$$k(x \rightarrow x') = \sum_i \rho(i) k((x_i, \bar{x}_i) \rightarrow (x'_i, \bar{x}_i)) = \sum_i \rho(i) P(x_i | \bar{x}_i)$$

Correctness of Gibbs Sampling

Theorem The previously defined kernel $k(\mathbf{x}' \rightarrow \mathbf{x})$ for Gibbs sampling has a stationary distribution equal to $P(\mathbf{x} | \mathbf{e})$, the true posterior distribution on the nonevidence variables.

Proof It suffices to show that, with $\pi(\mathbf{x}) = P(\mathbf{x} | \mathbf{e})$, the following condition for k in detailed balance is satisfied:

$$\pi(\mathbf{x})k(\mathbf{x} \rightarrow \mathbf{x}') = \pi(\mathbf{x}')k(\mathbf{x}' \rightarrow \mathbf{x}) \quad \text{for all states } \mathbf{x}, \mathbf{x}'.$$

Then it follows that k implies the stationarity distribution $P(\mathbf{x} | \mathbf{e})$.

- ◆ In the first and third cases, where \mathbf{x} and \mathbf{x}' differ in ≥ 2 variables and $\mathbf{x} = \mathbf{x}'$, respectively, detailed balance can be easily shown to be satisfied.
- ◆ In the second case, where \mathbf{x} and \mathbf{x}' differ in one variable x_i , we have

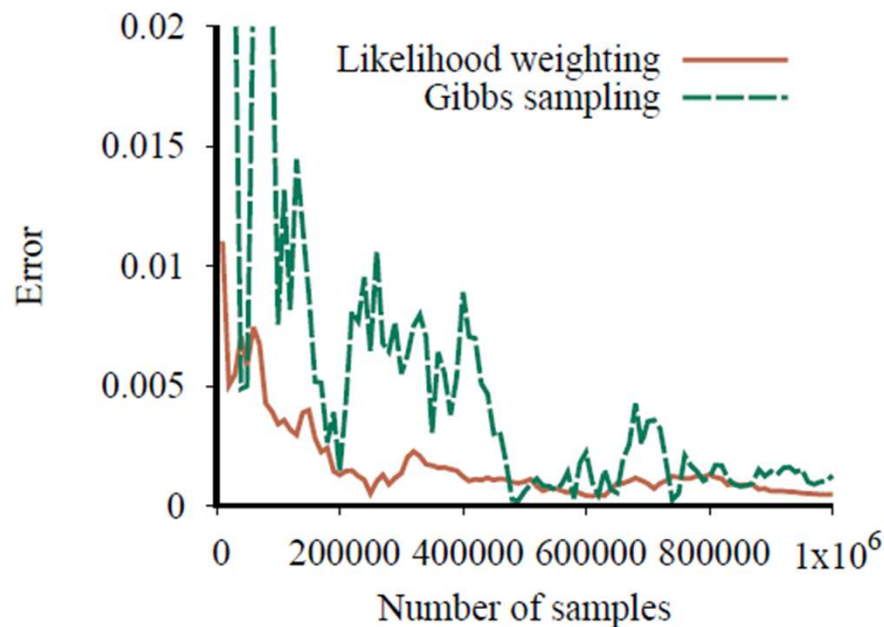
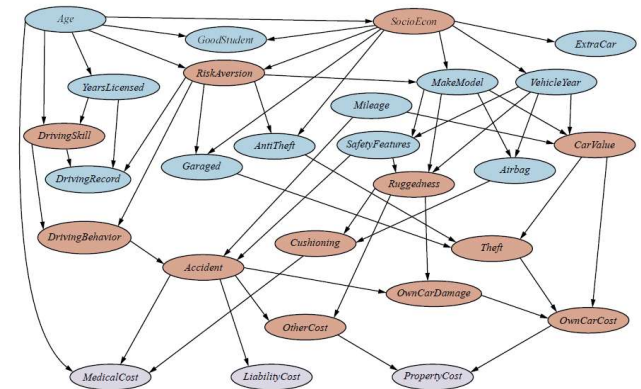
$$\begin{aligned} \pi(\mathbf{x})k(\mathbf{x} \rightarrow \mathbf{x}') &= P(\mathbf{x} | \mathbf{e}) \rho(i)P(x'_i | \bar{\mathbf{x}}_i, \mathbf{e}) = \rho(i) P(x_i, \bar{\mathbf{x}}_i | \mathbf{e})P(x'_i | \bar{\mathbf{x}}_i, \mathbf{e}) \\ &= \rho(i)P(x_i | \bar{\mathbf{x}}_i, \mathbf{e})P(\bar{\mathbf{x}}_i | \mathbf{e})P(x'_i | \bar{\mathbf{x}}_i, \mathbf{e}) \\ &= \rho(i)P(x_i | \bar{\mathbf{x}}_i, \mathbf{e}) P(x'_i, \bar{\mathbf{x}}_i | \mathbf{e}) = \rho(i)P(x'_i, \bar{\mathbf{x}}_i | \mathbf{e})P(x_i | \bar{\mathbf{x}}_i, \mathbf{e}) \\ &= \pi(\mathbf{x}')k(\mathbf{x}' \rightarrow \mathbf{x}) \end{aligned}$$



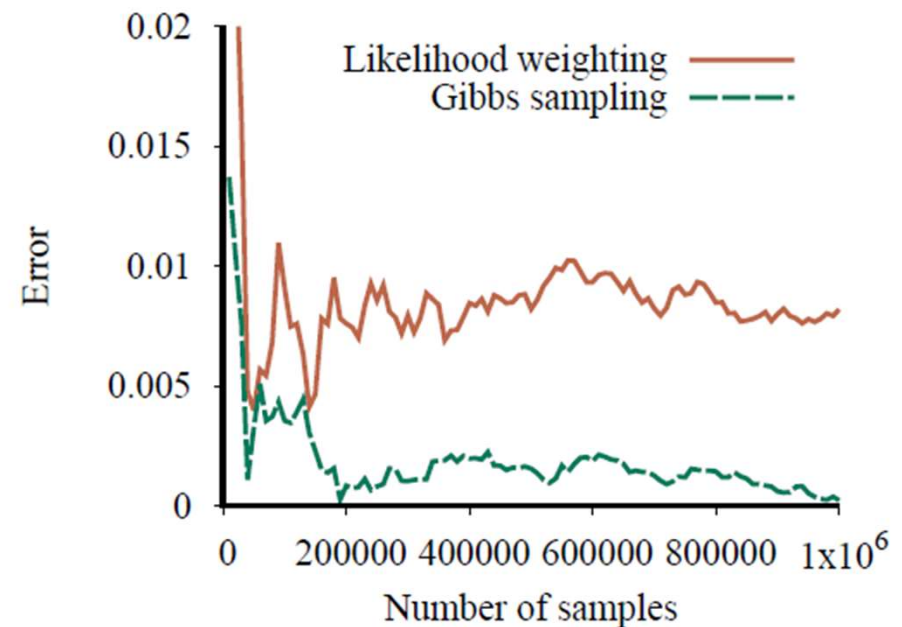
Performance of Gibbs Sampling

Gibbs sampling is expected to outperform likelihood weighting when evidence is downstream.

On the car insurance network:



Query on *PropertyCost*



Query on *Age* (with output observed)

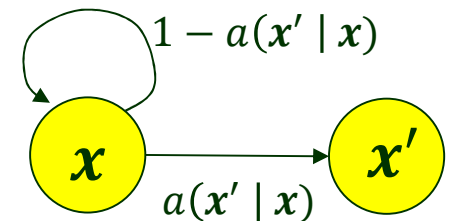
III. Metropolis-Hastings (MH) Sampling

- ♦ The most broadly applicable Markov chain Monte Carlo algorithm.
- ♦ MH generates samples x according to a target probability distribution $\pi(x)$ (in a BN, $\pi(x) = P(x | e)$).

The transition kernel $k(x \rightarrow x')$ is defined as follows:

- ♣ At the current state x , sample a new state x' from a **proposal distribution** $q(x' | x)$.
- ♣ Accept or reject x' according to the **acceptance probability**:

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



- ♣ The state transitions from x to x' in the case of acceptance, and stays at x in the case of rejection.

Proposal Distribution for MH

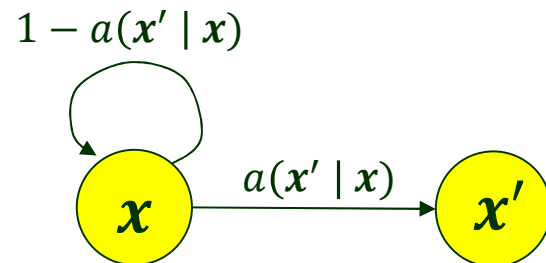
The *proposal distribution* $q(\mathbf{x}' | \mathbf{x})$ is responsible for proposing a new state \mathbf{x}' .

Example $q(\mathbf{x}' | \mathbf{x})$ could be defined as follows:

- With probability 0.95, perform a Gibbs sampling step to generate \mathbf{x}' .
- With probability 0.05, use likelihood weighting to generate \mathbf{x}' .
- ◆ This proposal distribution causes MH to do about 19 steps of Gibbs sampling and then generates a new state from scratch.
- ◆ It gets around the problem of Gibbs sampling getting stuck in one part of the state space.

Convergence of MH

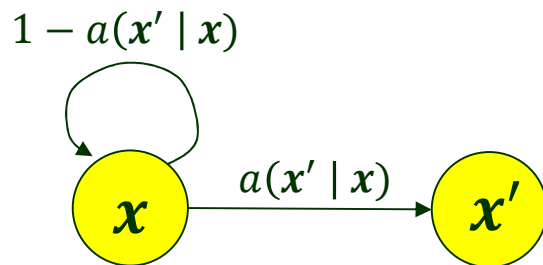
MH converges to the correct stationary distribution for any proposal distribution $q(x' | x)$, provided it results in an ergodic transition kernel.



- ◆ The self-loop with $x = x'$ automatically satisfies the detailed balance condition:

$$\pi(x)k(x \rightarrow x') = \pi(x')k(x' \rightarrow x)$$

Convergence of MH (cont'd)



- ◆ In the case $x \neq x'$, the transition can occur only if the proposal of x' is accepted.

$$k(x \rightarrow x') = q(x' | x)a(x' | x)$$

We can show that the flow from x to x' equals that from x' to x (i.e., $k(x \rightarrow x')$ is in detailed balance with $\pi(x)$) as follows:

$$\begin{aligned}\pi(x)k(x \rightarrow x') &= \pi(x)q(x' | x)a(x' | x) \\ &= \pi(x)q(x' | x) \min\left(1, \frac{\pi(x')q(x | x')}{\pi(x)q(x' | x)}\right) \\ &= \min(\pi(x)q(x' | x), \pi(x')q(x | x')) \\ &= \pi(x')q(x | x') \min\left(\frac{\pi(x)q(x' | x)}{\pi(x')q(x | x')}, 1\right) \\ &= \pi(x')k(x' \rightarrow x)\end{aligned}$$