

Lecture 24: 18 April, 2023

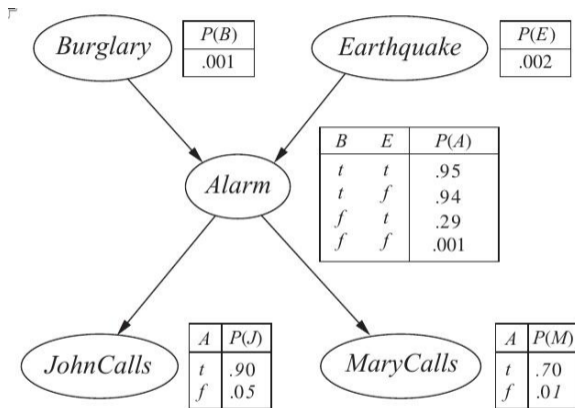
Madhavan Mukund

<https://www.cmi.ac.in/~madhavan>

Data Mining and Machine Learning
January–April 2023

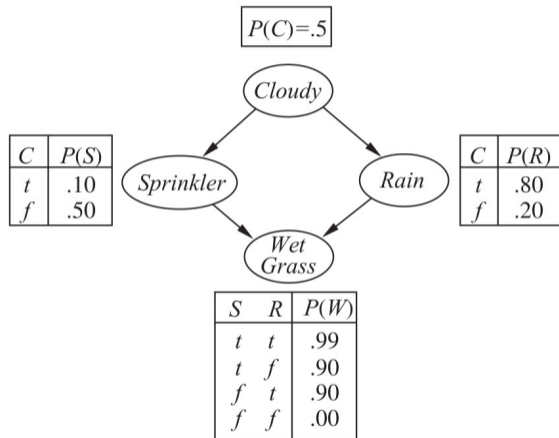
Approximate inference

- Exact inference is NP-complete
- Generate random samples, count to estimate probabilities
- Respect conditional probabilities — generate in topological order
- Suppose we are interested in $P(b | j, m)$
- Samples with $\neg j$ or $\neg m$ are useless
- Can we sample more efficiently?



Rejection sampling

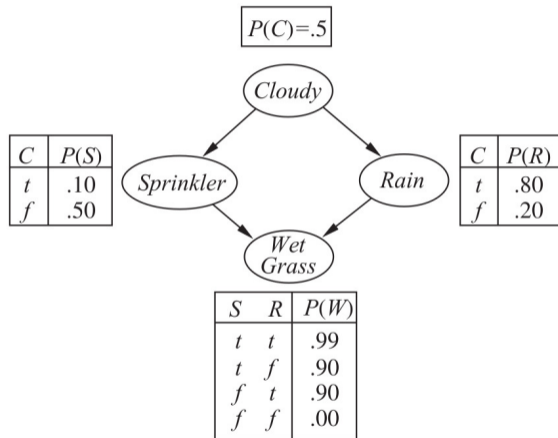
- $P(\text{Rain} \mid \text{Cloudy}, \text{Wet Grass})$
- If we start with $\neg \text{Cloudy}$, sample is useless
- Immediately stop and reject this sample — **rejection sampling**
- General problem with low probability situation — many samples are rejected



Likelihood weighted sampling

- $P(\text{Rain} \mid \text{Cloudy}, \text{Wet Grass})$
- Fix **evidence** *Cloudy*, *Wet Grass* true
- Then generate the other variables
- Compute likelihood of evidence
- Samples s_1, s_2, \dots, s_N with weights w_1, w_2, \dots, w_N

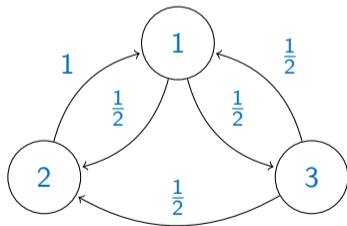
- $$P(r \mid c, w) = \frac{\sum_{s_i \text{ has rain}} w_i}{\sum_{1 \leq j \leq N} w_j}$$



Approximate inference using Markov chains

Markov chains

- Finite set of states, with transition probabilities between states
- For us, a state will be an assignment of values to variables
- A three state Markov Chain



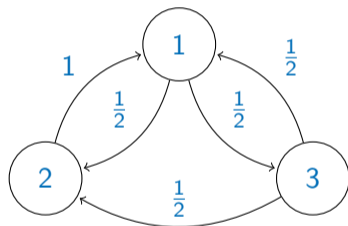
- Represent using a **transition matrix** — stochastic

$$A = \begin{bmatrix} 0 & \frac{1}{2} & \frac{1}{2} \\ 1 & 0 & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 \end{bmatrix}$$

- $P[j]$ is probability of being in state j

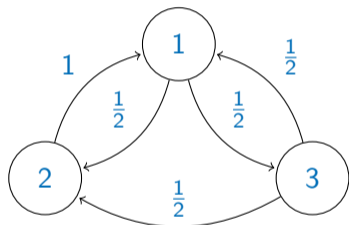
Ergodicity

- Markov chain A is **ergodic** if there is some t_0 such that for every P , for all $t > t_0$, for every j , $(P^\top A^t)[j] > 0$.
- Ergodic Markov chain has a stationary distribution π^* , $(\pi^*)^\top A = \pi^*$
- For *any* starting distribution P , $\lim_{t \rightarrow \infty} P^\top A^t = \pi^*$
- Stationary distribution represents fraction of visits to each state in a long enough execution
- Sufficient conditions for ergodicity
 - Irreducible (strong connected)
 - Aperiodic (paths of all lengths between states)



Approximate inference using Markov chains

- Bayesian network has variables V_1, V_2, \dots, V_n
- Each assignment of values to the variables is a state
- Set up a Markov chain based on these states
- Stationary distribution should assign to state s the probability $P(s)$ in the Bayesian network
- How to reverse engineer the transition probabilities to achieve this?



Reversible Markov chains

- Ergodic Markov chain with stationary distribution π

Reversible Markov chains

- Ergodic Markov chain with stationary distribution π^*
- Transition matrix A , write p_{jk} for $A[j][k]$
 - Probability of transition from state j to state k

Reversible Markov chains

- Ergodic Markov chain with stationary distribution π^*
- Transition matrix A , write p_{jk} for $A[j][k]$
 - Probability of transition from state j to state k
- **Reversibility** — in steady state, probability of going from j to k should equal probability of going from k to j

Reversible Markov chains

- Ergodic Markov chain with stationary distribution π^*
- Transition matrix A , write p_{jk} for $A[j][k]$
 - Probability of transition from state j to state k
- **Reversibility** — in steady state, probability of going from j to k should equal probability of going from k to j
- Given an evolution $x_1 x_2 \dots$, for large n , $P[x_n = k \mid x_{n-1} = j] = P[x_{n-1} = j \mid x_n = k]$

$$\leftarrow A(j,k) = p_{jk}$$

$j \rightarrow k$ $j \leftarrow k$

Reversible Markov chains

- Ergodic Markov chain with stationary distribution π^*
- Transition matrix A , write p_{jk} for $A[j][k]$
 - Probability of transition from state j to state k
- **Reversibility** — in steady state, probability of going from j to k should equal probability of going from k to j
- Given an evolution $x_1 x_2 \dots$, for large n , $P[x_n = k \mid x_{n-1} = j] = P[x_{n-1} = j \mid x_n = k]$
- $P[x_{n-1} = j \mid x_n = k] = P[x_n = k \mid x_{n-1} = j] \cdot \frac{P[x_{n-1} = j]}{P[x_n = k]}$
 $\pi_k^* = \pi_j$

Reversible Markov chains

- Ergodic Markov chain with stationary distribution π^*
- Transition matrix A , write p_{jk} for $A[j][k]$
 - Probability of transition from state j to state k
- **Reversibility** — in steady state, probability of going from j to k should equal probability of going from k to j
- Given an evolution $x_1 x_2 \dots$, for large n , $P[x_n = k \mid x_{n-1} = j] = P[x_{n-1} = j \mid x_n = k]$
- $P[x_{n-1} = j \mid x_n = k] = P[x_n = k \mid x_{n-1} = j] \cdot \frac{\pi_j}{\pi_k}$, in steady state

$\leftarrow P_{jk}$

$\pi^k \leftrightarrow \pi$

Reversible Markov chains

- Ergodic Markov chain with stationary distribution π^*
 - Transition matrix A , write p_{jk} for $A[j][k]$
 - Probability of transition from state j to state k
 - **Reversibility** — in steady state, probability of going from j to k should equal probability of going from k to j
 - Given an evolution $x_1 x_2 \dots$, for large n , $P[x_n = k \mid x_{n-1} = j] = P[x_{n-1} = j \mid x_n = k]$
 - $P[x_{n-1} = j \mid x_n = k] = P[x_n = k \mid x_{n-1} = j] \cdot \frac{\pi_j}{\pi_k}$, in steady state
 - $p_{kj} = p_{jk} \frac{\pi_j}{\pi_k}$
-

Reversible Markov chains

- Ergodic Markov chain with stationary distribution π^*
- Transition matrix A , write p_{jk} for $A[j][k]$
 - Probability of transition from state j to state k
- **Reversibility** — in steady state, probability of going from j to k should equal probability of going from k to j
- Given an evolution $x_1 x_2 \dots$, for large n , $P[x_n = k \mid x_{n-1} = j] = P[x_{n-1} = j \mid x_n = k]$
- $P[x_{n-1} = j \mid x_n = k] = P[x_n = k \mid x_{n-1} = j] \cdot \frac{\pi_j}{\pi_k}$, in steady state
- $p_{kj} = p_{jk} \frac{\pi_j}{\pi_k}$
- $\pi_j \cdot p_{jk} = \pi_k \cdot p_{kj}$

Alleged explain

— Treats this as the defn we want to use

- Ergodic Markov chain

Reversible Markov chains

- Ergodic Markov chain
- Suppose $a^T = (a_1, a_2, \dots, a_n)$ satisfies reversibility condition for all j, k
 - $a_j \cdot p_{jk} = a_k \cdot p_{kj}$

$$\pi_j p_{jk} = \pi_k p_{kj}$$

Reversible Markov chains

- Ergodic Markov chain
- Suppose $a^T = (a_1, a_2, \dots, a_n)$ satisfies reversibility condition for all j, k
 - $a_j \cdot p_{jk} = a_k \cdot p_{kj}$
- $$\sum_k a_j \cdot p_{jk} = \sum_k a_k \cdot p_{kj}$$

Reversible Markov chains

- Ergodic Markov chain
- Suppose $a^T = (a_1, a_2, \dots, a_n)$ satisfies reversibility condition for all j, k

- $a_j \cdot p_{jk} = a_k \cdot p_{kj}$ — for every j & k

- $\sum_k a_j \cdot p_{jk} = \sum_k a_k \cdot p_{kj}$

- $a_j \sum_k p_{jk} = \sum_k a_k \cdot p_{kj}$

↓
row j in A

$a_{j1} \cdot p_{j1} = a_1 \cdot p_{1j}$
 $a_{j2} \cdot p_{j2}$

↓

↓

Reversible Markov chains

- Ergodic Markov chain
- Suppose $a^T = (a_1, a_2, \dots, a_n)$ satisfies reversibility condition for all j, k
 - $a_j \cdot p_{jk} = a_k \cdot p_{kj}$
- $\sum_k a_j \cdot p_{jk} = \sum_k a_k \cdot p_{kj}$
- $a_j \sum_k p_{jk} = \sum_k a_k \cdot p_{kj}$
- $a_j \cdot 1 = \sum_k a_k \cdot p_{kj}$

Reversible Markov chains

$$a_j = a_j p_{1j} + a_2 p_{2j} + \dots + a_n p_{nj}$$

- Ergodic Markov chain
- Suppose $a^T = (a_1, a_2, \dots, a_n)$ satisfies reversibility condition for all j, k

- $a_j \cdot p_{jk} = a_k \cdot p_{kj}$

- $\sum_k a_j \cdot p_{jk} = \sum_k a_k \cdot p_{kj}$

- $a_j \sum_k p_{jk} = \sum_k a_k \cdot p_{kj}$

- $a_j \cdot 1 = \sum_k a_k \cdot p_{kj}$

- $a^T = a^T A$, so a^T is the stationary distribution of A

Handwritten diagram illustrating the reversibility condition. It shows a red bracketed expression $[\dots a_j \dots] =$ followed by a blue double-headed arrow $[\leftarrow \rightarrow]$ and a blue bracketed expression $[a^T]$ below it. To the right is a red bracketed expression $[]$ containing a vertical blue double-headed arrow and a blue j above it.

- State of a Bayesian network is a valuation of variables (V_1, V_2, \dots, V_n)

Gibbs sampling

- State of a Bayesian network is a valuation of variables $(V_1, V_2, \dots, V_n$
- Move probabilistically from $s_j = (x_1, x_2, \dots, x_n)$ to $s_k = (y_1, y_2, \dots, y_n)$

Gibbs sampling

- State of a Bayesian network is a valuation of variables (V_1, V_2, \dots, V_n)
- Move probabilistically from $s_j = (x_1, x_2, \dots, x_n)$ to $s_k = (y_1, y_2, \dots, y_n)$
- Allow such a move only when s_j, s_k differ at exactly one position
 - $s_j = (x_1, x_2, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n)$
 - $s_k = (x_1, x_2, \dots, x_{i-1}, y_i, x_{i+1}, \dots, x_n)$

Gibbs sampling

- State of a Bayesian network is a valuation of variables (V_1, V_2, \dots, V_n)
- Move probabilistically from $s_j = (x_1, x_2, \dots, x_n)$ to $s_k = (y_1, y_2, \dots, y_n)$
- Allow such a move only when s_j, s_k differ at exactly one position
 - $s_j = (x_1, x_2, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n)$
 - $s_k = (x_1, x_2, \dots, x_{i-1}, y_i, x_{i+1}, \dots, x_n)$
- Sampling algorithm
 - Current state is $s_j = (x_1, x_2, \dots, x_n)$
 - Choose i uniformly in $[1, n]$
 - Resample x_i given current values $(x_1, x_2, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$

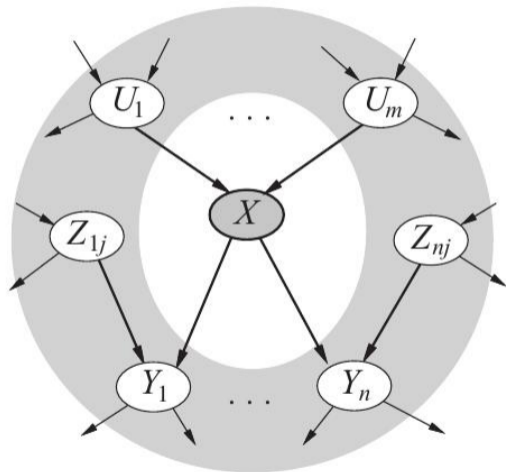
Gibbs sampling

- State of a Bayesian network is a valuation of variables (V_1, V_2, \dots, V_n)
- Move probabilistically from $s_j = (x_1, x_2, \dots, x_n)$ to $s_k = (y_1, y_2, \dots, y_n)$
- Allow such a move only when s_j, s_k differ at exactly one position
 - $s_j = (x_1, x_2, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n)$
 - $s_k = (x_1, x_2, \dots, x_{i-1}, y_i, x_{i+1}, \dots, x_n)$
- Sampling algorithm
 - Current state is $s_j = (x_1, x_2, \dots, x_n)$
 - Choose i uniformly in $[1, n]$
 - Resample x_i given current values $(x_1, x_2, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$
- Need to compute $P[y_i \mid x_1, x_2, \dots, x_{i-1}, x_{i+1}, \dots, x_n]$

could be x_i

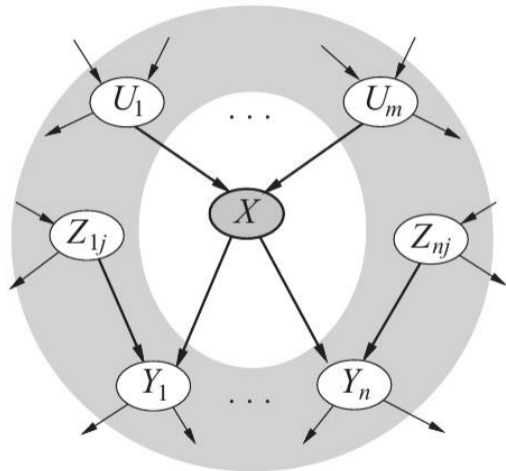
Markov blanket

- Recall $MB(X)$ — Markov blanket of X
 - $Parents(X)$
 - $Children(X)$
 - $Parents\ of\ Children(X)$
- $X \perp \neg MB(X) \mid MB(X)$



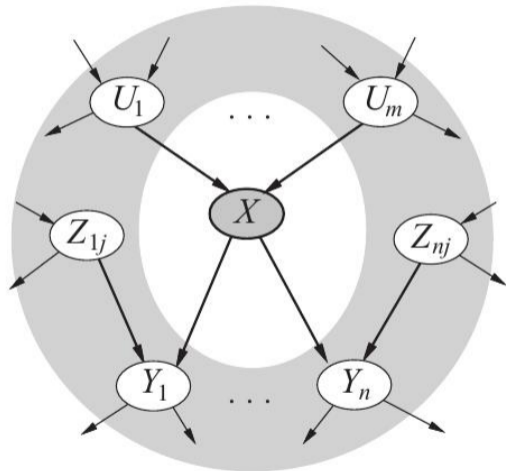
Markov blanket

- Recall $MB(X)$ — Markov blanket of X
 - $Parents(X)$
 - $Children(X)$
 - $Parents\ of\ Children(X)$
- $X \perp \neg MB(X) \mid MB(X)$
- Need to compute $P[y_i \mid x_1, x_2, \dots, x_{i-1}, x_{i+1}, \dots, x_n]$



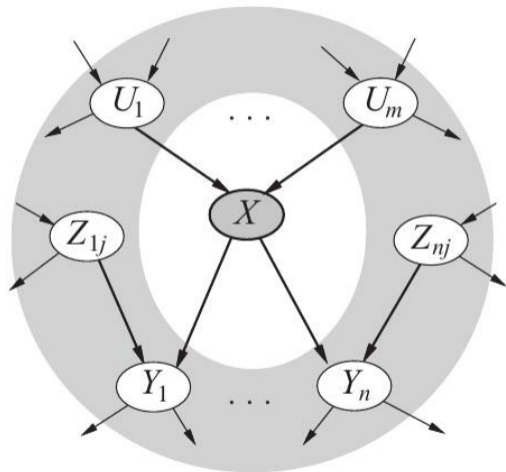
Markov blanket

- Recall $MB(X)$ — Markov blanket of X
 - $Parents(X)$
 - $Children(X)$
 - $Parents\ of\ Children(X)$
- $X \perp \neg MB(X) \mid MB(X)$
- Need to compute $P[y_i \mid x_1, x_2, \dots, x_{i-1}, x_{i+1}, \dots, x_n]$
- $x_1, x_2, \dots, x_{i-1}, x_{i+1}, \dots, x_n$ fixes $MB(V_i)$



Markov blanket

- Recall $MB(X)$ — Markov blanket of X
 - $Parents(X)$
 - $Children(X)$
 - $Parents\ of\ Children(X)$
- $X \perp \neg MB(X) \mid MB(X)$
- Need to compute $P[y_i \mid x_1, x_2, \dots, x_{i-1}, x_{i+1}, \dots, x_n]$
- $x_1, x_2, \dots, x_{i-1}, x_{i+1}, \dots, x_n$ fixes $MB(V_i)$
- Can compute $P[y_i \mid x_1, x_2, \dots, x_{i-1}, x_{i+1}, \dots, x_n]$ given conditional probability tables in the network



Gibbs sampling

- Move from $s_j = (x_1, x_2, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n)$ to $s_k = (x_1, x_2, \dots, x_{i-1}, y_i, x_{i+1}, \dots, x_n)$

Gibbs sampling

- Move from $s_j = (x_1, x_2, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n)$ to $s_k = (x_1, x_2, \dots, x_{i-1}, y_i, x_{i+1}, \dots, x_n)$
- Let $\bar{x} = (x_1, x_2, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$

Gibbs sampling

- Move from $s_j = (x_1, x_2, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n)$ to $s_k = (x_1, x_2, \dots, x_{i-1}, y_i, x_{i+1}, \dots, x_n)$
- Let $\bar{x} = (x_1, x_2, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$
- $P_{jk} = \frac{1}{n} P[y_i | \bar{x}]$

↑
Pick $i \in [1, n]$

Gibbs sampling

- Move from $s_j = (x_1, x_2, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n)$ to $s_k = (x_1, x_2, \dots, x_{i-1}, y_i, x_{i+1}, \dots, x_n)$
- Let $\bar{x} = (x_1, x_2, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$
- $P_{jk} = \frac{1}{n} P[y_i | \bar{x}] = \frac{1}{n} \frac{P(s_k)}{P(\bar{x})}$

Gibbs sampling

- Move from $s_j = (x_1, x_2, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n)$ to $s_k = (x_1, x_2, \dots, x_{i-1}, y_i, x_{i+1}, \dots, x_n)$
- Let $\bar{x} = (x_1, x_2, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$
- $$P_{jk} = \frac{1}{n} P[y_i | \bar{x}] = \frac{1}{n} \frac{P(s_k)}{P(\bar{x})} = \frac{1}{n} \frac{\pi_k}{P(\bar{x})}$$

Gibbs sampling

- Move from $s_j = (x_1, x_2, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n)$ to $s_k = (x_1, x_2, \dots, x_{i-1}, y_i, x_{i+1}, \dots, x_n)$
- Let $\bar{x} = (x_1, x_2, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$
- $P_{jk} = \frac{1}{n} P[y_i | \bar{x}] = \frac{1}{n} \frac{P(s_k)}{P(\bar{x})} = \frac{1}{n} \frac{\pi_k}{P(\bar{x})}$ ← ←
- Likewise $P_{kj} = \frac{1}{n} P[x_i | \bar{x}] = \frac{1}{n} \frac{P(s_j)}{P(\bar{x})} = \frac{1}{n} \frac{\pi_j}{P(\bar{x})}$

Gibbs sampling

- Move from $s_j = (x_1, x_2, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n)$ to $s_k = (x_1, x_2, \dots, x_{i-1}, y_i, x_{i+1}, \dots, x_n)$

- Let $\bar{x} = (x_1, x_2, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$

- $P_{jk} = \frac{1}{n} P[y_i | \bar{x}] = \frac{1}{n} \frac{P(s_k)}{P(\bar{x})}$

- Likewise $P_{kj} = \frac{1}{n} P[x_i | \bar{x}] = \frac{1}{n} \frac{P(s_j)}{P(\bar{x})}$

- Therefore, $\frac{P_{jk}}{P_{kj}} = \frac{P(s_k)}{P(s_j)} \Rightarrow P(s_j) \cdot P_{jk} = P(s_k) \cdot P_{kj}$

Gibbs sampling

- Move from $s_j = (x_1, x_2, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n)$ to $s_k = (x_1, x_2, \dots, x_{i-1}, y_i, x_{i+1}, \dots, x_n)$
- Let $\bar{x} = (x_1, x_2, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$
- $P_{jk} = \frac{1}{n} P[y_i | \bar{x}] = \frac{1}{n} \frac{P(s_k)}{P(\bar{x})} = \frac{1}{n} \frac{\pi_k}{P(\bar{x})}$
- Likewise $P_{kj} = \frac{1}{n} P[x_i | \bar{x}] = \frac{1}{n} \frac{P(s_j)}{P(\bar{x})} = \frac{1}{n} \frac{\pi_j}{P(\bar{x})}$
- Therefore, $\frac{P_{jk}}{P_{kj}} = \frac{\pi_k}{\pi_j}$
- Hence, $\pi_j \cdot P_{jk} = \pi_k \cdot P_{kj}$

Gibbs sampling

- Move from $s_j = (x_1, x_2, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n)$ to $s_k = (x_1, x_2, \dots, x_{i-1}, y_i, x_{i+1}, \dots, x_n)$
- $\pi_j \cdot P_{jk} = \pi_k \cdot P_{kj}$
- We have created a reversible Markov chain whose stationary distribution provides the true probabilities of the original Bayesian network!

Gibbs sampling

- Move from $s_j = (x_1, x_2, \dots, x_{i-1}, x_i, x_{i+1}, \dots, x_n)$ to $s_k = (x_1, x_2, \dots, x_{i-1}, y_i, x_{i+1}, \dots, x_n)$
- $\pi_j \cdot P_{jk} = \pi_k \cdot P_{kj}$
- We have created a reversible Markov chain whose stationary distribution provides the true probabilities of the original Bayesian network!
- Gibbs sampling is a special case of the more general **Metropolis-Hastings** algorithm

Gibbs sampling

- Since we are dealing with steady state probabilities, it is not necessary to change just one variable at a time

Gibbs sampling

- Since we are dealing with steady state probabilities, it is not necessary to change just one variable at a time
 - Generate an entirely new sample state (y_1, y_2, \dots, y_n)

Gibbs sampling

- Since we are dealing with steady state probabilities, it is not necessary to change just one variable at a time
 - Generate an entirely new sample state (y_1, y_2, \dots, y_n)
 - First generate y_1 , given x_2, x_3, \dots, x_n

Gibbs sampling

- Since we are dealing with steady state probabilities, it is not necessary to change just one variable at a time
 - Generate an entirely new sample state (y_1, y_2, \dots, y_n)
 - First generate y_1 , given x_2, x_3, \dots, x_n
 - Then generate y_2 , given y_1, x_3, \dots, x_n

Gibbs sampling

- Since we are dealing with steady state probabilities, it is not necessary to change just one variable at a time
 - Generate an entirely new sample state (y_1, y_2, \dots, y_n)
 - First generate y_1 , given x_2, x_3, \dots, x_n
 - Then generate y_2 , given y_1, x_3, \dots, x_n
 - ...
 - Then generate y_n , given y_1, y_2, \dots, y_{n-1}

Gibbs sampling

- Since we are dealing with steady state probabilities, it is not necessary to change just one variable at a time
 - Generate an entirely new sample state (y_1, y_2, \dots, y_n)
 - First generate y_1 , given x_2, x_3, \dots, x_n
 - Then generate y_2 , given y_1, x_3, \dots, x_n
 - ...
 - Then generate y_n , given y_1, y_2, \dots, y_{n-1}
- **Standard Gibbs sampler** — again a reversible Markov chain

MC MC

Markov Chain Monte Carlo

Approximate inference using Markov chains

- Bayesian network has variables V_1, V_2, \dots, V_n
- Use Gibbs sampling to set up a reversible Markov chain
- Stationary distribution will assign to state s the probability $P(s)$ in the Bayesian network

