



BITS Pilani
Pilani | Dubai | Goa | Hyderabad

PROBABILISTIC GRAPHICAL MODEL SESSION # 12: APPROXIMATE INFERENCE

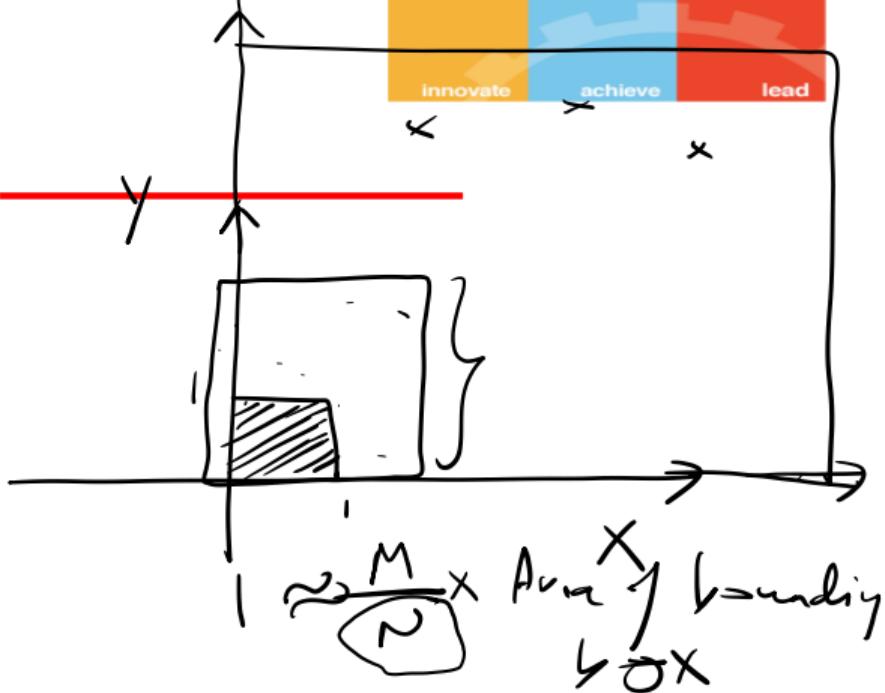
SEETHA PARAMESWARAN
seetha.p@pilani.bits-pilani.ac.in

Table of Contents

1 [Approximate Inference](#)

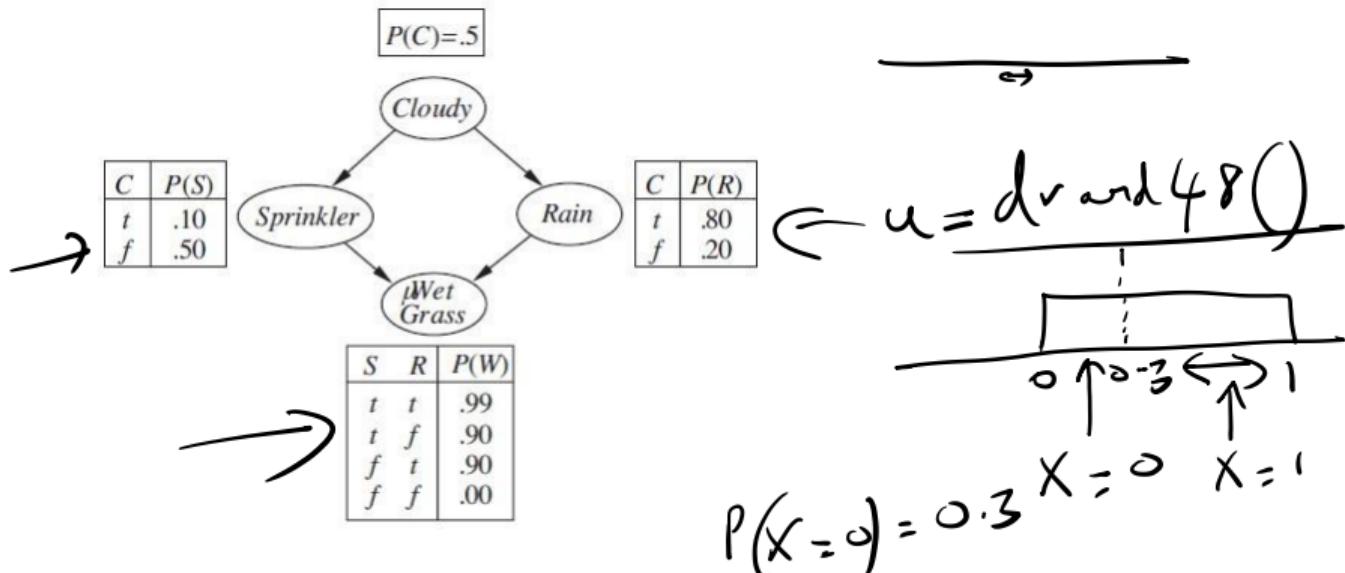
2 [Propagation-Based Approximation](#)

3 [Markov Chain Monte Carlo Simulation](#)



Approximate Inference Methods

- Given the intractability of exact inference in large, multiply connected networks, it is essential to consider approximate inference methods.

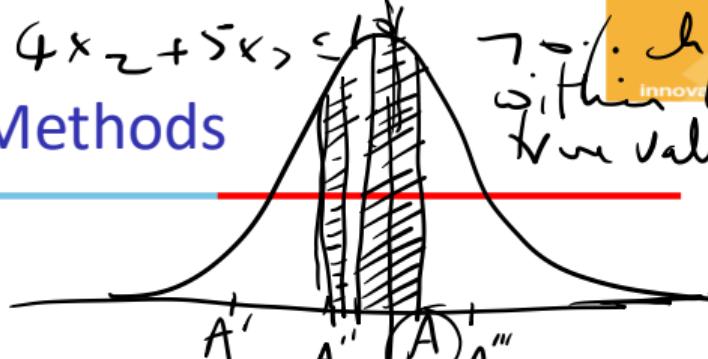


$$a_1^T x \leq b_1$$

$$3x_1 + 4x_2 + 5x_3 \leq b_2$$

Approximate Inference Methods

$$\frac{\sigma^2}{N} \rightarrow 0$$



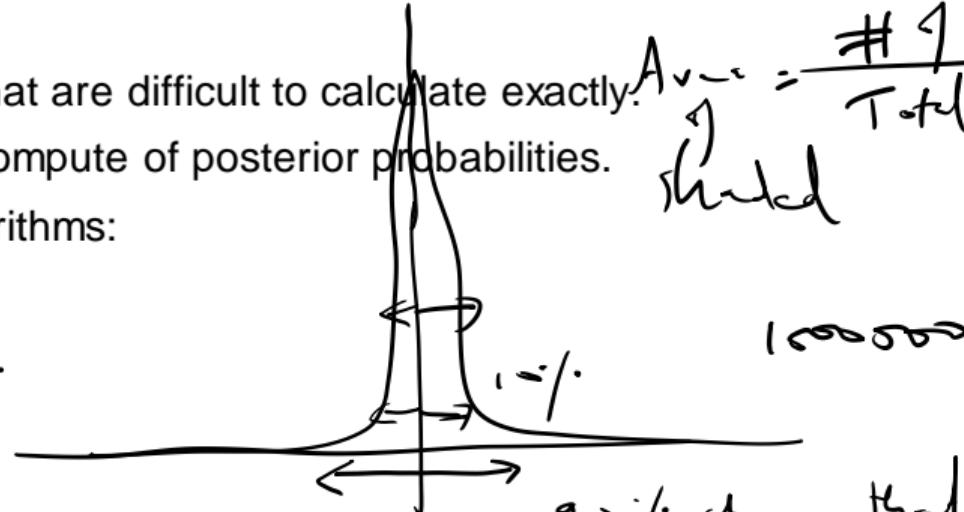
- Provide approximate answers whose accuracy depends on the number of samples generated.
- Used to estimate quantities that are difficult to calculate exactly.
- Sampling can be applied to compute of posterior probabilities.
- Two families of sampling algorithms:

1 Direct sampling

2 Markov chain sampling

$$Ave = \frac{\# \text{ of } X \text{ in sample}}{\text{Total } \# \text{ of } X \text{ s}}$$

shaded



are within 90% chance that we are within 1.5 standard deviations of the mean

$$f(x) \rightarrow \text{pdf} \quad P(X_i = 0) = 0.2 \quad P(X_i = 1) = 0.8$$

Direct Sampling in Bayesian Networks

$F(x) \rightarrow \text{cdf}$
The primitive

The primitive element in any sampling algorithm is the generation of samples from a known probability distribution.

Given a source of random numbers uniformly distributed in the range $[0, 1]$, it is a simple matter to sample any distribution on a single variable, whether discrete or continuous.

- PRIOR-SAMPLE generates samples from the prior joint distribution specified by the network.

function PRIOR-SAMPLE(bn) **returns** an event sampled from the prior specified by bn

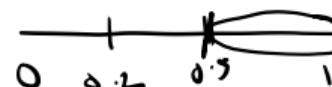
inputs: bn , a Bayesian network specifying joint distribution $\mathbf{P}(X_1, \dots, X_n)$

$\mathbf{x} \leftarrow$ an event with n elements

foreach variable X_i ~~in~~ X_1, \dots, X_n **do**

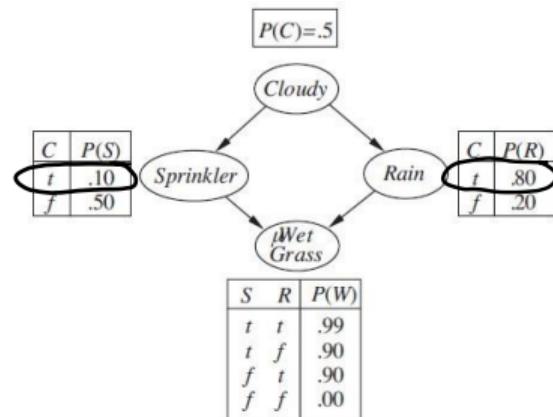
$\mathbf{x}[i] \leftarrow$ a random sample from $\mathbf{P}(X_i \mid parents(X_i))$

return x



Direct Sampling - Example

- Assume ordering [Cloudy, Sprinkler, Rain, WetGrass]
- Sample from $P(\text{Cloudy}) = < 0.5, 0.5 >$.
- Sampled value is true.
- Sample from $P(\text{Sprinkler} | \text{Cloudy} = \text{true}) = < 0.1, 0.9 >$. Sampled value is false.
- Sample from $P(\text{Rain} | \text{Cloudy} = \text{true}) = < 0.8, 0.2 >$.
- Sampled value is true.
- Sample from $P(\text{WetGrass} | \text{Sprinkler} = \text{false}, \text{Rain} = \text{true}) = < 0.9, 0.1 >$. Sampled value is true.
- PRIOR-SAMPLE returns the event
[true, false, true, true].

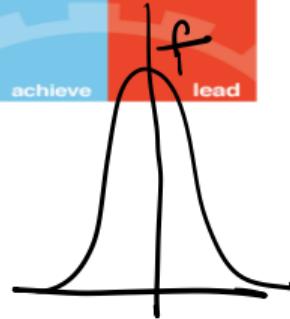


$$F^{-1}(u) = y$$

Rejection Sampling in Bayesian Networks



$$\int f(x) dx \stackrel{?}{=} \int f(x) = F(f) \stackrel{?}{=} cdf$$



- Produce samples from a hard-to-sample distribution given an easy-to-sample distribution.
- Used to compute conditional probabilities $P(X | e)$.
- Rejection sampling produces a consistent estimate of the true probability.

Rejection Sampling in Bayesian Networks

- First, it generates samples from the prior distribution specified by the network. Then, it rejects all those that do not match the evidence. Finally, the estimate $\hat{P}(X | e)$ is obtained by counting how often $X = x$ occurs in the remaining samples.

```
function REJECTION-SAMPLING( $X, e, bn, N$ ) returns an estimate of  $P(X|e)$ 
    inputs:  $X$ , the query variable
             $e$ , observed values for variables  $E$ 
             $bn$ , a Bayesian network
             $N$ , the total number of samples to be generated
    local variables:  $N$ , a vector of counts for each value of  $X$ , initially zero
    for  $j = 1$  to  $N$  do
         $x \leftarrow \text{PRIOR-SAMPLE}(bn)$ 
        if  $x$  is consistent with  $e$  then
             $N[x] \leftarrow N[x] + 1$  where  $x$  is the value of  $X$  in  $x$ 
    return NORMALIZE( $N$ )
```

10000

950
samples
survived
 $N(x)$
950

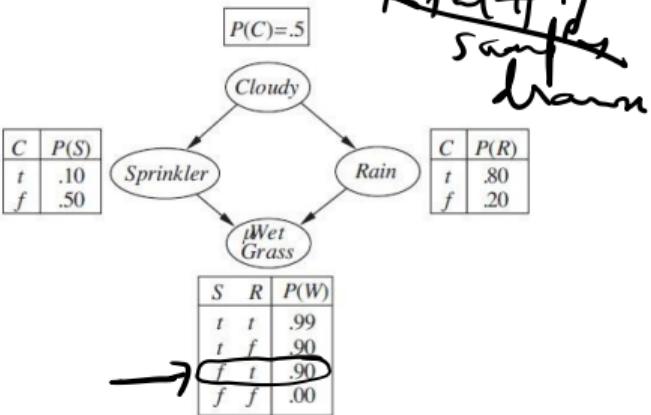
Rejection Sampling - Example

$$P(Q/e) = \frac{P(Q, e)}{P(e)}$$

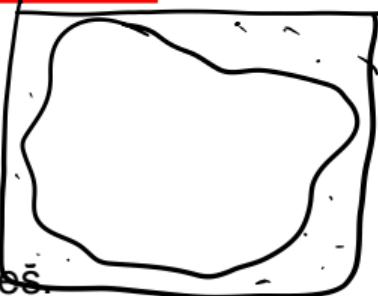
samples
 rain & sprinkler
 / total # samples
 = $\frac{\# Q \text{ samples having } e}{\# Q \text{ samples}}$
 width

- Assume that we wish to estimate $P(Rain | Sprinkler = true)$, using 100 samples.
- Of the 100 that we generate, suppose that 73 have $Sprinkler = false$ and are rejected, while 27 have $Sprinkler = true$;
- Of the 27 accepted samples, 8 have $Rain = true$ and 19 have $Rain = false$. Hence,
 $P(Rain | Sprinkler = true) \approx \text{Normalize}(< 8, 19 >) = < 0.296, 0.704 >$.

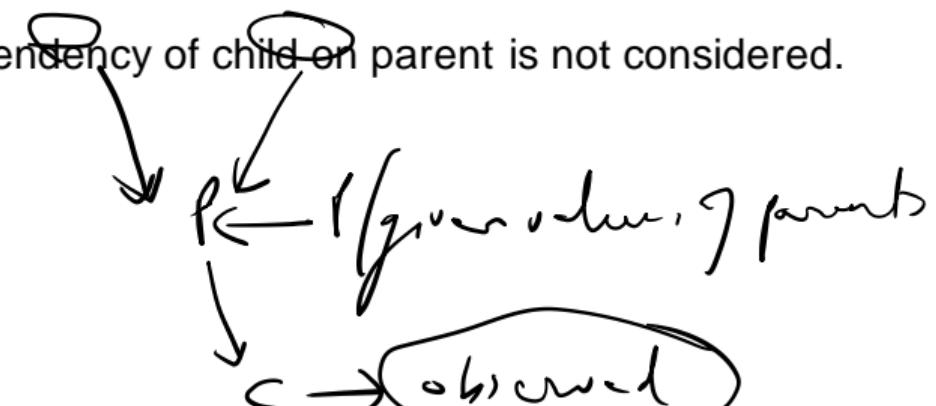
$$< \frac{8}{8+19} \mid \frac{19}{8+19} >$$



Rejection Sampling - Issues



- Rejection sampling is expensive, as it rejects many samples.
- Unusable for complex problems.
- When the child is observed, the dependency of child on parent is not considered.



Likelihood Weighting in Bayesian Networks

- Likelihood weighting avoids the inefficiency of rejection sampling by generating only events that are consistent with the evidence e .
- It is a particular instance of the general statistical technique of **importance sampling**, tailored for inference in Bayesian networks.
- LIKELIHOOD WEIGHTING fixes the values for the evidence variables E and samples only the nonevidence variables. This guarantees that each event generated is consistent with the evidence.
- Each event is weighted by the likelihood that the event accords to the evidence, as measured by the product of the conditional probabilities for each evidence variable, given its parents. Events in which the actual evidence appears unlikely should be given less weight.

Likelihood Weighting in Bayesian Networks

```

function LIKELIHOOD-WEIGHTING( $X, e, bn, N$ ) returns an estimate of  $P(X|e)$ 
  inputs:  $X$ , the query variable
   $e$ , observed values for variables  $E$ 
   $bn$ , a Bayesian network specifying joint distribution  $P(X_1, \dots, X_n)$ 
   $N$ , the total number of samples to be generated
  local variables:  $W$ , a vector of weighted counts for each value of  $X$ , initially zero
  for  $j = 1$  to  $N$  do
     $x, w \leftarrow \text{WEIGHTED-SAMPLE}(bn, e)$ 
     $W[x] \leftarrow W[x] + w$  where  $x$  is the value of  $X$  in  $x$ 
  return NORMALIZE(W)

```

```

function WEIGHTED-SAMPLE( $bn, e$ ) returns an event and a weight
   $w \leftarrow 1; x \leftarrow$  an event with  $n$  elements initialized from  $e$ 
  foreach variable  $X_i$  in  $X_1, \dots, X_n$  do
    if  $X_i$  is an evidence variable with value  $x_i$  in  $e$ 
    then  $w \leftarrow w \times P(X_i = x_i | \text{parents}(X_i))$ 
    else  $x[i] \leftarrow$  a random sample from  $P(X_i | \text{parents}(X_i))$ 
  return  $x, w$ 

```

*justify ~
for draws*

Likelihood Weighting - Example

- Assume Query is $P(\text{Rain} | \text{Cloudy} = \text{true}, \text{WetGrass} = \text{true})$ and the ordering is $[\text{Cloudy}, \text{Sprinkler}, \text{Rain}, \text{WetGrass}]$.
- First, the weight w is set to 1.0.
- Then an event is generated:

1 Cloudy is an evidence variable with value true. Therefore,

$$w \leftarrow w \times P(\text{Cloudy} = \text{true}) = 0.5$$

2 Sprinkler is not an evidence variable, so sample from

$$P(\text{Sprinkler} | \text{Cloudy} = \text{true}) = <0.1, 0.9>; \text{ suppose this returns false.}$$

$$P(\text{Cloudy} = \text{true})$$

3 Sample from $P(\text{Rain} | \text{Cloudy} = \text{true}) = <0.8, 0.2>$; suppose this returns true.

4 WetGrass is an evidence variable with value true. Therefore,

$$w \leftarrow w \times P(\text{WetGrass} = \text{true} | \text{Sprinkler} = \text{false}, \text{Rain} = \text{true}) = 0.45. = 0.5 \times 0.9$$

- WEIGHTED-SAMPLE returns the event $[\text{true}, \text{false}, \text{true}, \text{true}]$ with weight 0.45, and this is tallied under Rain = true.



$$\frac{1}{M} \sum_{i=1}^M f(x_i) \underbrace{P(x_i)}_{\text{Probability}} dx; M$$

$$\int f(x) P(x) dx$$

innovate

achieve

lead

Some theory about Importance Sampling

We would like to estimate th. expectation of some function f with respect to a distribution $P(x)$

$$E_P(f) = \frac{1}{M} \sum_{m=1}^M f(x_m)$$

What if P is not known or is very difficult to sample from?

We use another distribution Q and adjust for it

Some theory about Importance Sampling

$$E_{P(x)}[f(x)] = E_{Q(x)} \left[f(x) \frac{P(x)}{Q(x)} \right] = \int f(x) \frac{P(x)}{Q(x)} Q(x) dx$$

We use the standard estimator for expectations relative to Q . Let $D = \{x[1], x[2], \dots, x[M]\}$ be a set of samples drawn from Q

$$\hat{E}_D(f) = \frac{1}{M} \sum_{m=1}^M f(x[m]) \frac{P(x[m])}{Q(x[m])}$$

Some theory about Importance Sampling

This is the unnormalized Importance Sampling estimator

Proposition 12.1: For data sets D sampled from ϕ
 we have $E_D(E_D(F)) = E_{\phi(x)}[f(x)w(x)] = E_P(x)(f(x))$

Proof: $E_P\left[\frac{1}{M} \sum_{m=1}^M f(x[m]) \frac{P(x[m])}{Q(x[m])}\right] = \frac{1}{M} \sum_{m=1}^M E\left(f(x[m]) \frac{P(x[m])}{Q(x[m])}\right)$

Some theory on importance sampling

We need to take the expectation over the joint distribution on M-samples. The probability of selecting a particular sample $x_c[1], x_c[2] \dots x_c[M]$ is $\prod_{i=1}^M Q(x_c[i])$ and therefore

$$E_D(\hat{D}(f)) = \frac{1}{M} \sum_{x_c[1]} \sum_{x_c[2]} \dots \sum_{x_c[M]} \left(\frac{f(x_c[1])P(x_c[1])}{Q(x_c[1])} + \dots + \frac{f(x_c[M])P(x_c[M])}{Q(x_c[M])} \right) \prod_{i=1}^M Q(x_c[i])$$

We can rewrite this expression as in the next slide

Some theory on importance Sampling

$$\begin{aligned}
 &= \frac{1}{M} \sum_{x \in [1]} \sum_{x \in [2]} \dots \sum_{x \in [M]} f(x[i]) P(i \in [i]) \prod_{i=1, i \neq i}^M Q(x[i]) \\
 &\quad + \frac{1}{M} \sum_{x \in [1]} \sum_{x \in [2]} \dots \sum_{x \in [M]} f(x[2]) P(i \in [2]) \prod_{i=1, i \neq 2}^M Q(x[i]) \\
 &\quad + \dots + \frac{1}{M} \sum_{x \in [1]} \sum_{x \in [2]} \dots \sum_{x \in [M]} f(x[M]) P(i \in [M]) \prod_{i=1, i \neq M}^M Q(x[i])
 \end{aligned}$$

There are M terms in the above sum which are all very similar \rightarrow they each evaluate to the same quantity.

Some theory on importance sampling

Each term can be simplified to a term like the following

$$\left[\sum_{x[1]} \frac{1}{M} f(x[1]) p(x[1]) \right] \underbrace{\left[\sum_{x[2]} \sum_{x[3]} \dots \sum_{x[M]} \frac{1}{M} q(x[i]) \right]}_{\text{This evaluates to 1}}$$

We are left with $\sum_{x[1]} \frac{1}{M} f(x[1]) p(x[1]) = \frac{1}{M} E_p(x)(f(x))$

Thus we have $E_p(x)(f(x))$ overall

Note: this derivation is strictly valid for discrete

$$\text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y)$$



Some theory about Importance Sampling

Letting $\varepsilon_D = \hat{E}_D[f] - E_p[f(x)]$ we see that

the distribution of ε_D is $N(0, \frac{\sigma_Q^2}{M})$ where

$$\sigma_Q^2 = E_Q(x) [(f(x)w(x))^2] - E_Q(x) [f(x)w(x)]^2 = \text{Var}(f(x)w(x))$$

Why is this true?

$$\sigma_Q^2 = \text{Var}(\hat{E}_D(f)) = \frac{1}{M^2} M \cdot \text{Var}[f(x)w(x)] \text{ since we}$$

$\hat{E}_D(f)$ is the average of M independent samples $f(x_i)w(x_i)$

$$\sum_x P_B(x, e) = P_B(e) \neq 1$$

Normalized Importance Sampling

We do not have access to a distribution $P(x)$, only $\tilde{P}(x)$ which is not normalized

$$Z P(x) = \tilde{P}(x) \text{ where } Z \text{ is a normalizing constant}$$

In a Bayesian network $P(x)$ could be the posterior distribution $P_B(x|e)$ and $\tilde{P}(x)$ could be $P_B(x, e)$

$$\text{Let us define } \omega(x) = \frac{\tilde{P}(x)}{q(x)}$$

$$\begin{aligned} \omega(x) &= \frac{P(e) P_B(x|e)}{q(x)} \\ &= P_B(x, e) \end{aligned}$$

Normalized Importance Sampling

$$E_Q(\omega(x)) = \int \frac{\tilde{P}(x)}{Q(x)} Q(x) dx \quad \text{or} \quad \sum \frac{\tilde{P}(x)}{Q(x)} Q(x)$$

$$= \int \tilde{P}(x) dx \quad \text{or} \quad \sum \tilde{P}(x) \rightarrow \text{both are equal to } z$$

$$\begin{aligned} E_{P(x)}[f(x)] &= \sum_x f(x) P(x) = \sum_x f(x) \frac{P(x)}{Q(x)} Q(x) \\ &= \frac{1}{z} \sum_x \frac{f(x) \tilde{P}(x)}{Q(x)} Q(x) = \frac{1}{z} \underbrace{E_{Q(x)}[f(x) \omega(x)]}_{\text{---}} \end{aligned}$$

Normalized Importance Sampling

$$\underline{E}_{P(x)}[f(x)] = \sum_{m=1}^M \frac{f(x[m]) w(x[m])}{\sum_{m=1}^M w(x[m])} \rightarrow ?$$

where

$$w(x) = \frac{p(x)}{q(x)}$$

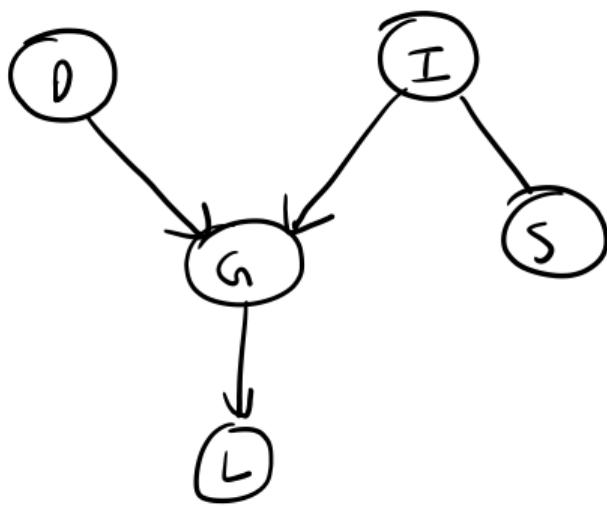
:

$x[1], x[2], \dots, x[m]$ are all samples drawn

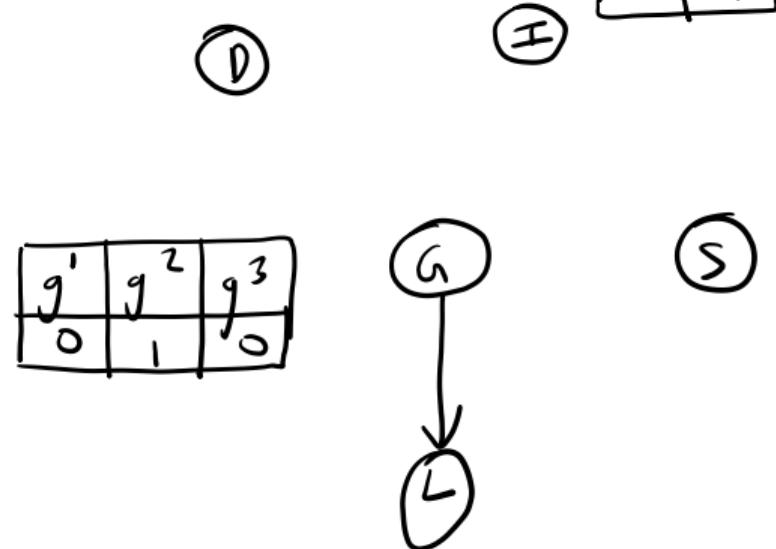
from $q(x)$

The bias of the normalized importance sampling estimator is not zero

Link to Likelihood Weighting



Original network B



Mutilated Network $B_{z=3}$

Link to Likelihood Weighting

Let ξ be a sample generated by the Likelihood weighting algorithm and let w be its weight. Then the distribution over ξ is as defined by the network

$$B_{Z=2} \quad \text{and} \quad w(\xi) = \frac{P_B(\xi)}{P_{B_{Z=2}}(\xi)} = \frac{\tilde{p}}{q}$$

Proof: Note that here $P_B \rightarrow \tilde{p}$ in the previous slides
 and $P_{B_{Z=2}}(\xi) \rightarrow q$

Link to Likelihood Weighting

Proof Continued

The only difference between sampling from $P_{B_{Z=2}}(\xi)$ and $P_B(\xi)$ concerns what happens when we draw a sample $z_i \in \text{Evidence}_{\text{Var}(y_i/\text{parent}(z_i))}$

In $P_{B_{Z=2}}(\xi)$, we draw this sample z_i with probability

$= 1$
In $P_B(\xi)$ we draw this sample with probability $P_B(z_i/\text{parent}(z_i))$

Link to Likelihood weighting

The weight in the Algorithm gets multiplied by

$$\frac{P_B(z_i | \text{parent}(z_i))}{1}$$

when z_i is not in evidence, $P_B(z_i | \text{parent}(z_i)) = P_{B_{z=3}}(z_i | \text{parent}(z_i))$

So we have $\frac{1}{1}$.

Multiplying the weights together, we get $w(\xi) = \frac{P_B(\xi)}{P_{B_{z=3}}(\xi)}$

Sampling-based Approximate Methods

- The methods using instantiations are generally known as **particle-based methods**.
- Each instantiation is known as a **particle**.
- Either create particles using a deterministic process, or sample particles from some distribution.
- Complete assignments to all of the network variables is commonly known as **full particles**.
- Disadvantage of full particle is that each particle covers only a very small part of the space.
- A collapsed particle specifies an assignment w only to some subset of the variables W , associating with it the conditional distribution $P(X|w)$.
- Assignments only to a subset $P(X|w)$ of variables of the network representing the conditional probability are commonly known as **collapsed particles**.

Table of Contents

- 1 [Approximate Inference](#)
- 2 [Propagation-Based Approximation](#)
- 3 [Markov Chain Monte Carlo Simulation](#)

Propagation-Based Approximation

- Use the same message propagation as in exact inference.
- Use Cluster graph instead of clique tree.
- Use **Loopy belief propagation**.
- Message propagation process may not converge in two passes, since information from one pass will circulate and affect the next round.
- In some cases, the propagation of beliefs may not converge at all.
- Use Bethe cluster graph to eliminate the loops.

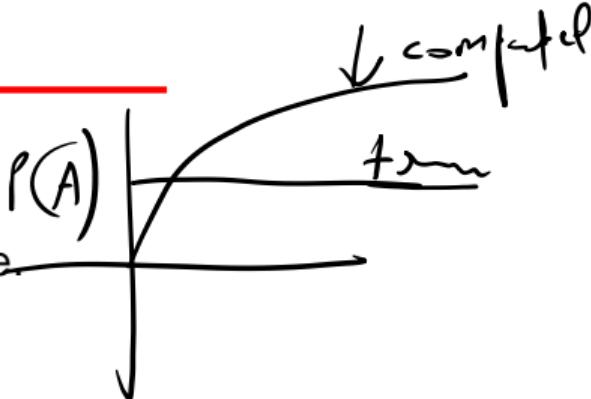
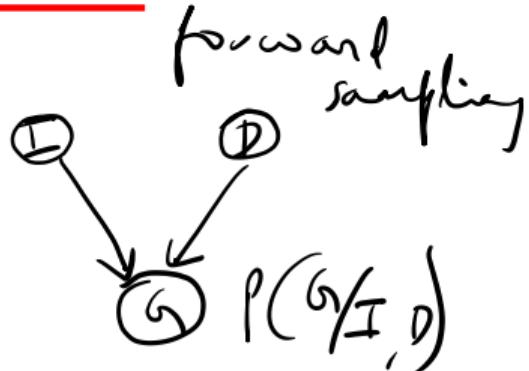


Table of Contents

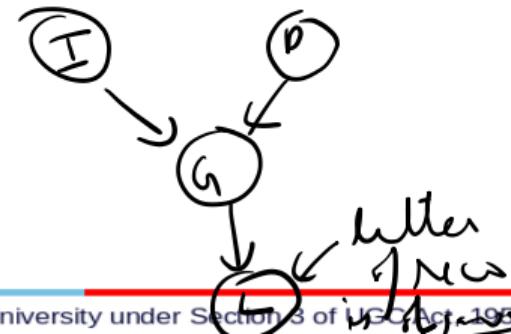
1 Approximate Inference

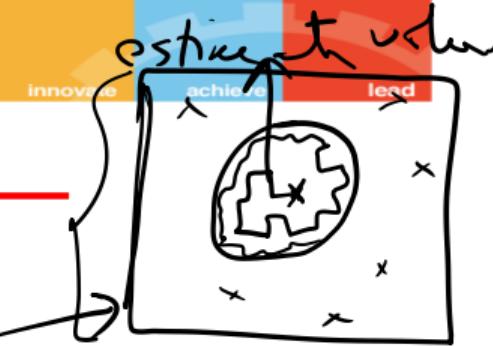
2 Propagation-Based Approximation

3 Markov Chain Monte Carlo Simulation



$$P(G) = f(G/I)$$





Markov Chain Monte Carlo Methods

- Generate a **sequence** of samples.
- Instead of generating each sample from scratch, MCMC algorithms ~~generate each~~ sample by making a random change to the preceding sample.
- Think of an MCMC algorithm as being in a particular **current state** specifying a value for every variable and generating a **next state** by making random changes to the current state.
- Markov chain methods apply equally well to directed and to undirected models.
- Works for both Bayesian and Markov networks.
- A particular form of MCMC is called Gibbs sampling, which is especially well suited for Bayesian networks.

Gibbs Sampling Algorithm

- Starts out by generating a sample of the unobserved variables from some initial distribution.
- Starting from initial sample, iterate over each of the unobserved variables, sampling a new value for each variable given our current sample for all other variables.
- This allows information to **flow** across the network as each variable is sampled.
- The Gibbs sampling algorithm for Bayesian networks starts with an arbitrary state (with the evidence variables fixed at their observed values) and generates a next state by randomly sampling a value for one of the nonevidence variables X_i .
- The sampling for X_i is done conditioned on the current values of the variables in the Markov blanket of X_i .
- The algorithm therefore wanders randomly around the state space flipping one variable at a time, but keeping the evidence variables fixed.

Gibbs Sampling Algorithm

```
function GIBBS-ASK( $X, \mathbf{e}, bn, N$ ) returns an estimate of  $\mathbf{P}(X|\mathbf{e})$ 
    local variables:  $\mathbf{N}$ , a vector of counts for each value of  $X$ , initially zero
     $\mathbf{Z}$ , the nonevidence variables in  $bn$ 
     $\mathbf{x}$ , the current state of the network, initially copied from  $\mathbf{e}$ 

    initialize  $\mathbf{x}$  with random values for the variables in  $\mathbf{Z}$ 
    for  $j = 1$  to  $N$  do
        for each  $Z_i$  in  $\mathbf{Z}$  do
            set the value of  $Z_i$  in  $\mathbf{x}$  by sampling from  $\mathbf{P}(Z_i|mb(Z_i))$ 
             $\mathbf{N}[x] \leftarrow \mathbf{N}[x] + 1$  where  $x$  is the value of  $X$  in  $\mathbf{x}$ 
    return NORMALIZE( $\mathbf{N}$ )
```

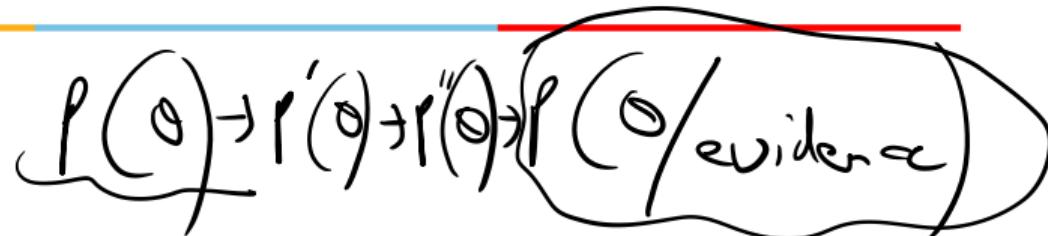
not from for each Z_i in Z do

Do this - k times

Gibbs Sampling - Example

- Assume Query is $P(Rain | Cloudy = \text{true}, WetGrass = \text{true})$ and the ordering is $[Cloudy, Sprinkler, Rain, WetGrass]$.
- The evidence variables *Sprinkler* and *WetGrass* are fixed to their observed values,
- The nonevidence variables *Cloudy* and *Rain* are initialized randomly; say to *true* and *false* respectively. The initial state is $[\text{true}, \text{true}, \text{false}, \text{true}]$.
- The nonevidence variables are sampled repeatedly in an arbitrary order.
 - ① *Cloudy* is sampled from $P(Cloudy | \underline{\text{Sprinkler} = \text{true}}, \underline{\text{Rain} = \text{false}})$. Suppose the result is *Cloudy = false*. Then the new current state is *false, true, false, true*. *new sample*
 - ② *Rain* is sampled from $P(Rain | \text{Cloudy} = \text{false}, \text{Sprinkler} = \text{true}, \text{WetGrass} = \text{true})$. Suppose this yields *Rain = true*. The new current state is $[\text{false}, \text{true}, \text{true}, \text{true}]$. *new sample*
- Each state visited during this process is a sample that contributes to the estimate for the query variable *Rain*.
- If the process visits 20 states where *Rain is true* and 60 states where *Rain is false*, then $\text{Normalize}(< 20, 60 >) = < 0.25, 0.75 >$.

Markov Chain Monte Carlo (MCMC)



- MCMC framework generates samples from the posterior distribution, where we cannot efficiently sample from the posterior directly.
- Construct an iterative process that gradually samples from distributions that are closer and closer to the posterior.
- The number of iterations is to be determined.



Markov Chain

$$\{0, 1\} \xrightarrow{\{0, 0.25\}} X_1 \\ \xrightarrow{\{0.25, 0.75\}} X_2 \\ \xrightarrow{> 0.75} X_3$$

lead

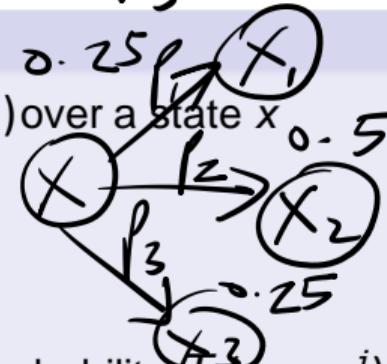
- A Markov chain is defined in terms of a graph of states over which the sampling algorithm takes a random walk.

$$p_1 + p_2 + p_3 = 1$$

Definition

Markov Chain defines a probabilistic transition model $T(x \rightarrow x^j)$ over a state x

$$T(T \leftarrow H) = 0.25$$

$$\forall x: \sum_{x^j} T(x \rightarrow x^j) = 1$$


The **transition model** T specifies for each pair of state x, x^j the probability $T(x \rightarrow x^j)$ of going from x to x^j . This transition probability applies whenever the chain is in state x .

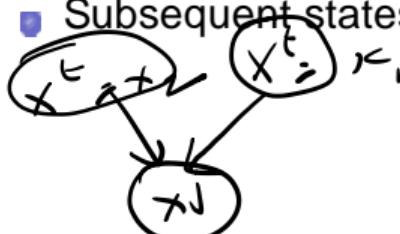
- MCMC is defined on **homogeneous** systems, where the system dynamics do not change over time.

$P(x_j \text{ at time } t+1) \leq P(x \text{ at } t \& x_j \text{ at time } t+1)$

Markov Chain



- Visualize the state space as a graph, with probability-weighted directed edges corresponding to transitions between different states.
- Generate a random sequence of states $x^{(0)}, x^{(1)}, x^{(2)}, \dots$
- The state of the process at step t is random variable $X^{(t)}$.
- Initial state $X^{(0)}$ is distributed according to some initial state distribution $P^{(0)}(X^{(0)})$.
- Subsequent states are defined by distributions $P^{(1)}(X^{(1)}), P^{(2)}(X^{(2)}), \dots$

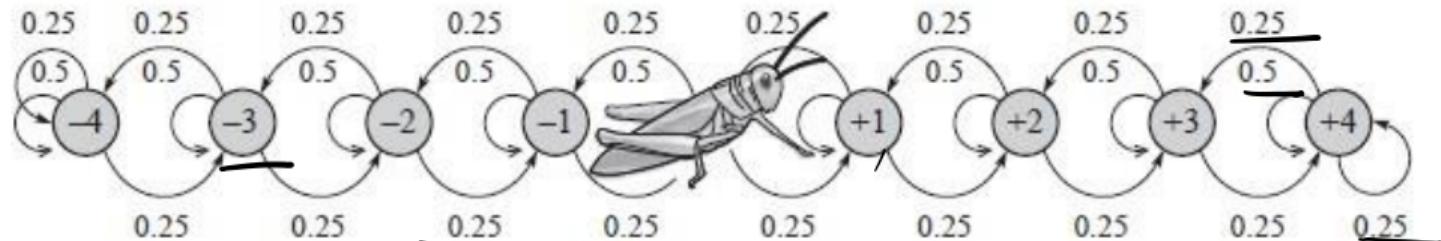


$$P^{(t+1)}(X^{t+1} = x^j) = \sum_{x \in Val(x)} P^{(t)}(X^t = x) P(x \rightarrow x^j)$$

$$P(X = T) + P(X = H) = 1$$

$$P\left(\frac{x^j}{x \text{ at } t}\right) P(x \text{ at } t)$$

Grasshopper Example



Time	-2	-1	0	1	2
$P^{(0)}$	0	0	1	0	0
$P^{(1)}$	0	<u>0.25</u>	<u>0.5</u>	<u>0.25</u>	0
$P^{(2)}$	<u>0.25^2</u>	<u>$2 * 0.25$</u>	<u>$0.5^2 + 2 * 0.25$</u>	<u>$2 * 0.25$</u>	<u>0.25^2</u>
	$=0.0625$	$=0.25$	$=0.375$	$=0.25$	$=0.0625$

Grasshopper Example

$$P^2(-1) = P'(-1)T(-1 \rightarrow -1) + P'(1)T(1 \rightarrow -1) + P'(0)T(0 \rightarrow -1)$$

(all other probabilities in P' are 0)

$$\text{Now } T(1 \rightarrow -1) = 0$$

$$\begin{aligned} \therefore P^2(-1) &= (0.5)(0.25) + (0.5)(0.25) \\ &= \underline{\underline{0.25}} \end{aligned}$$

Markov Chain Monte Carlo (MCMC)

- Markov chain Monte carlo (MCMC) sampling is a process that mirrors the dynamics of the Markov chain.
- The sample $X^{(t)}$ is drawn from the distribution $P^{(t)}$.
- Find out whether $P^{(t)}$ converges, and if so, to what limit.

$$P^t(X^t = x)$$

$$P^{(t+1)}(X^{t+1} = x^j) = \sum_{x \in Val(x)} P^{(t)}(X^t = x) T(x \rightarrow x^j)$$

- Apply on uniform distribution.
- Expected time required for a chain to reach boundaries of interval $[-K, K]$ is K^2 steps.



Markov Chain Monte Carlo (MCMC)

Algorithm 12.5 Generating a Markov chain trajectory

```
Procedure MCMC-Sample (
     $P^{(0)}(\mathbf{X})$ , // Initial state distribution
     $\mathcal{T}$ , // Markov chain transition model
     $T$  // Number of time steps
)
1   Sample  $\mathbf{x}^{(0)}$  from  $P^{(0)}(\mathbf{X})$ 
2   for  $t = 1, \dots, T$ 
3       Sample  $\mathbf{x}^{(t)}$  from  $\mathcal{T}(\mathbf{x}^{(t-1)} \rightarrow \mathbf{X})$ 
4   return  $\mathbf{x}^{(0)}, \dots, \mathbf{x}^{(T)}$ 
```

Stationary Distributions

- As MCMC process converges, $P^{(t+1)}$ will be close to $P^{(t)}$.

$$P^{(t)}(x^j) \approx P^{(t+1)}(X^{t+1} = x^j) = \sum_{x \in \text{Val}(x)} P^{(t)}(X^t = x)^T (x \rightarrow x^j)$$

The resulting distribution is $\pi(x)$.

- At equilibrium, for any state x^j , $P^{(t+1)}$ will be almost equal to $P^{(t)}$.
- Example : Grasshopper Markov Chain

Stationary Distributions

Definition (Stationary Distribution)

A distribution $\pi(x)$ is a stationary distribution for a Markov chain T if it satisfies

$$\pi(X = x^j) = \sum_{x \in Val(x)} \pi(X=x) T(x \rightarrow x^j)$$

A stationary distribution is also called an invariant distribution.

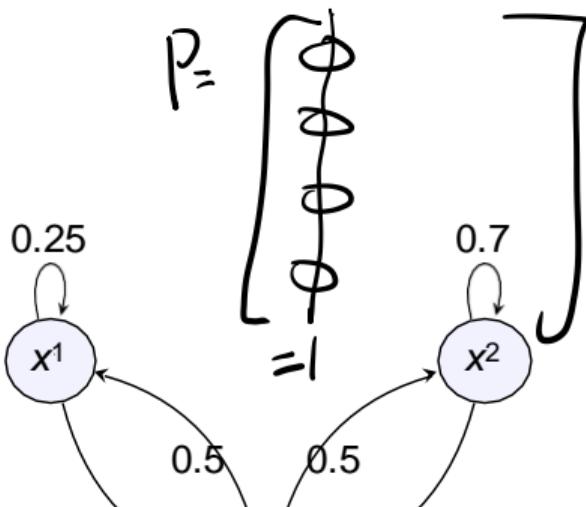
$P = \text{Markov} -$

~~Stationary Dist.~~ $P_u = \lambda u$

$P_u = u$

$P^2 u = P(Pu) \equiv Pu = u$

innovate achieve lead



Less uniform distribution.

$$\begin{aligned}\Pi(x^1) &= 0.25\Pi(x^1) + 0.5\Pi(x^3) \\ \Pi(x^2) &= 0.7\Pi(x^2) + 0.5\Pi(x^3) \\ \Pi(x^3) &= 0.75\Pi(x^1) + 0.3\Pi(x^2) \\ \Pi(x^1) + \Pi(x^2) + \Pi(x^3) &= 1\end{aligned}$$

3 variables
4 equations

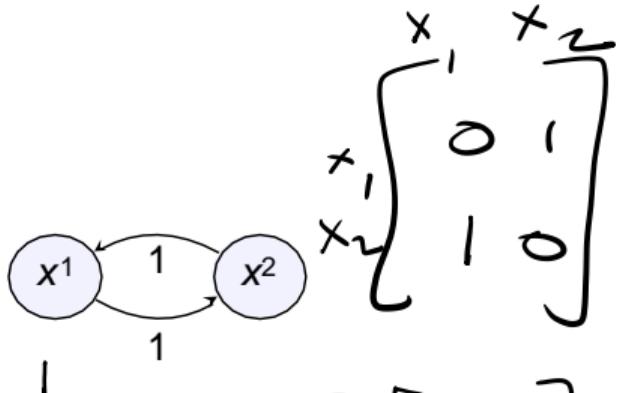
Unique solution

$$\begin{aligned}u_{K+1} &= P^K u \\ u_K &\quad \Pi(x^1) = 0.2 \\ &\quad \Pi(x^2) = 0.5 \\ u_0 &= \text{The eigenvector } \lambda_1 \text{ for eigenvalue } 1\end{aligned}$$

Periodic Markov Chains

$$P^t(x^1) = 1$$

$$P^{t+1}(x^1) = 0$$



Markov Chain with two states x^1 and x^2

$$T(x^1 \rightarrow x^2) = 1$$

$$T(x^2 \rightarrow x^1) = 1$$

Let $P^{(0)}(x^1) = 1$

$$P^{(t)}(x^1) = 1 \quad \left. \begin{array}{l} \\ \end{array} \right\}$$

$$P^{(t)}(x^2) = 0 \quad \left. \begin{array}{l} \\ \end{array} \right\}$$

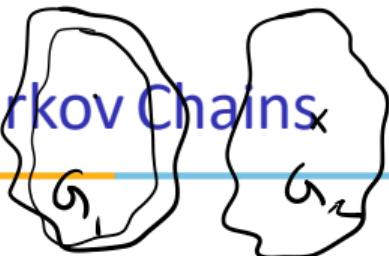
$P^{(t)}(x^1) = 1$ if t is even

$P^{(t)}(x^2) = 1$ if t is odd

$$\begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

- No convergence. So not a stationary distribution.
- Markov chains that exhibit a fixed cyclic behavior, are called periodic Markov chains.

Regular Markov Chains



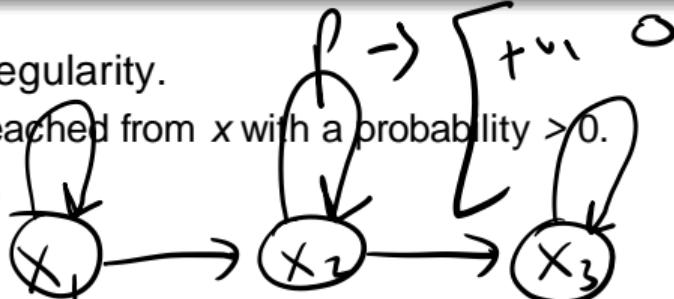
- A Markov chain is said to be regular if there exists some number k such that, for every $x, x^j \in \text{Val}(x)$, the probability of getting from x to x^j in exactly k steps is > 0 .
- Example: Grasshopper Markov chain $\rightarrow k = 9$.

$$P \rightarrow \begin{bmatrix} v & v \\ v & v \end{bmatrix} \quad (k = 1)$$

Theorem

If a finite state Markov chain T is regular, then it has a unique stationary distribution.

- Two conditions that together guarantee regularity.
 - Every state is connected. x^j can be reached from x with a probability > 0 .
 - For each state, there is a self-transition.



Using Markov Chains

- Goal: Compute $P(x \in S)$
 -) P is too hard to sample from directly.
- Construct a regular Markov Chain T whose unique stationary distribution is P .
- Sample $x^{(0)}$ from some arbitrary distribution Q .
- For $t = 0, 1, 2, \dots$
 -) Generate $x^{(t+1)}$ from $T(x^{(t)} \rightarrow x^j)$.

Using Samples

- Once the MCMC finds equilibrium, all samples $x^{(t)}$ are from stationary distribution $\pi(x)$.
- Collect and use samples $x^{(t)}$.
- Ideally collect every 100th sample, as nearby samples are correlated.

MCMC Algorithm I

- we are talking about C different Markov chains*
- For $c = 1, 2, \dots, C$
 -) Sample $x^{(0)}$ from arbitrary distribution Q (or Gibbs distribution can be used).
 - Repeat until equilibrium
 -) For $c = 1, 2, \dots, C$
 - ç Generate $x^{(c, t+1)}$ from $T(x^{(c, t)} \rightarrow x^t)$.
 -) Check for equilibrium
 -) $t := t + 1$

MCMC Algorithm II

- Repeat until sufficient samples
 -) Dataset $D := \text{empty}$
 -) For $c = 1, 2, \dots, C$ *We have C different Markov chains*
 - ↳ Generate $x^{(c, t+1)}$ from $T(x^{(c, t)} \rightarrow x^i)$.
 - ↳ $D := D \cup x^{(c, t+1)}$
 -) $t := t + 1$
- Compute Expectation
 -) Let $D := \{x[1], x[2], \dots, x[m]\}$
 -) Estimate expectation

$$\hat{E}_{\mathbf{P}}[f] = \frac{1}{M} \sum_{i=1}^M f(x[i])$$