



CMU 15-781

Lecture 9:

Bayesian Networks II

Teacher:

Gianni A. Di Caro

PROBABILISTIC INFERENCE TASKS

- Simple queries: $P(X_i | \mathbf{E} = \mathbf{e})$
- Conjunctive queries: $P(X_i, X_j | \mathbf{E} = \mathbf{e}) = P(X_i | \mathbf{E} = \mathbf{e}) P(X_j | X_i, \mathbf{E} = \mathbf{e})$
- Optimal decisions: decision networks include utility information; probabilistic inference is required for $P(\text{outcome} | \text{action}, \text{evidence})$
- Value of information: which evidence to seek next?
- Sensitivity analysis: which variables are most critical?



GENERAL INFERENCE PROCEDURE

Let's partition the set of random variables in the model in:

- Evidence variables \mathbf{E} , and be \mathbf{e} the list of observed values from them
- Remaining unobserved / hidden variables \mathbf{Y}
- Query variables X (let's consider single query for simplicity)

An inference query is $P(X | \mathbf{e})$? and can be evaluated as:

$$P(X | \mathbf{e}) = P(X, \mathbf{e}) / P(\mathbf{e}) = \alpha P(X, \mathbf{e}) = \alpha \sum_y P(X, \mathbf{e}, \mathbf{y})$$

(from prob recall: marginal of a subset of variables + normalization constants)

→ Given the full joint distribution, any probabilistic query can then be answered



INFERENCE WITH BNs

An inference query $P(X \mid \mathbf{e})$ can be evaluated as:

$$P(X \mid \mathbf{e}) = P(X, \mathbf{e}) / P(\mathbf{e}) = \alpha P(X, \mathbf{e}) = \alpha \sum_{\mathbf{y}} P(X, \mathbf{e}, \mathbf{y})$$

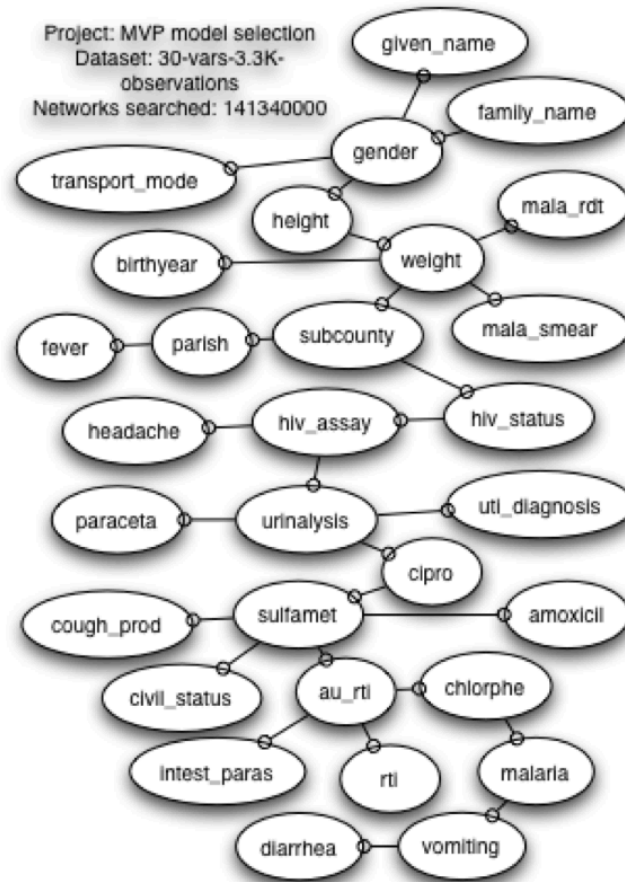
→ Given the **full joint distribution**, any probabilistic query can then be answered

A BN is a compact way to represent a joint distribution, where the terms in the joint distribution are written as products of conditional probabilities from the network

$$P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i \mid \text{parents}(X_i))$$

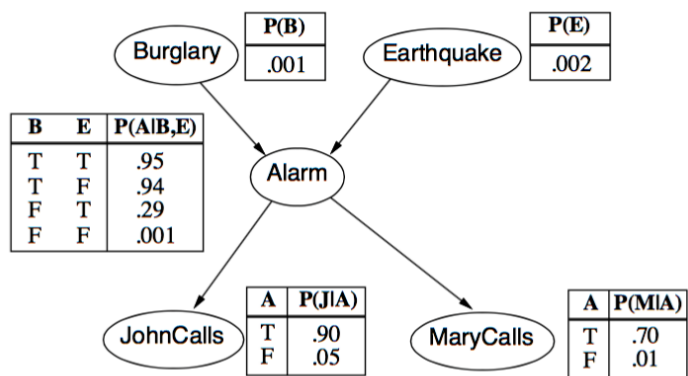


(Exact) Inference by enumeration?



EXACT INFERENCE BY ENUMERATION

$$P(X | e) = \alpha P(X, e) = \alpha \sum_y P(X, e, y) = \prod_{i=1}^n P(x_i | \text{parents}(X_i))$$



$$P(B, | J=T, M=T)$$

Evidence: J, M

Hidden: E, A

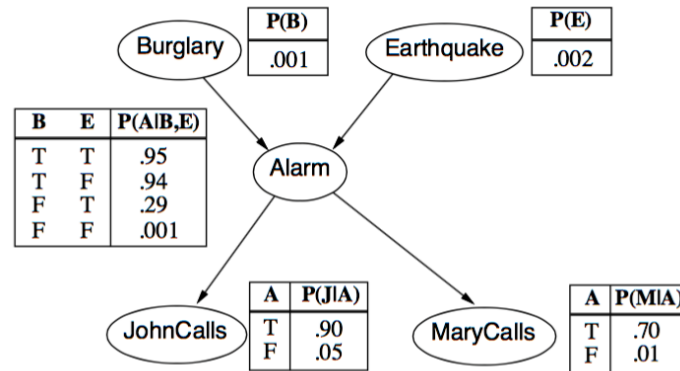
Query: B

$$P(B | j, m) = \alpha P(B, j, m) = \alpha \sum_e \sum_a P(B, e, a, j, m)$$

$$P(B | j, m) = \alpha \sum_e \sum_a P(B)P(e)P(a | b, e)P(j | a)P(m | a)$$



INFERENCE BY ENUMERATION



$$P(B \mid j, m) = \alpha \sum_e \sum_a P(B)P(e)P(a \mid b, e)P(j \mid a)P(m \mid a)$$

What is the complexity of this calculation? $O(n2^n)$



INFERENCE BY ENUMERATION

- $P(B=b)$ is a constant and can be moved outside the sums
- $P(e)$ can be moved outside the summation over a

$$P(B | j, m) = \alpha P(B) \sum_e P(e) \sum_a P(a | B, e) P(j | a) P(m | a)$$

$$\begin{aligned} P(b | j, m) &= \alpha P(b) \sum_e P(e) \sum_a P(a | b, e) P(j | a) P(m | a) \\ &= 0.001\alpha \sum_e P(e) \left[P(\neg a | b, e) P(j | \neg a) P(m | \neg a) + P(a | b, e) P(j | a) P(m | a) \right] \\ &= 0.001\alpha \sum_e P(e) [0.598525] \\ &= 0.001\alpha (0.02 \cdot 0.598525) + (0.998 \cdot 0.598525) = 0.000602\alpha \end{aligned}$$

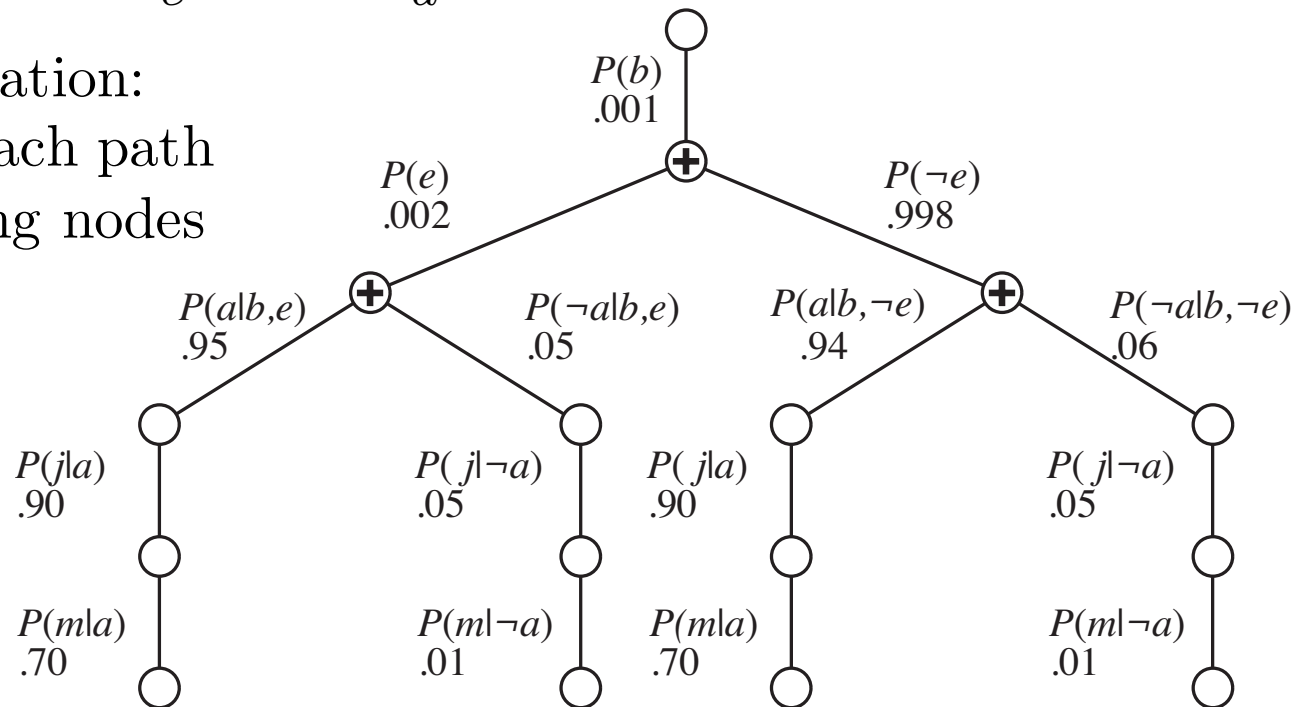


INFERENCE BY ENUMERATION

$$P(B \mid j, m) = \alpha P(B) \sum_e P(e) \sum_a P(a \mid B, e) P(j \mid a) P(m \mid a)$$

Top-down DF evaluation:

- \times Values along each path
- $+$ at the branching nodes

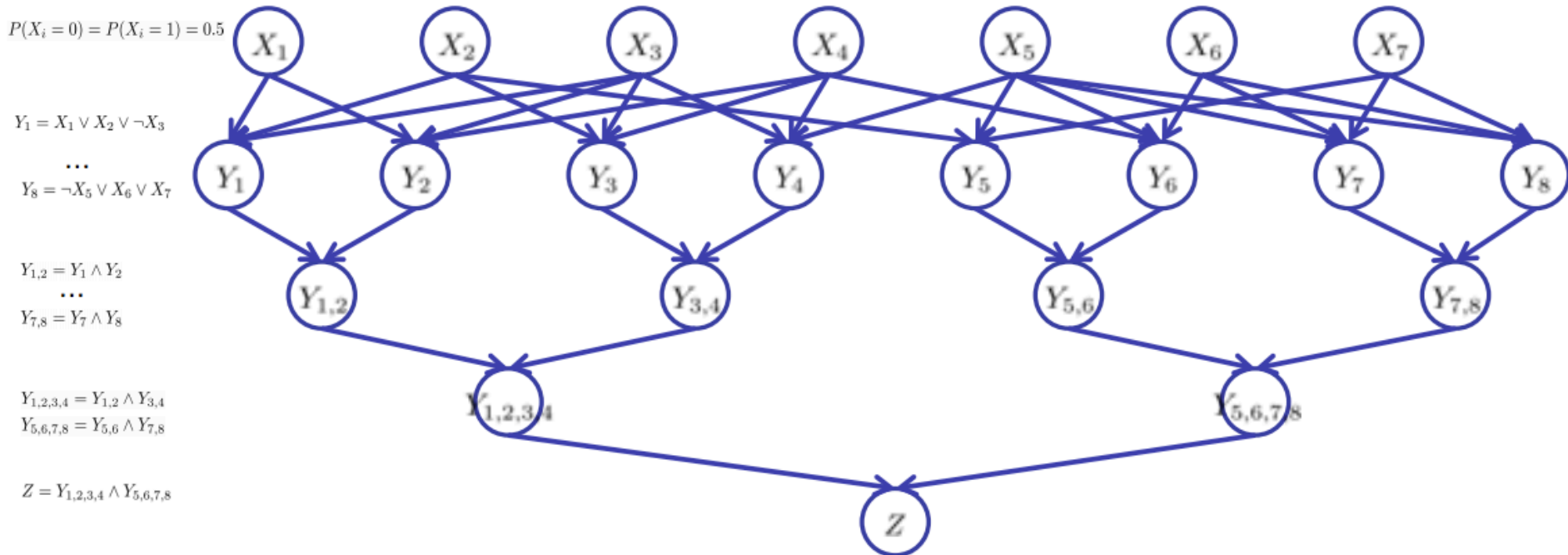


WORST CASE: NP-HARD

- Consider the 3-SAT clause:

$$(x_1 \vee x_2 \vee \neg x_3) \wedge (\neg x_1 \vee x_3 \vee \neg x_4) \wedge (x_2 \vee \neg x_2 \vee x_4) \wedge (\neg x_3 \vee \neg x_4 \vee \neg x_5) \wedge (x_2 \vee x_5 \vee x_7) \wedge (x_4 \vee x_5 \vee x_6) \wedge (\neg x_5 \vee x_6 \vee \neg x_7) \wedge (\neg x_5 \vee \neg x_6 \vee x_7)$$

which can be encoded by the following Bayes' net:



If we can answer $P(z)$ equal to zero or not, we answered whether the 3-SAT problem has a solution.

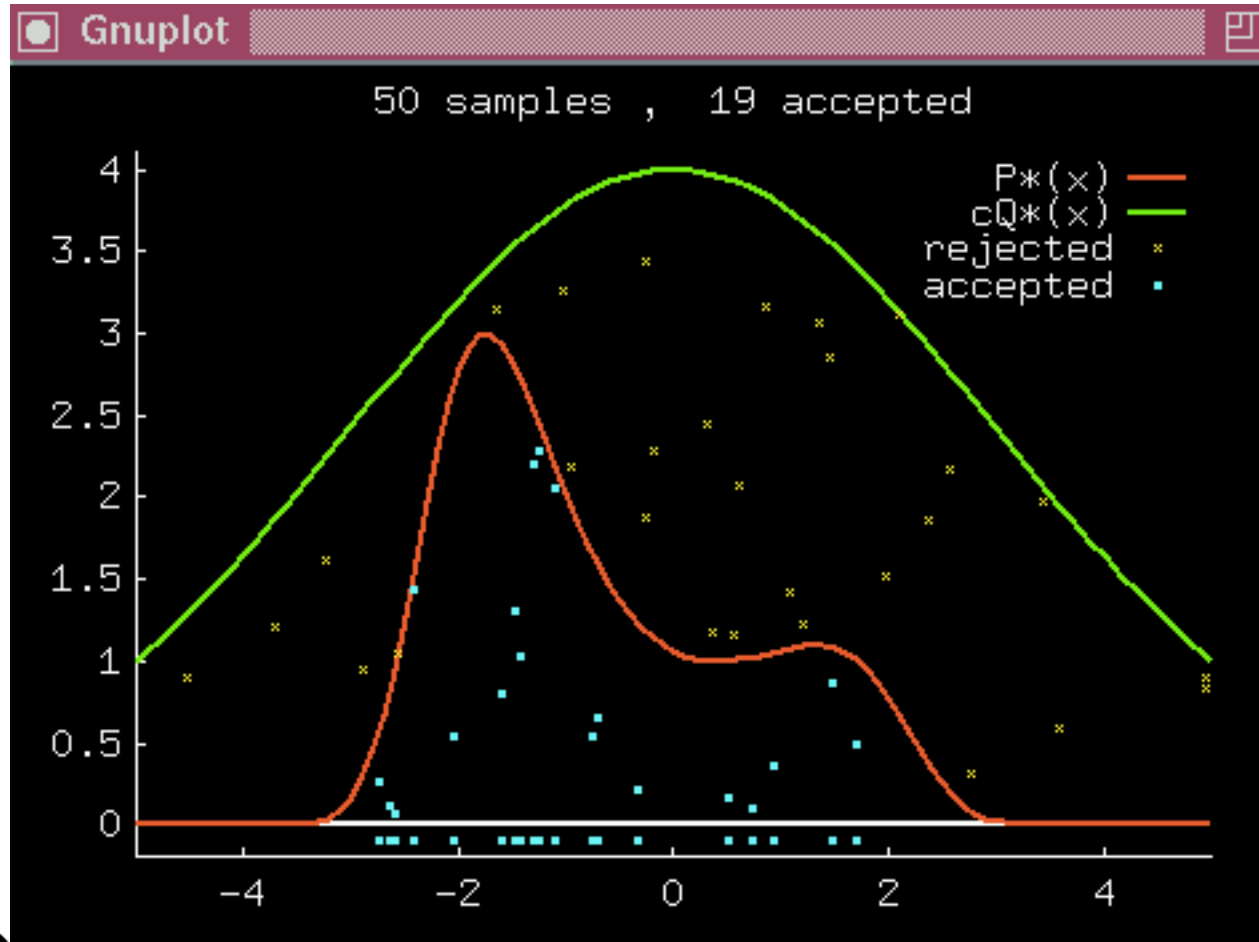


APPROXIMATE INFERENCE IN BN

- Given the general intractability of exact inference in large, multiply connected networks, it is essential to consider **approximate inference methods**
- **Monte Carlo algorithms:** randomized sampling methods ... we have already seen one example of them!
- **Basic idea:** we express the quantity we want to know as the expected value of a random variable \mathbf{X} , such as $\mu = E(\mathbf{X})$. Then we generate values $\mathbf{X}_1, \dots, \mathbf{X}_n$ independently and randomly from the distribution of \mathbf{X} and take their sample average as the estimate of μ (\rightarrow *law of large numbers*)
- **Problems:** it might be difficult to sample from \mathbf{X} 's distribution; a large number of samples might be needed



MONTE CARLO SAMPLING



SAMPLING FOR INFERENCE

- Have some method for generating samples given a known probability distribution (e.g., uniform in $[0,1]$)
- A sample is an assignment of values to each variable in the network
- Use samples to approximately compute posterior probabilities
- Queries can be issued after finish sampling

P



X_1	X_2	X_3	X_4	X_5
T	T	F	T	T
F	F	F	F	T
T	T	T	F	F
T	F	F	F	T
T	T	T	T	T

Prob (T,T,F,F,T)?

$\#(T,T,F,F,T) / \#\text{Samples}$



DIRECT SAMPLING METHODS

- Generate events from a network with no evidence
- Each variable is sampled in turn, in *topological order*
- The probability distribution from which the value of a variable is sampled is conditioned on the values already assigned to the variable's parents

```
function PRIOR-SAMPLE(bn) returns an event sampled from bn  
  inputs: bn, a belief network specifying joint distribution  $\mathbf{P}(X_1, \dots, X_n)$   
  
  x ← an event with n elements  
  for i = 1 to n do  
    xi ← a random sample from  $\mathbf{P}(X_i \mid \text{parents}(X_i))$   
      given the values of Parents(Xi) in x  
  return x
```

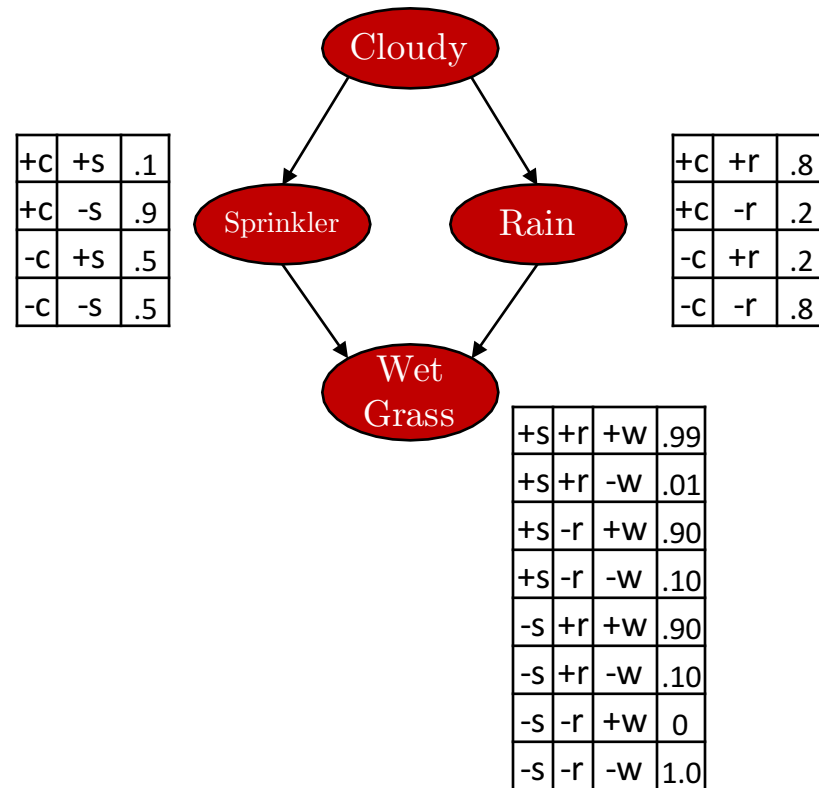


DIRECT SAMPLING

- Sample $\text{Pr}[C] = (.5, .5)$
 \Rightarrow true

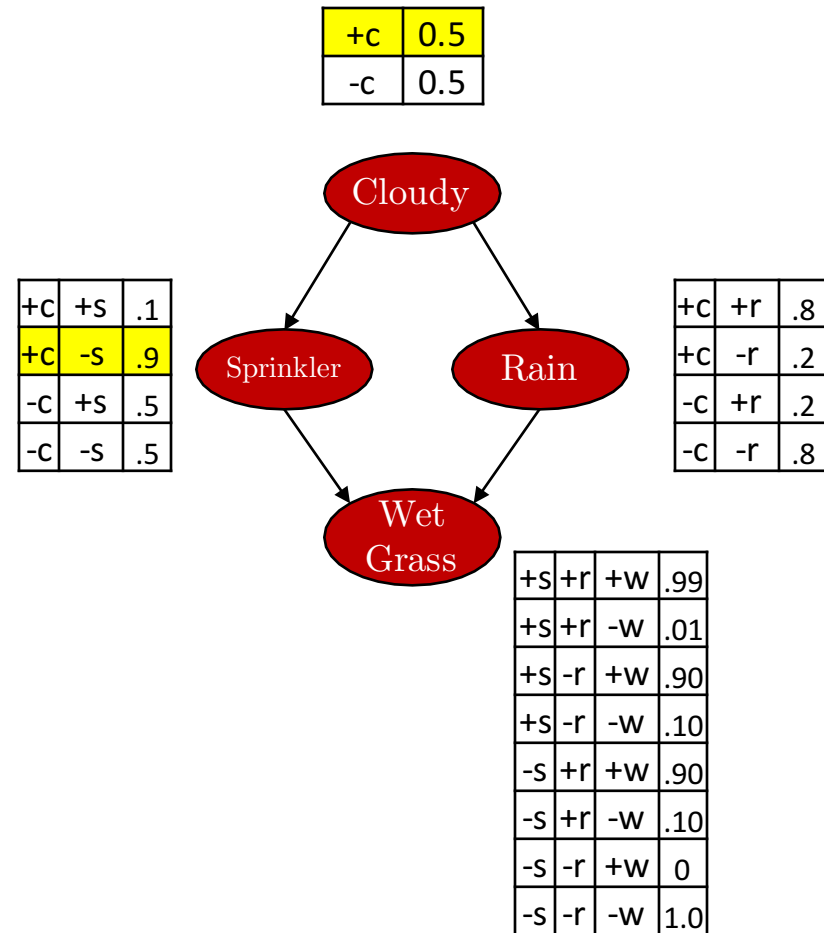
Order: C, S, R, W

+c	0.5
-c	0.5



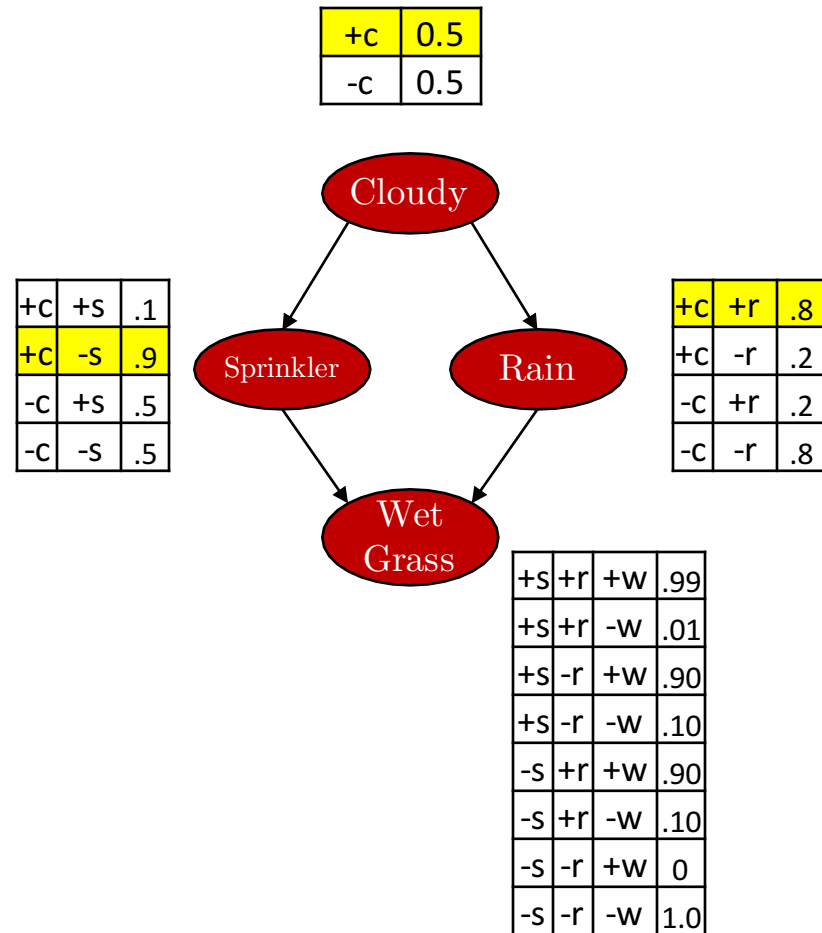
DIRECT SAMPLING

- Sample $\Pr[C]=(.5,.5)$
 \Rightarrow true
- Sample $\Pr[S|C=t]=(.1,.9)$
 \Rightarrow false



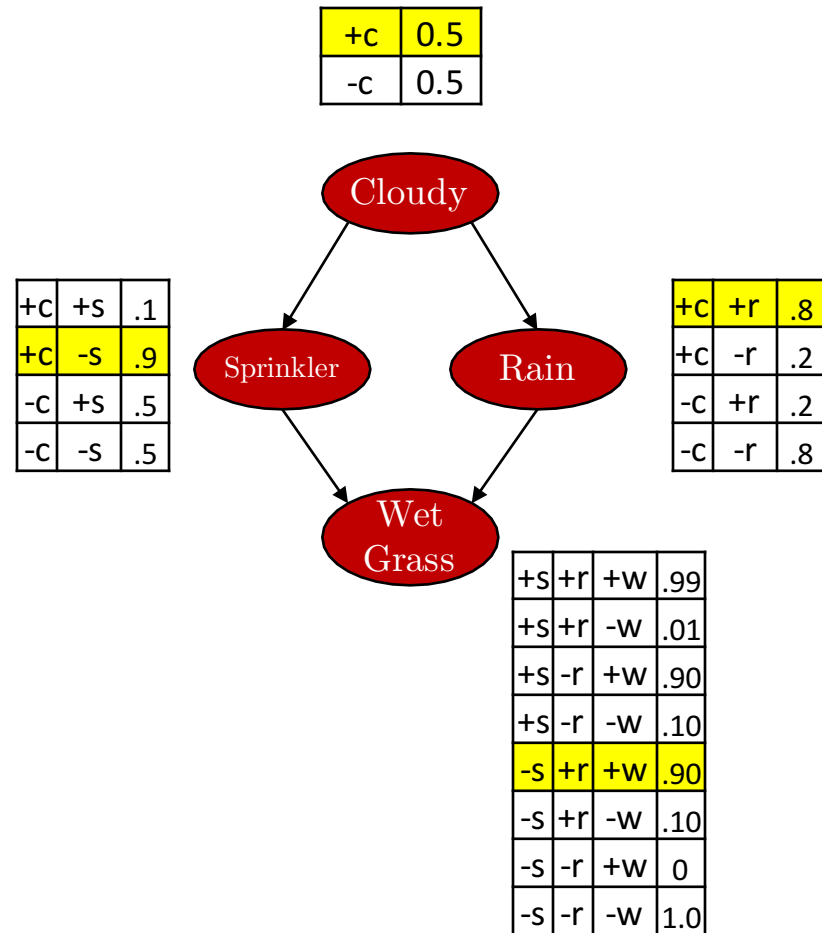
DIRECT SAMPLING

- Sample $\Pr[C] = (.5, .5)$
 \Rightarrow true
- Sample $\Pr[S|C=t] = (.1, .9)$
 \Rightarrow false
- Sample $\Pr[R|C=t] = (.8, .2)$
 \Rightarrow true



DIRECT SAMPLING

- Sample $\Pr[C]=(.5,.5)$
 \Rightarrow true
- Sample $\Pr[S|C=t]=(.1,.9)$
 \Rightarrow false
- Sample $\Pr[R|C=t]=(.8,.2)$
 \Rightarrow true
- Sample $\Pr[W|S=f,R=t]=(.9,.1)$
 \Rightarrow true
- Sampled [t,f,t,t]



DIRECT SAMPLING

Probability that PRIORSAMPLE generates a particular event

$$S_{PS}(x_1 \dots x_n) = \prod_{i=1}^n P(x_i | \text{parents}(X_i)) = P(x_1 \dots x_n)$$

i.e., the true prior probability

For large N, 32.4% of events (t,f,f,t) are expected

E.g., $S_{PS}(t, f, t, t) = 0.5 \times 0.9 \times 0.8 \times 0.9 = 0.324 = P(t, f, t, t)$ ←

Let $N_{PS}(x_1 \dots x_n)$ be the number of samples generated for event x_1, \dots, x_n

Then we have

$$\begin{aligned} \lim_{N \rightarrow \infty} \hat{P}(x_1, \dots, x_n) &= \lim_{N \rightarrow \infty} N_{PS}(x_1, \dots, x_n) / N \\ &= S_{PS}(x_1, \dots, x_n) \\ &= P(x_1 \dots x_n) \end{aligned}$$

That is, estimates derived from PRIORSAMPLE are **consistent**

Shorthand: $\hat{P}(x_1, \dots, x_n) \approx P(x_1 \dots x_n)$



REJECTION SAMPLING

- What about $P(X|\mathbf{e})$, i.e., when we have **evidence**?

function REJECTION-SAMPLING(X, \mathbf{e}, bn, N) **returns** an estimate of $P(X|\mathbf{e})$

local variables: \mathbf{N} , a vector of counts over X , initially zero

for $j = 1$ to N **do**

$\mathbf{x} \leftarrow$ PRIOR-SAMPLE(bn)

if \mathbf{x} is consistent with \mathbf{e} **then**

$\mathbf{N}[x] \leftarrow \mathbf{N}[x] + 1$ where x is the value of X in \mathbf{x}

return NORMALIZE($\mathbf{N}[X]$)

Similar to estimation of conditional
Probabilities directly from the real world

Problem: try to estimate $\Pr[\text{Rain} \mid \text{RedSkyAtNight}=t]$!



REJECTION SAMPLING

- Want to estimate $\Pr[\text{Rain}=t \mid \text{Sprinkler}=t]$
- 100 direct samples (no evidence included) are generated
- 73 have $S=f$, of which 12 have $R=t$
- 27 have $S=t$, of which 8 have $R=t$

$$\hat{P}(\text{Rain}|\text{Sprinkler} = \text{true}) = \text{NORMALIZE}(\langle 8, 19 \rangle) = \langle 0.296, 0.704 \rangle$$

- Error goes as $1/\sqrt{n}$, n = useful samples
- The estimate is consistent
- Too many samples thrown away! (because they are generated with direct sampling)



SOLUTION: LIKELIHOOD WEIGHTING

- Generate only samples that **agree with evidence**
- Fix the evidence vars and sample the nonevidence only
- → Each generated event is consistent with evidence
- **Weight each generated event according to likelihood that the event accords to evidence**
- The likelihood is measured as the product of the conditional probabilities for each evidence variable given its parents
- Event unlikely according to current evidence should weight less



LIKELIHOOD WEIGHTING

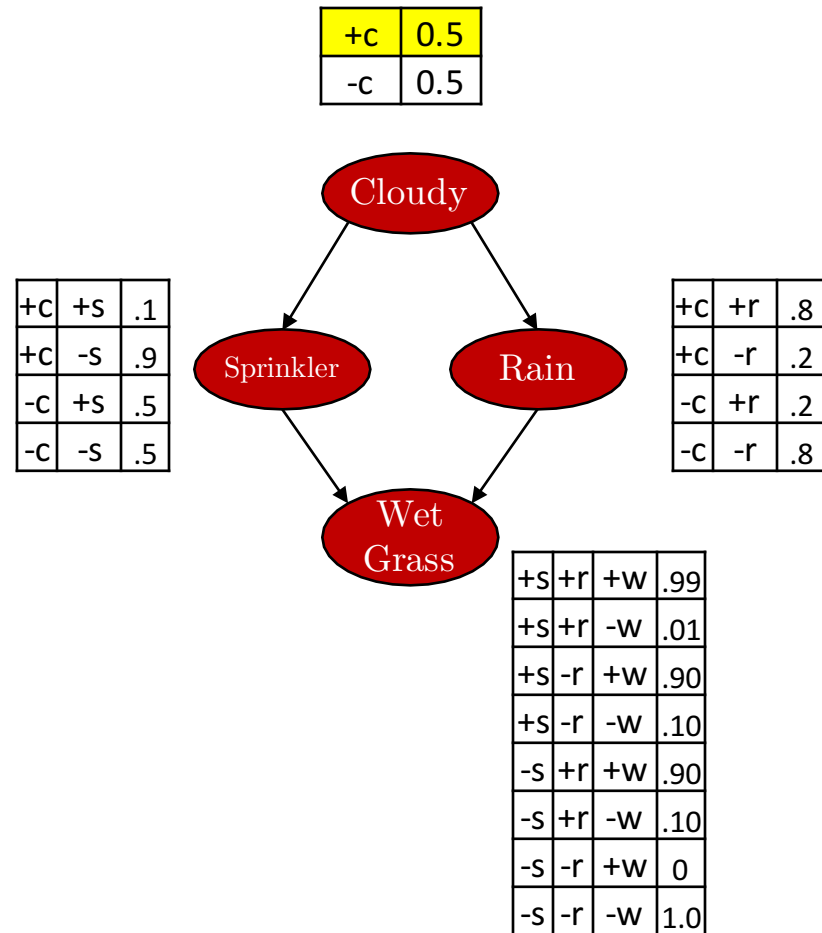
function LIKELIHOOD-WEIGHTING(X, \mathbf{e}, bn, N) **returns** an estimate of $P(X|\mathbf{e})$
local variables: \mathbf{W} , a vector of weighted counts over X , initially zero
for $j = 1$ to N **do**
 $\mathbf{x}, w \leftarrow$ WEIGHTED-SAMPLE(bn)
 $\mathbf{W}[x] \leftarrow \mathbf{W}[x] + w$ where x is the value of X in \mathbf{x}
return NORMALIZE($\mathbf{W}[X]$)

function WEIGHTED-SAMPLE(bn, \mathbf{e}) **returns** an event and a weight
 $\mathbf{x} \leftarrow$ an event with n elements; $w \leftarrow 1$
for $i = 1$ to n **do**
 if X_i has a value x_i in \mathbf{e}
 then $w \leftarrow w \times P(X_i = x_i \mid \text{parents}(X_i))$
 else $x_i \leftarrow$ a random sample from $\mathbf{P}(X_i \mid \text{parents}(X_i))$
return \mathbf{x}, w



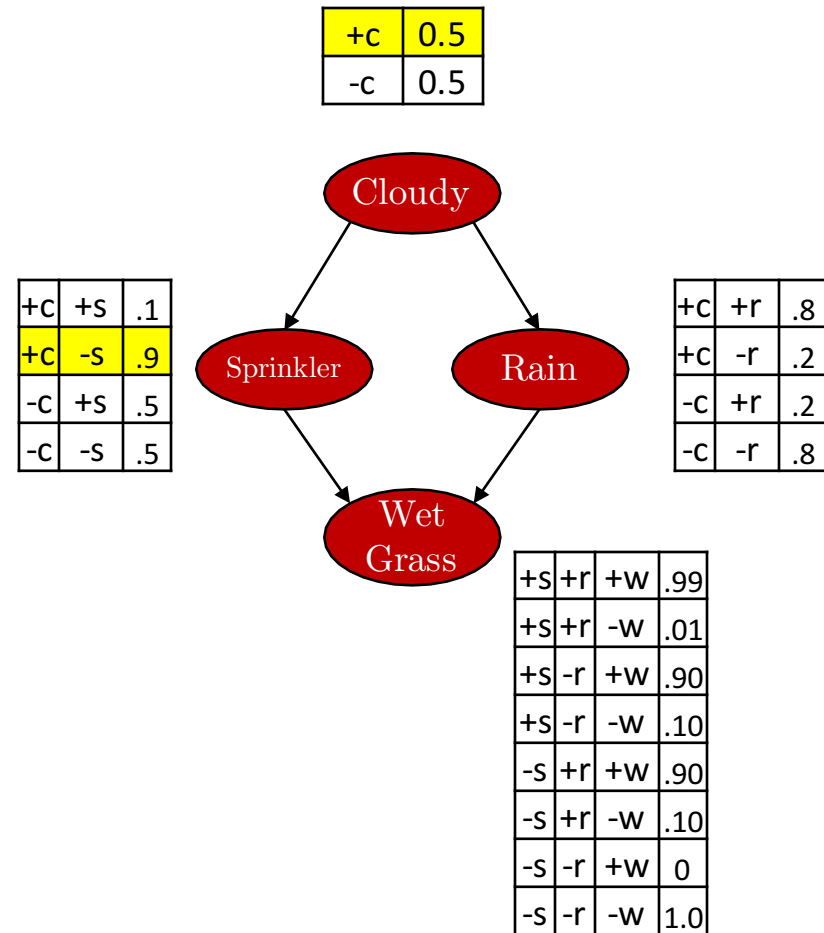
LIKELIHOOD WEIGHTING

- Evidence: $C=t, W=t$
- C is evidence var
 $\Rightarrow w = 1 \cdot \Pr[C=t] = 0.5$



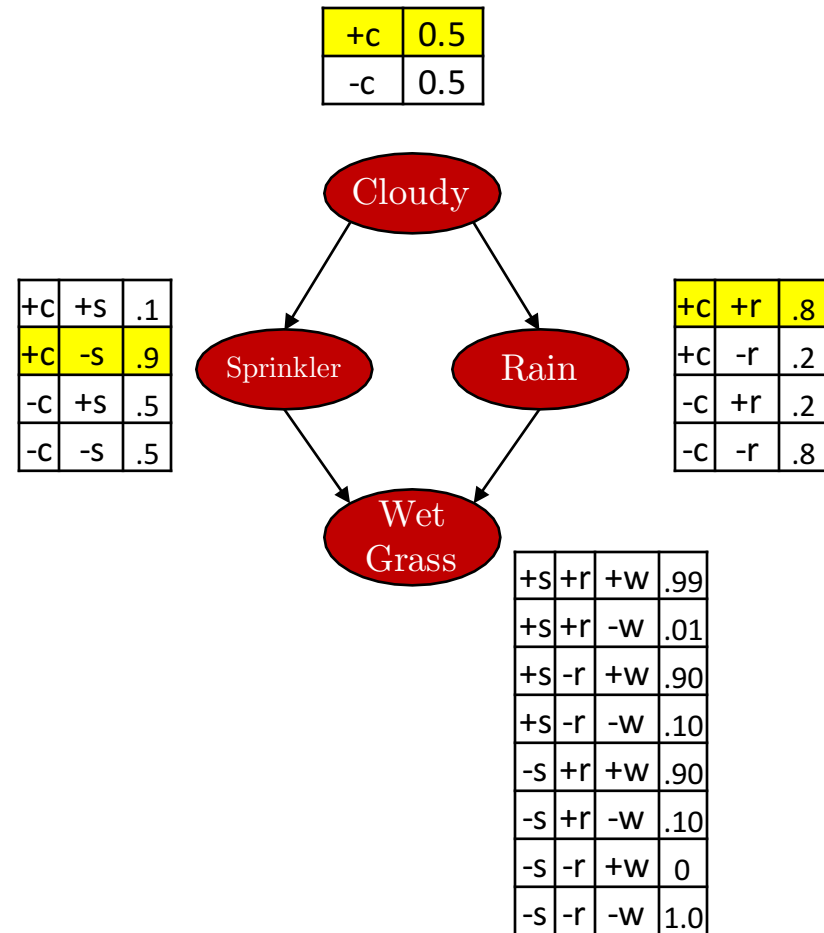
LIKELIHOOD WEIGHTING

- Evidence: $C=t, W=t$
- C is evidence var
 $\Rightarrow w = 1 \cdot \Pr[C=t] = 0.5$
- Sample $\Pr[S|C=t] = (.1, .9)$
 \Rightarrow false



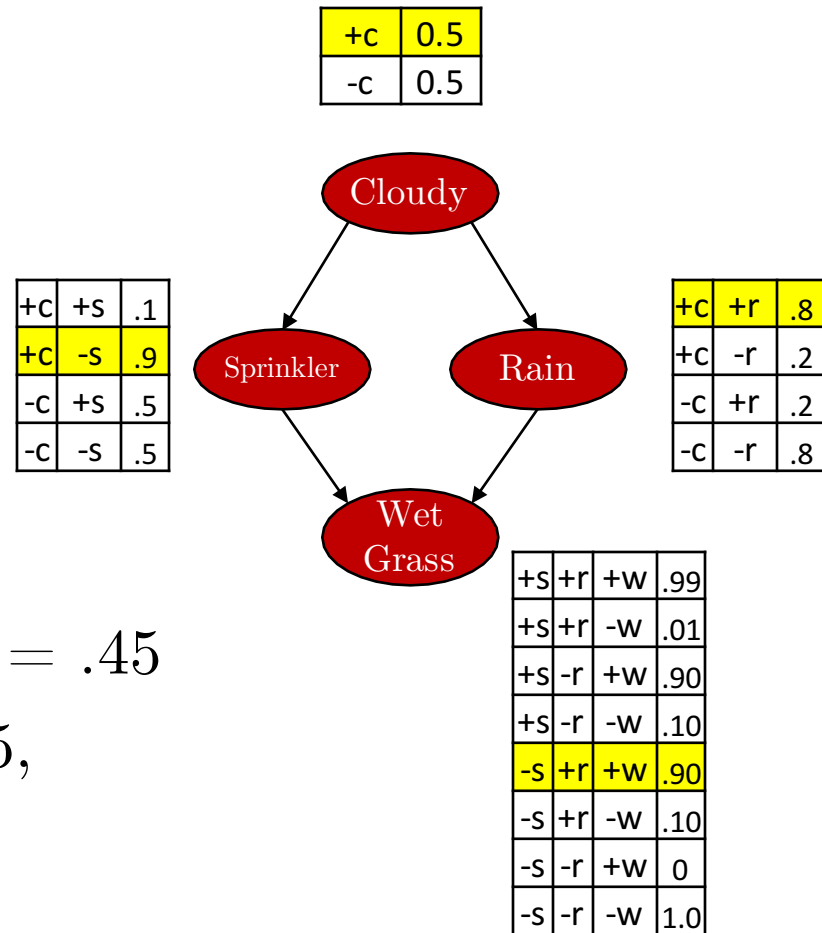
LIKELIHOOD WEIGHTING

- Evidence: $C=t, W=t$
- C is evidence var
 $\Rightarrow w = 1 \cdot \Pr[C=t] = 0.5$
- Sample $\Pr[S|C=t] = (.1, .9)$
 \Rightarrow false
- Sample $\Pr[R|C=t] = (.8, .2)$
 \Rightarrow true



LIKELIHOOD WEIGHTING

- Evidence: $C=t, W=t$
- C is evidence var
 $\Rightarrow w = 1 \cdot \Pr[C=t] = 0.5$
- Sample $\Pr[S|C=t] = (.1, .9)$
 \Rightarrow false
- Sample $\Pr[R|C=t] = (.8, .2)$
 \Rightarrow true
- W is evidence var
 $\Rightarrow w = 0.5 \cdot \Pr[W=t | S=f, R=t] = .45$
- Sampled $[t, f, t, t]$ with weight .45,
 tallied under $R=t$



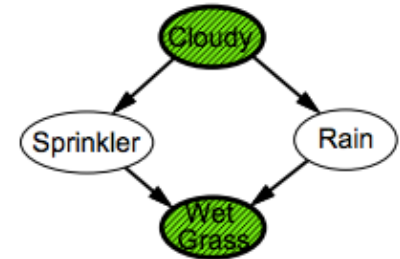
DISCUSSION

Sampling probability for WEIGHTEDSAMPLE is

$$S_{WS}(\mathbf{z}, \mathbf{e}) = \prod_{i=1}^l P(z_i | \text{parents}(Z_i))$$

Note: pays attention to evidence in **ancestors** only

⇒ somewhere “in between” prior and posterior distribution



Weight for a given sample \mathbf{z}, \mathbf{e} is

$$w(\mathbf{z}, \mathbf{e}) = \prod_{i=1}^m P(e_i | \text{parents}(E_i))$$

Weighted sampling probability is

$$\begin{aligned} & S_{WS}(\mathbf{z}, \mathbf{e}) w(\mathbf{z}, \mathbf{e}) \\ &= \prod_{i=1}^l P(z_i | \text{parents}(Z_i)) \prod_{i=1}^m P(e_i | \text{parents}(E_i)) \\ &= P(\mathbf{z}, \mathbf{e}) \text{ (by standard global semantics of network)} \end{aligned}$$

Hence likelihood weighting returns consistent estimates
but performance still degrades with many evidence variables
because a few samples have nearly all the total weight



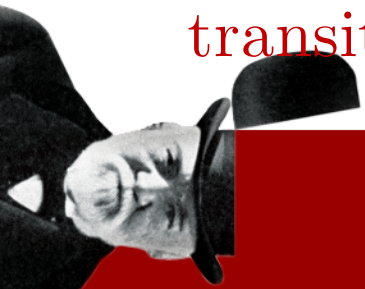
DISCUSSION

- Because of topological order, when sampling S and R the evidence $W=t$ is ignored \Rightarrow samples with $S=f$ and $R=f$ although evidence rules this out
- Weight makes up for this difference: above weight would be 0
- If we have 100 samples with $R=t$ and total weight 1, and 400 samples with $R=f$ and total weight 2, what is estimate of $R=t$?
- Problem: bad if evidence variables occur later in ordering



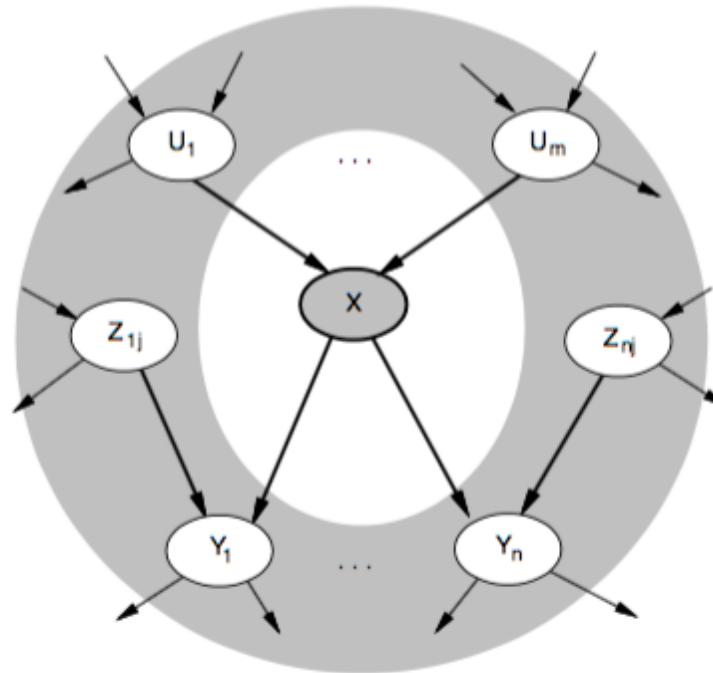
GIBBS SAMPLING

- **Markov Chain Monte Carlo (MCMC)**: each sample is generated by making a random change to the preceding sample
- MCMC are algorithms with a *state*, the next state is generated from the current one
- Specific MCMC: **Gibbs sampling**, the sampling process settles into a “dynamic equilibrium” in which the long-run fraction of time spent in each state is exactly proportional to its posterior probability
- The states are generated given the *Markov blanket*: state transitions are defined by the conditional distribution



MARKOV BLANKET

Each node is conditionally independent of all others given its Markov blanket: parents + children + children's parents



GIBBS SAMPLING

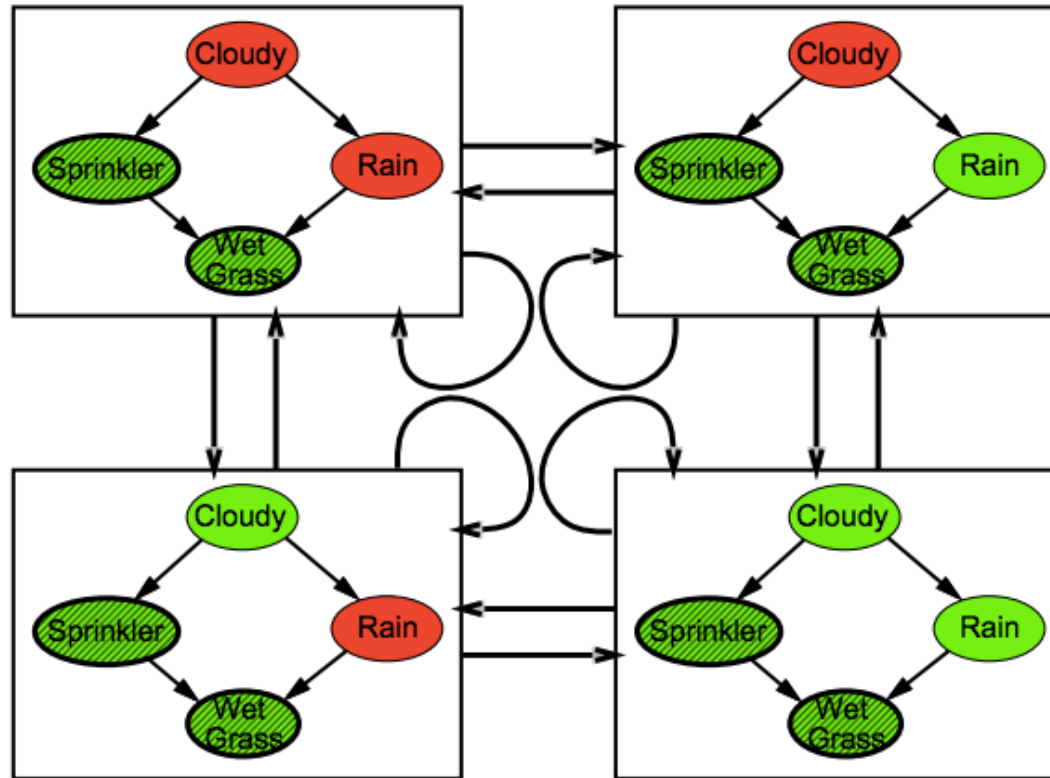
function GIBBS-SAMPLING(X, \mathbf{e}, bn, N) **returns** an estimate of $P(X|\mathbf{e})$
local variables: $\mathbf{N}[X]$, a vector of counts over X , initially zero
 \mathbf{Z} , the nonevidence variables in bn
 \mathbf{x} , the current state of the network, initially copied from \mathbf{e}

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



GIBBS SAMPLING

With $Sprinkler = true$, $WetGrass = true$, there are four states:



Wander about for a while, average what you see

GIBBS SAMPLING

Estimate $P(\textit{Rain} | \textit{Sprinkler} = \textit{true}, \textit{WetGrass} = \textit{true})$

Sample *Cloudy* or *Rain* given its Markov blanket, repeat.
Count number of times *Rain* is true and false in the samples.

E.g., visit 100 states

31 have *Rain* = true, 69 have *Rain* = false

$$\hat{P}(\textit{Rain} | \textit{Sprinkler} = \textit{true}, \textit{WetGrass} = \textit{true}) \\ = \text{NORMALIZE}(\langle 31, 69 \rangle) = \langle 0.31, 0.69 \rangle$$

Theorem: chain approaches **stationary distribution**:
long-run fraction of time spent in each state is exactly
proportional to its posterior probability

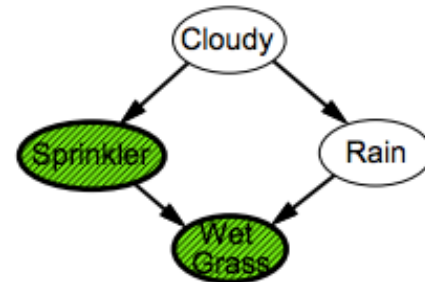


MB SAMPLING

Markov blanket of *Cloudy* is
Sprinkler and *Rain*

Markov blanket of *Rain* is

Cloudy, *Sprinkler*, and *WetGrass*



Probability given the Markov blanket is calculated as follows:

$$P(x'_i | mb(X_i)) = P(x'_i | parents(X_i)) \prod_{Z_j \in Children(X_i)} P(z_j | parents(Z_j))$$

Easily implemented in message-passing parallel systems, brains

Main computational problems:

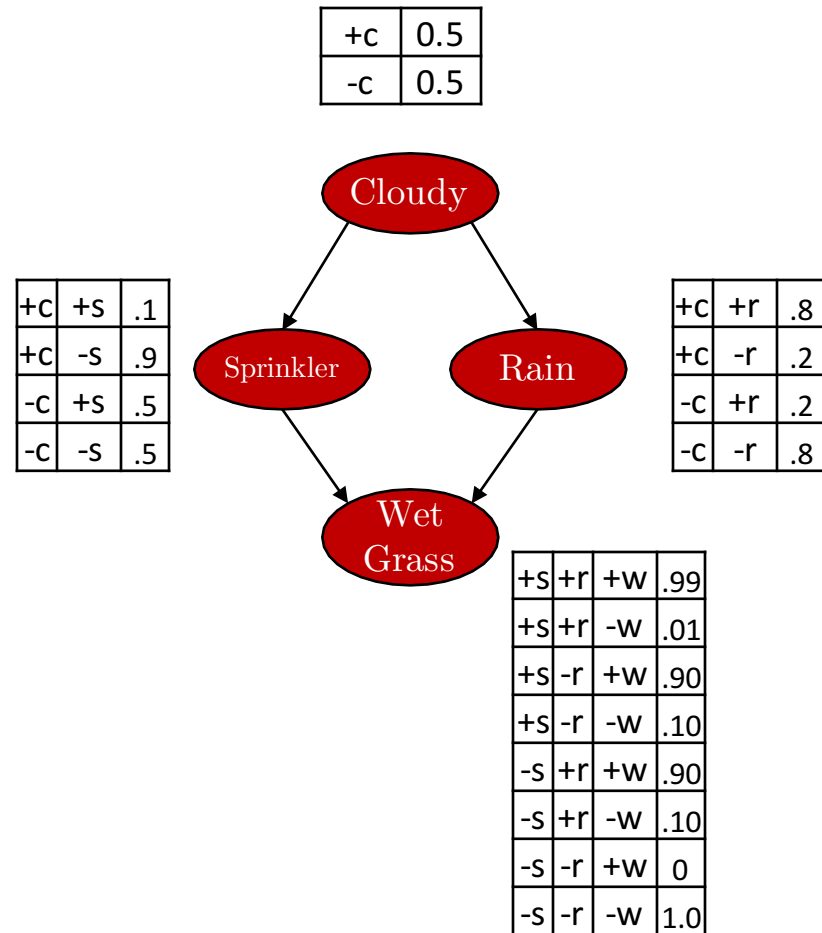
- 1) Difficult to tell if convergence has been achieved
- 2) Can be wasteful if Markov blanket is large:

$P(X_i | mb(X_i))$ won't change much (law of large numbers)



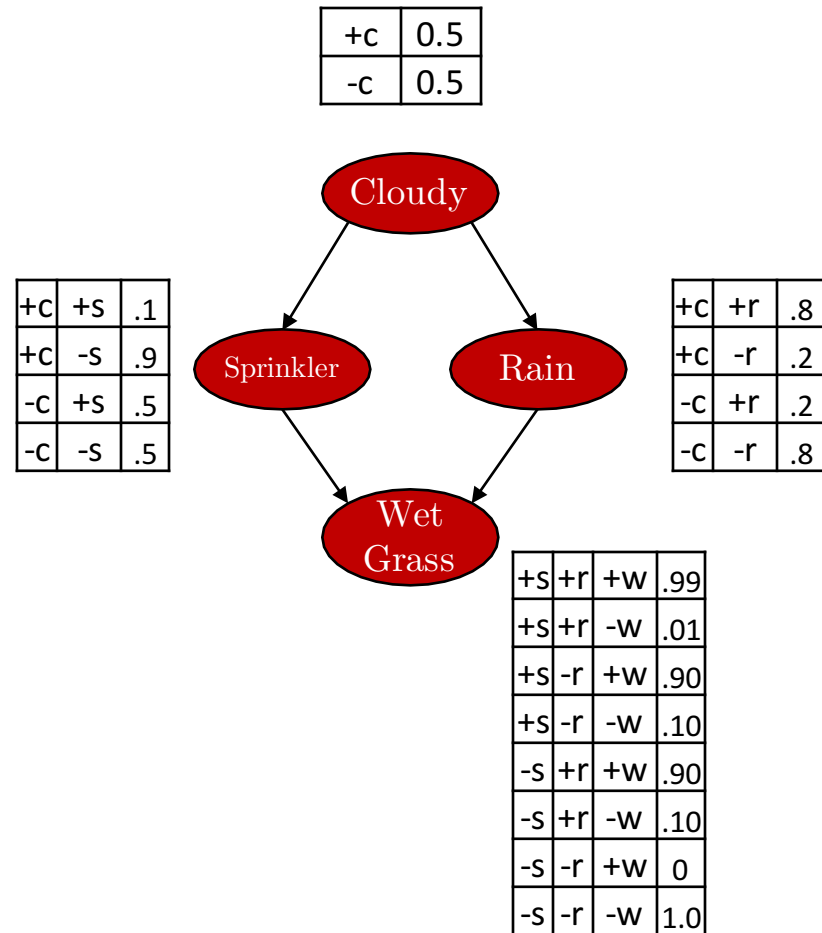
GIBBS SAMPLING EXAMPLE

- Want $\Pr(R|S=t, W=t)$
- Non-evidence variables are C & R
- Initialize randomly: C= t and R=f
- Initial state $(C,S,R,W) = [t,t,f,t]$
- Sample C given current values of its Markov Blanket



GIBBS SAMPLING EXAMPLE

- Want $\Pr(R|S=t, W=t)$
- Non-evidence variables are C & R
- Initialize randomly: C= t and R=f
- Initial state $(C,S,R,W) = [t,t,f,t]$
- Sample C given current values of its Markov Blanket
- Markov blanket is parents, children and children's parents: for C=S & R
- Sample C given $P(C|S=t, R=f)$
- First have to compute $P(C|S=t, R=f)$
- Use exact inference to do this



HOW DO WE COMPUTE $P(C \mid S=t, R=f)$? (RECALL OF EXACT INFERENCE)

- $P(C \mid S=t, R=f)$
- What is the probability $P(C=t \mid S=t, R=f)$?
 $= P(C=t, S=t, R=f) / (P(S=t, R=f))$

Proportional to $P(C=t, S=t, R=f)$

Use normalization trick, & compute the above for $C=t$ and $C=f$

$P(C=t, S=t, R=f) = P(C=t) P(S=t \mid C=t) P(R=f \mid C=t, S=t)$

product rule

$= P(C=t) P(S=t \mid C=t) P(R=f \mid C=t)$ (BN independencies)

$= 0.5 * 0.1 * 0.2 = 0.01$

$P(C=f, S=t, R=f) = P(C=f) P(S=t \mid C=f) P(R=f \mid C=f)$

$= 0.5 * 0.5 * 0.8 = 0.2$

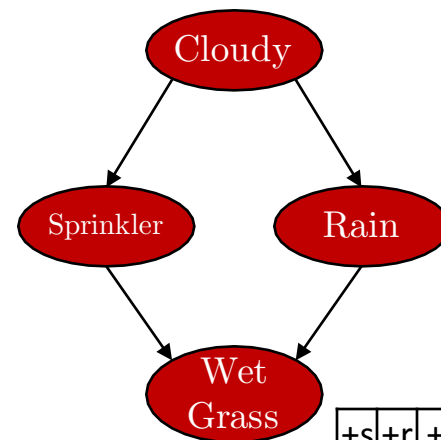
$(P(S=t, R=f))$ use sum rule $= P(C=f, S=t, R=f) + P(C=t, S=t, R=f)$

$P(C=t \mid S=t, R=f) = 0.01 / 0.21 = 0.0476$

$P(C=t \mid S=t, R=f) = 0.01 / 0.21 \sim 0.0476$

+c	0.5
-c	0.5

+c	+s	.1
+c	-s	.9
-c	+s	.5
-c	-s	.5



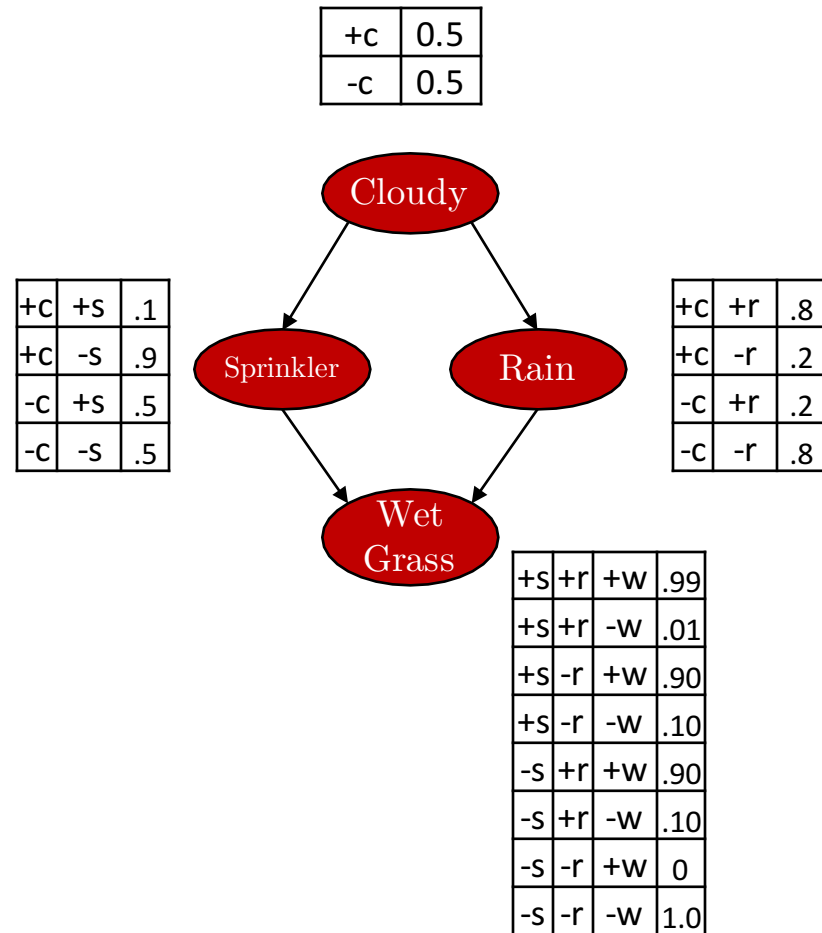
+c	+r	.8
+c	-r	.2
-c	+r	.2
-c	-r	.8

+s	+r	+w	.99
+s	+r	-w	.01
+s	-r	+w	.90
+s	-r	-w	.10
-s	+r	+w	.90
-s	+r	-w	.10
-s	-r	+w	0
-s	-r	-w	1.0



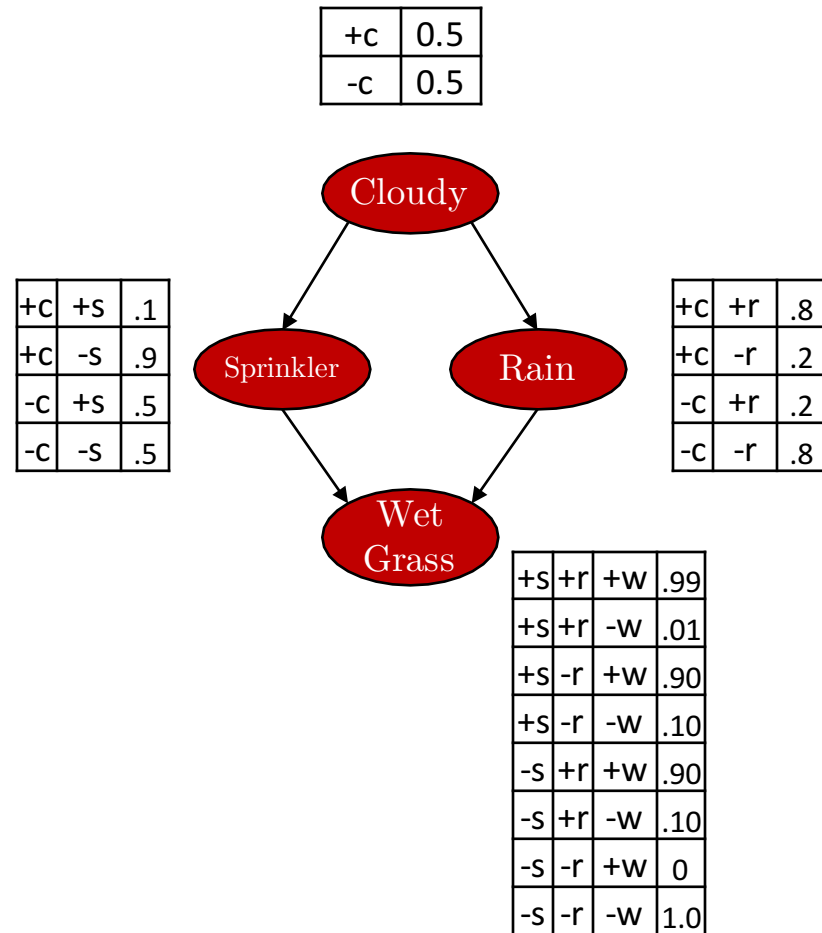
GIBBS SAMPLING EXAMPLE

- Want $\Pr(R|S=t, W=t)$
- Non-evidence variables are C & R
- Initialize randomly: $C = t$ and $R = f$
- Initial state $(C, S, R, W) = [t, t, f, t]$
- Sample C given current values of its Markov Blanket
- Markov blanket is parents, children and children's parents: for $C = S$ & R
- Exactly compute $P(C|S=t, R=f)$
- Sample C given $P(C|S=t, R=f)$
- Get $C = f$
- New state (f, t, f, t)



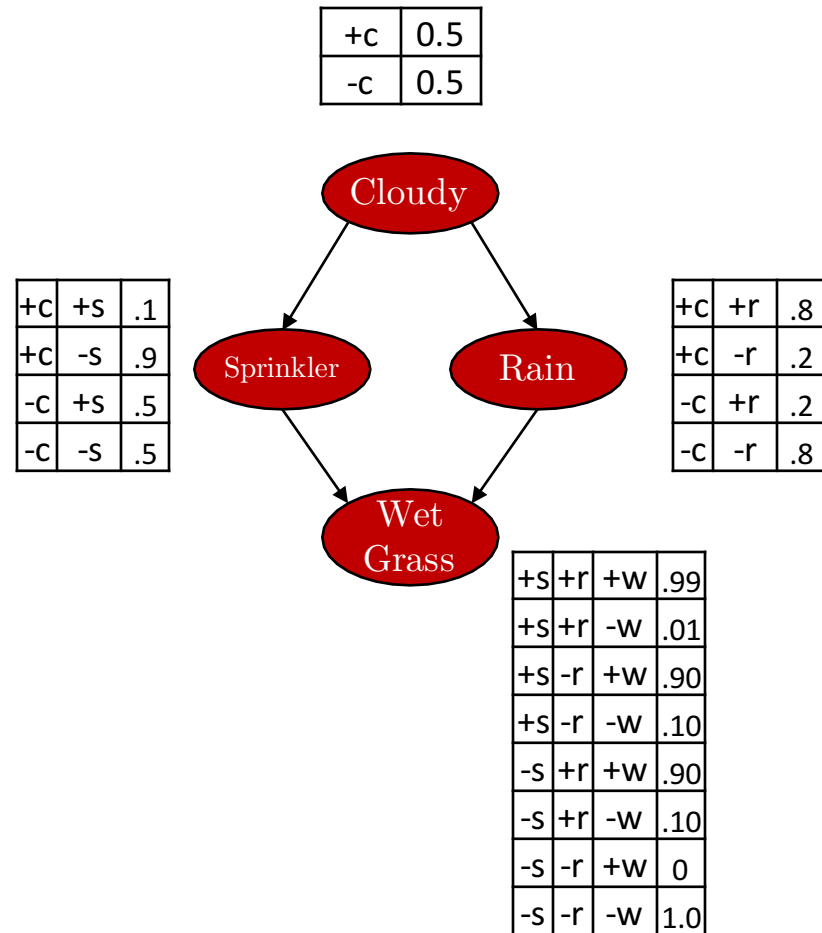
GIBBS SAMPLING EXAMPLE

- Want $\Pr(R|S=t, W=t)$
- Initialize non-evidence variables (C and R) randomly to t and f
- Initial state $(C, S, R, W) = [t, t, f, t]$
- Sample C given current values of its Markov Blanket, $p(C|S=t, R=f)$
- Suppose result is $C=f$
- New state (f, t, f, t)
- Sample Rain given its MB
- What is its Markov blanket?



GIBBS SAMPLING EXAMPLE

- Want $\Pr(R|S=t, W=t)$
- Initialize non-evidence variables (C and R) randomly to t and f
- Initial state $(C, S, R, W) = [t, t, f, t]$
- Sample C given current values of its Markov Blanket, $p(C|S=t, R=f)$
- Suppose result is $C=f$
- New state (f, t, f, t)
- Sample Rain given its MB, $p(R|C=f, S=t, W=t)$
- Suppose result is $R=t$
- New state (f, t, t, t)



POLL: GIBBS SAMPLING EX.

- Want $\Pr(R|S=t, W=t)$
- Initialize non-evidence variables (C and R) randomly to t and f
- Initial state $(C,S,R,W) = [t,t,f,t]$
- Current state (f,t,t,t)
- What is **not** a possible next state
 1. (f,t,t,t)
 2. (t,t,t,t)
 3. (f,t,f,t)
 4. (f,f,t,t)
 5. Not sure

