

15-381 / 781

**BAYESIAN NETS &
PROBABILISTIC
INFERENCE**

**EMMA BRUNSKILL (THIS TIME)
ARIEL PROCACCIA**

**WITH THANKS TO DAN KLEIN (BERKELEY), PERCY LIANG
(STANFORD) AND PAST 15-381 INSTRUCTORS FOR SOME SLIDE
CONTENT, AND RUSSELL & NORVIG**

WHAT YOU SHOULD KNOW

- Define probabilistic inference
- How to define a Bayes Net given a real example
- How joint can be used to answer any query
- Complexity of exact inference
- Approximation inference (direct, likelihood, Gibbs)
 - Be able to implement and run algorithm
 - Compare benefits and limitations of each



BAYESIAN NETWORK

- Compact representation of the joint distribution
- Conditional independence relationships explicit
 - Each var conditionally independent of all its non-descendants in the graph given the value of its parents



JOINT DISTRIBUTION EX.

- Variables: Cloudy, Sprinkler, Rain, Wet Grass
- Domain of each variable: 2 (true or false)
- Joint encodes probability of all combos of variables & values

+c	+s	+r	+w	.01
+c	+s	+r	-w	.01
+c	+s	-r	+w	.05
+c	+s	-r	-w	.1
+c	-s	+r	+w	#
+c	-s	+r	-w	#
+c	-s	-r	+w	#
+c	-s	-r	-w	#
-c	+s	+r	+w	#
-c	+s	+r	-w	#
-c	+s	-r	+w	#
-c	+s	-r	-w	#
-c	-s	+r	+w	#
-c	-s	+r	-w	#
-c	-s	-r	+w	#
-c	-s	-r	-w	#

$P(\text{Cloudy}=\text{false} \ \& \ \text{Sprinkler} = \text{true} \ \& \ \text{Rain} = \text{false} \ \& \ \text{WetGrass} = \text{True})$



JOINT AS PRODUCT OF CONDITIONALS (CHAIN RULE)

+c	+s	+r	+w	.01
+c	+s	+r	-w	.01
+c	+s	-r	+w	.05
+c	+s	-r	-w	.1
+c	-s	+r	+w	#
+c	-s	+r	-w	#
+c	-s	-r	+w	#
+c	-s	-r	-w	#
-c	+s	+r	+w	#
-c	+s	+r	-w	#
-c	+s	-r	+w	#
-c	+s	-r	-w	#
-c	-s	+r	+w	#
-c	-s	+r	-w	#
-c	-s	-r	+w	#
-c	-s	-r	-w	#

=

$P(\text{WetGrass}|\text{Cloudy},\text{Sprinkler},\text{Rain})^*$
 $P(\text{Rain}|\text{Cloudy},\text{Sprinkler})^*$
 $P(\text{Sprinkler}|\text{Cloudy})^*$
 $P(\text{Cloudy})$

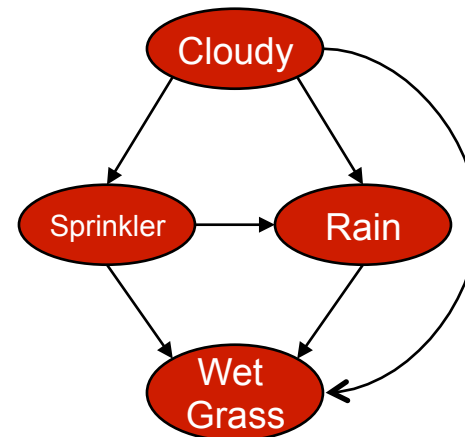


JOINT AS PRODUCT OF CONDITIONALS

+c	+s	+r	+w	.01
+c	+s	+r	-w	.01
+c	+s	-r	+w	.05
+c	+s	-r	-w	.1
+c	-s	+r	+w	#
+c	-s	+r	-w	#
+c	-s	-r	+w	#
+c	-s	-r	-w	#
-c	+s	+r	+w	#
-c	+s	+r	-w	#
-c	+s	-r	+w	#
-c	+s	-r	-w	#
-c	-s	+r	+w	#
-c	-s	+r	-w	#
-c	-s	-r	+w	#
-c	-s	-r	-w	#

=

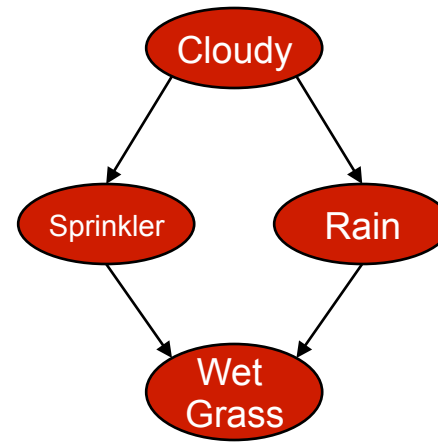
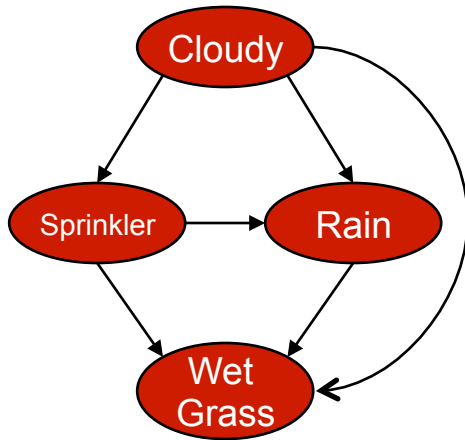
$P(\text{WetGrass}|\text{Cloudy},\text{Sprinkler},\text{Rain})^*$
 $P(\text{Rain}|\text{Cloudy},\text{Sprinkler})^*$
 $P(\text{Sprinkler}|\text{Cloudy})^*$
 $P(\text{Cloudy})$



...but there may be additional conditional independencies



WHAT IF SOME VARIABLES ARE CONDITIONALLY INDEP?

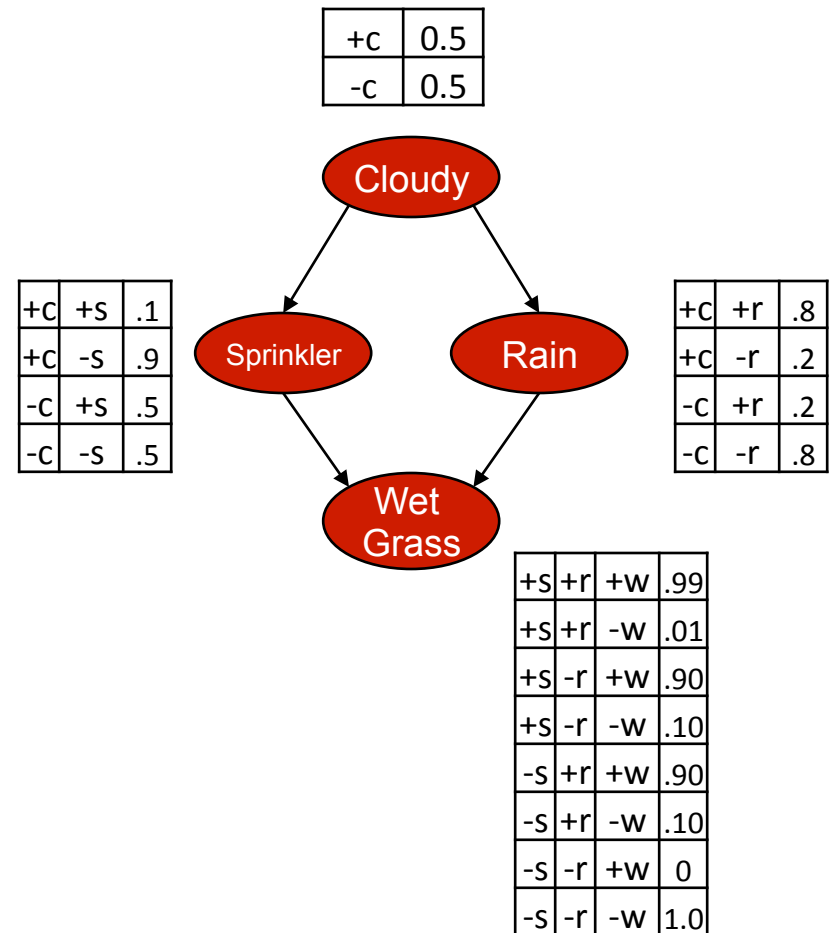


Explicitly shows any conditional independencies



CONDITIONAL INDEPENDENCIES

+c	+s	+r	+w	.01
+c	+s	+r	-w	.01
+c	+s	-r	+w	.05
+c	+s	-r	-w	.1
+c	-s	+r	+w	#
+c	-s	+r	-w	#
+c	-s	-r	+w	#
+c	-s	-r	-w	#
-c	+s	+r	+w	#
-c	+s	+r	-w	#
-c	+s	-r	+w	#
-c	+s	-r	-w	#
-c	-s	+r	+w	#
-c	-s	+r	-w	#
-c	-s	-r	+w	#
-c	-s	-r	-w	#



BAYESIAN NETWORK

- Compact representation of the joint distribution
- Conditional independence relationships explicit
- Still represents joint so can be used to answer any probabilistic query



PROBABILISTIC INFERENCE

- Compute probability of a **query** variable (or variables) taking on a value (or set of values) given some **evidence**
- $\Pr[Q \mid E_1=e_1, \dots, E_k=e_k]$



USING THE JOINT TO ANSWER QUERIES

- Joint distribution is sufficient to answer any probabilistic inference question involving variables described in joint
- Can take Bayes Net, construct full joint, and then look up entries where evidence variables take on specified values

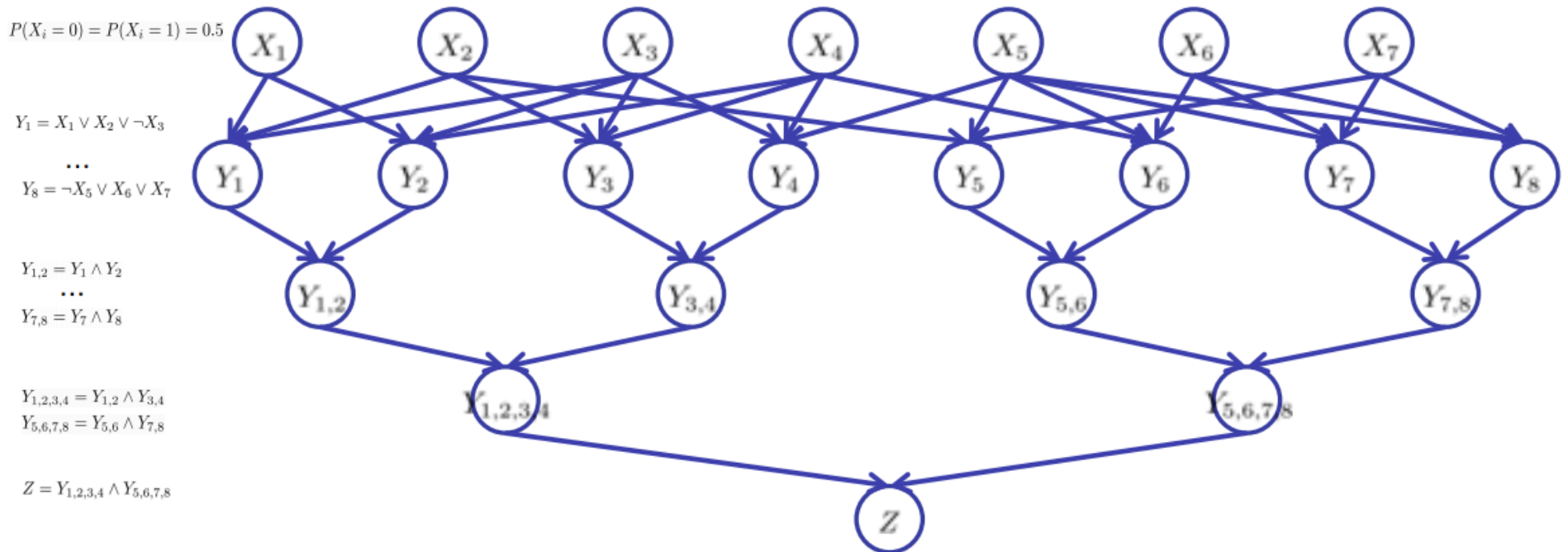


BUT CONSTRUCTING JOINT EXPENSIVE & EXACT INFERENCE IS NP-HARD

- 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.



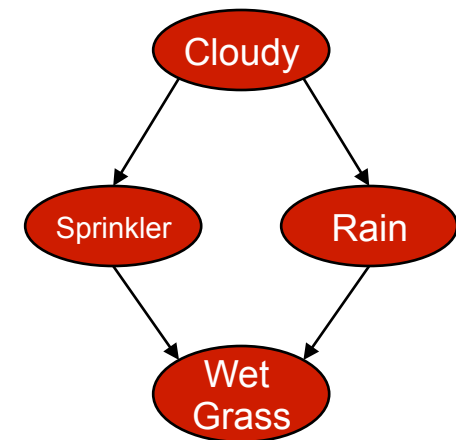
SOLN: APPROXIMATE INFERENCE

- Use samples to approximate posterior distribution $\Pr[Q \mid E_1=e_1, \dots, E_k=e_k]$
- Last time
 - Direct sampling
 - Likelihood weighting
- Today
 - Gibbs



POLL: WHICH ALGORITHM?

- Evidence: Cloudy=+c, Rain=+r
- Query variable: Sprinkler
- $P(\text{Sprinkler}|\text{Cloudy}=+c, \text{Rain}=+r)$
- Samples
 - +c,+s,+r,-w
 - +c,-s,-r,-w
 - +c,+s,-r,+w
 - +c,-s,+r,-w
- What algorithm could've generated these samples?
 - 1) Direct sampling
 - 2) Likelihood weighting
 - 3) Both
 - 4) No clue



DIRECT SAMPLING RECAP

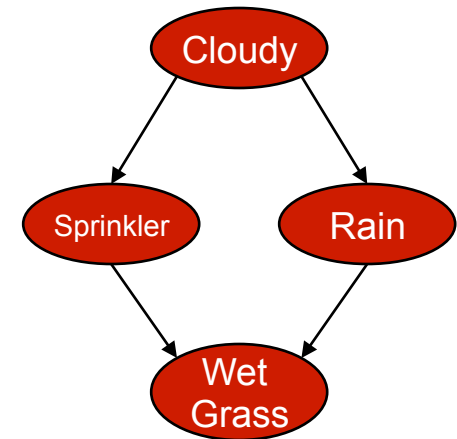
Algorithm:

1. Create a topological order of the variables in the Bayes Net



TOPOLOGICAL ORDER

- Any ordering in directed acyclic graph where a node can only appear after all of its ancestors in the graph
- E.g.
 - Cloudy, Sprinkler, Rain, WetGrass
 - Cloudy, Rain, Sprinkler, WetGrass



DIRECT SAMPLING RECAP

Algorithm:

1. Create a topological order of the variables in the Bayes Net
2. Sample each variable conditioned on the values of its parents
3. Use samples which match evidence variable values to estimate probability of query variable

e.g. $P(\text{Sprinkler}=+s|\text{Cloudy}=+c,\text{Rain}=+r) \sim \frac{\# \text{ samples with } +s,+c,+r}{\# \text{ samples with } +c,+r}$

- Consistent in limit of infinite samples
- Inefficient (why?)



CONSISTENCY

- In the limit of infinite samples, estimated $\Pr[Q \mid E_1=e_1, \dots, E_k=e_k]$ will converge to true posterior probability
- Desirable property (otherwise always have some error)



LIKELIHOOD WEIGHTING RECAP

1. Create array TotalWeights
 1. Initialize value of each array element to 0
2. For $j=1:N$
 1. $w_{tmp} = 1$
 2. Set evidence variables in sample $\mathbf{z}=\langle z_1, \dots, z_n \rangle$ to observed values
 3. For each variable z_i in topological order
 1. If x_i is an evidence variable
 1. $w_{tmp} = w_{tmp} * P(Z_i = e_i | \text{Parents}(Z) = \mathbf{x}(\text{Parents}(Z_i)))$
 2. Else
 1. Sample x_i conditioned on the values of its parents
 4. Update weight of resulting sample
 1. $\text{TotalWeights}[\mathbf{z}] = \text{TotalWeights}[\mathbf{z}] + w_{tmp}$

3. Use weights to compute probability of query variable

$$P(\text{Sprinkler}=+s | \text{Cloudy}=+c, \text{Rain}=+r) \sim \text{Sum}_{c,r,w} \text{TotalWeight}(+s, c, r, w) / \text{Sum}_{s,c,r,w} \text{TotalWeight}(s, c, r, w)$$



LW CONSISTENCY

- Probability of getting a sample (\mathbf{z}, \mathbf{e}) where \mathbf{z} is a set of values for the non-evidence variables and \mathbf{e} is the vals of evidence vars

Sampling
distribution for a
weighted sample
(WS)

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

- Is this the true posterior distribution $P(\mathbf{z}|\mathbf{e})$?
 - No, why?
 - Doesn't consider evidence that is not an ancestor...
 - Weights fix this!



WEIGHTED PROBABILITY

- Samples each non-evidence variable z according to

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

- Weight of a sample is

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

- Weighted probability of a sample is

$$\begin{aligned} S_{WS}(\mathbf{z}, \mathbf{e})w(\mathbf{z}, \mathbf{e}) &= \prod_{i=1}^l P(z_i | \text{parents}(Z_i)) \prod_{i=1}^m P(e_i | \text{parents}(E_i)) \\ &= P(\mathbf{z}, \mathbf{e}) \text{ From chain rule \& conditional indep} \end{aligned}$$



DOES LIKELIHOOD WEIGHTING PRODUCE CONSISTENT ESTIMATES? YES

$$P(X = x | e) = \frac{P(X = x, e)}{P(e)} \propto P(X = x, e)$$

X is query var(s)
E is evidence var(s)
Y is non-query vars

$$\tilde{P}(X = x | e) \propto \tilde{P}(X = x, e) = \sum_y N_{WS}(x, y, e) w(x, y, e)$$

of samples where query variables=x, non-query=y, Evidence=e

$$\approx \sum_y n * S_{WS}(x, y, e) w(x, y, e)$$

$$= \sum_y P(x, y, e)$$

$$= P(x, e)$$

as # samples n \rightarrow infinity



EXAMPLE

- When sampling S and R the evidence $W=t$ is ignored
 - 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 $R=t$?
 - = 1/3



LIMITATIONS OF LIKELIHOOD WEIGHTING

- Poor performance if evidence vars occur later in ordering
- Why?
- Not being used to influence samples!
- Yields samples with low weights



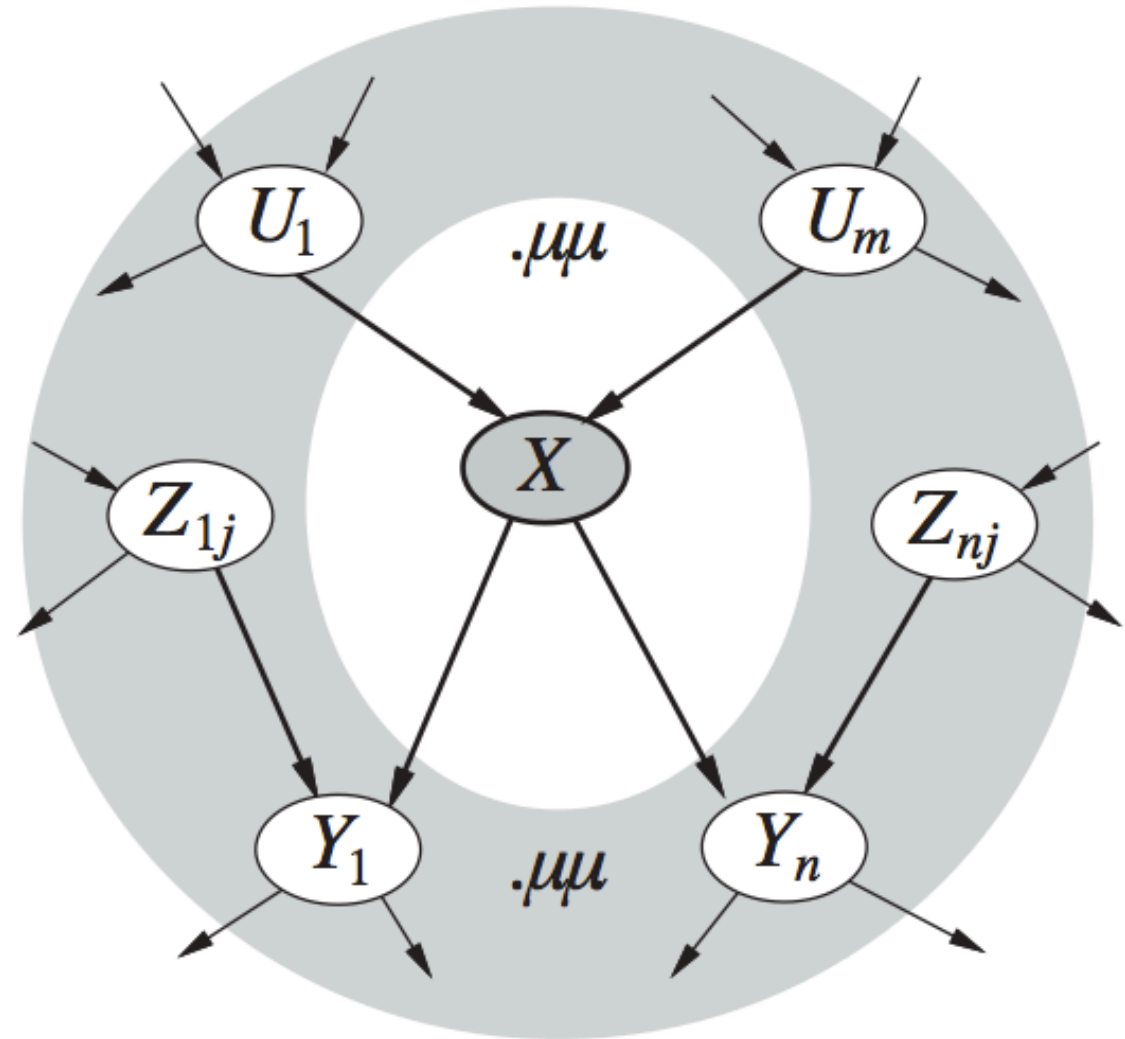
MARKOV CHAIN MONTE CARLO METHODS

- Prior methods generate each new sample from scratch
- MCMC generate each new sample by making a random change to preceding sample
- Can view algorithm as being in a particular state (assignment of values to each variable)



REVIEW: MARKOV BLANKET

- Markov blanket
 - Parents
 - Children
 - Children's parents
- Variable conditionally independent of all other nodes given its Markov Blanket



GIBBS SAMPLING: COMPUTE $P(X|e)$

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}

$mb(Z_i)$ = Markov Blanket of Z_i

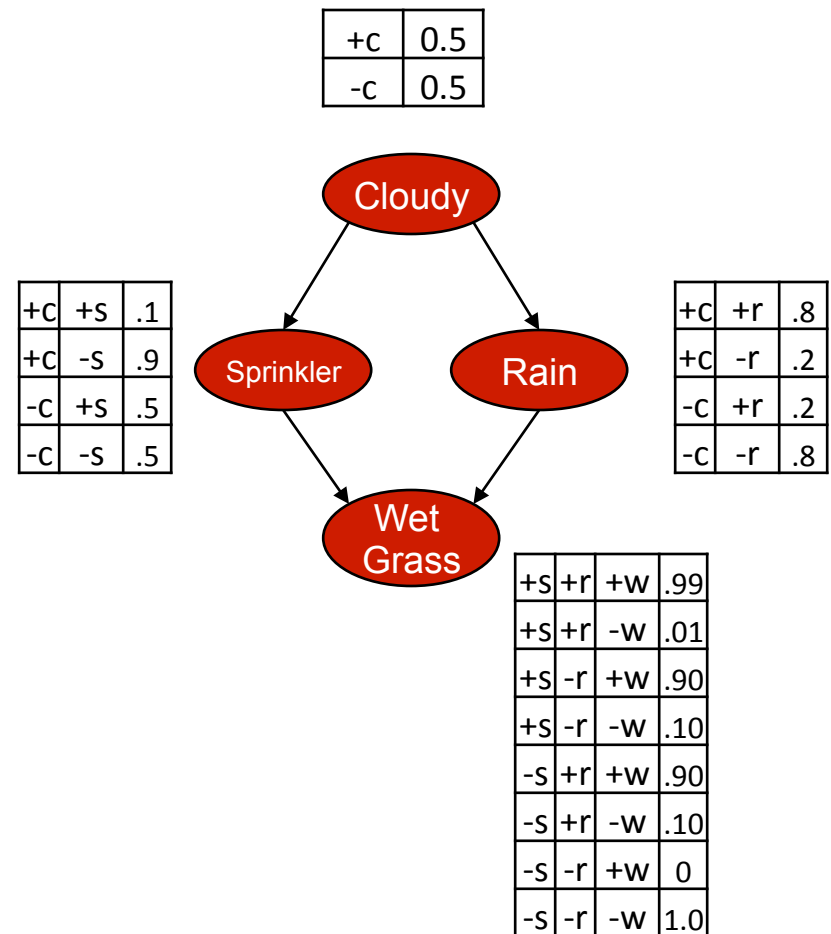


from Russell & Norvig

Carnegie Mellon University 27

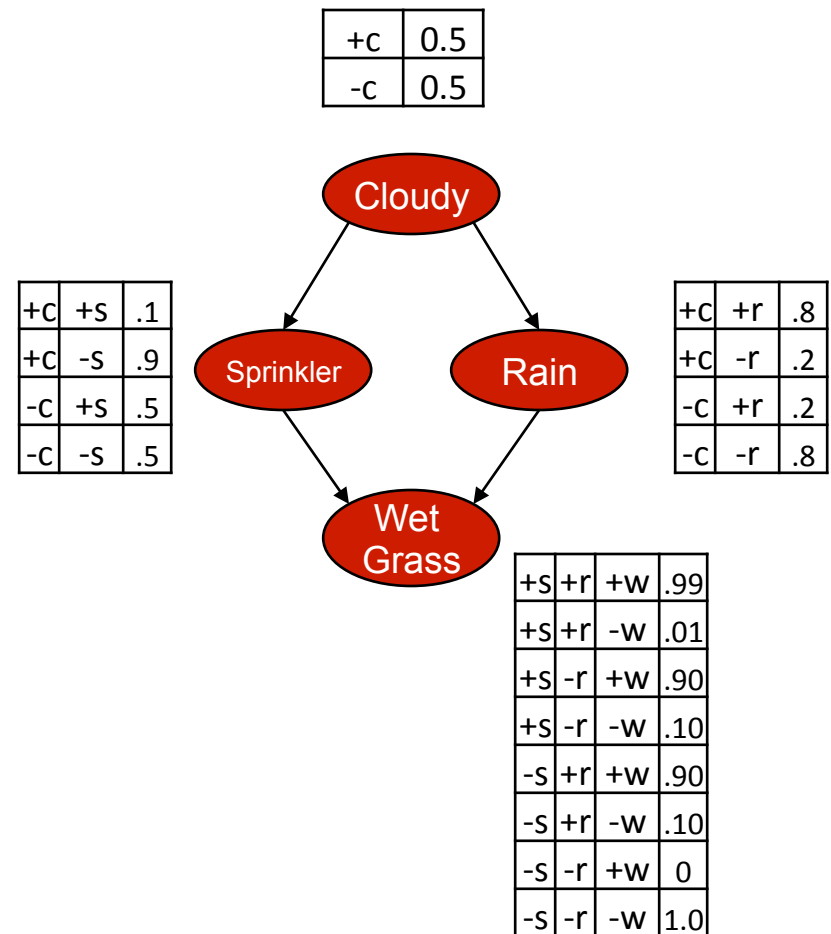
GIBBS SAMPLING EXAMPLE

- 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



EXERCISE: COMPUTE $P(C=T|S=T,R=F)$?

- Quick refresher

- Sum rule

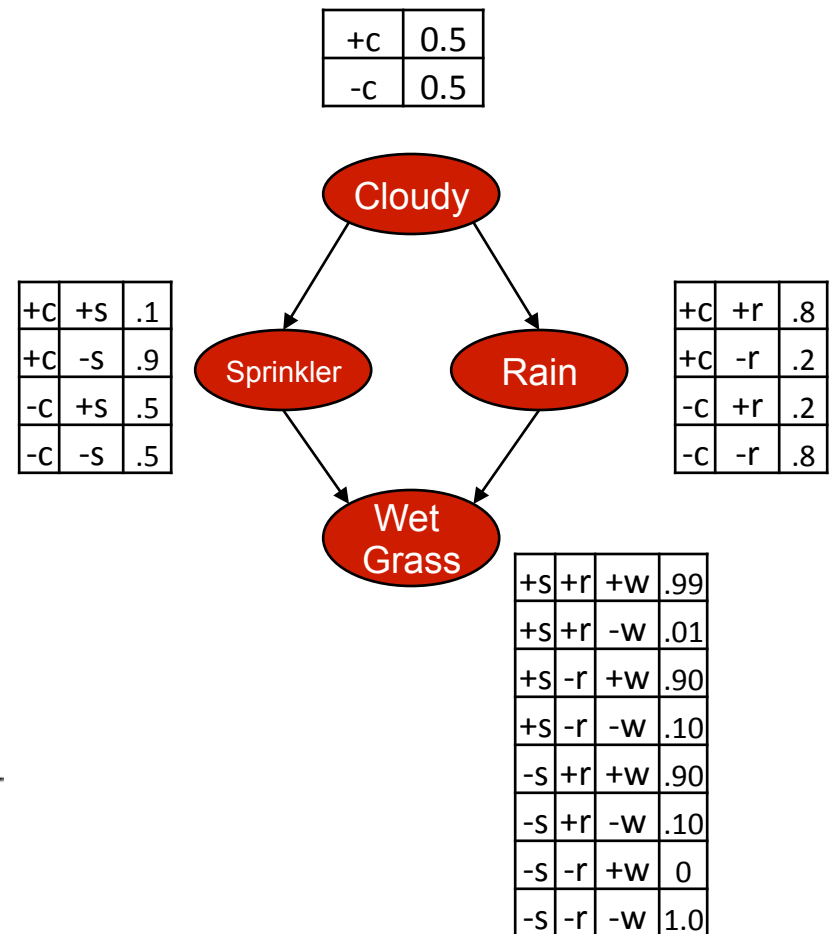
$$p(X) = \sum_Y p(X, Y)$$

- Product/Chain rule

$$p(X, Y) = p(Y|X)p(X)$$

- Bayes rule

$$p(X|Y) = \frac{p(Y|X)p(X)}{p(Y)}$$



EXACT INFERENCE EXERCISE

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

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

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

$P(C=t, S=t, R=f) = P(C=t) P(S=t|C=t) P(R=f | C=t, S=t)$ product rule

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

$$= 0.5 * 0.1 * 0.2 = 0.01$$

$P(C=f, S=t, R=f) = P(C=f) P(S=t|C=f) P(R=f|C=f)$

$$= 0.5 * 0.5 * 0.8 = 0.2$$

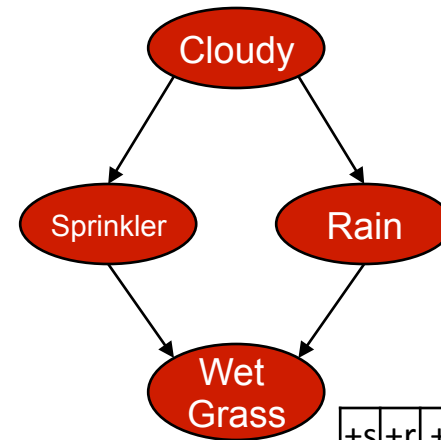
$(P(S=t,R=f))$ use sum rule $= P(C=f, S=t, R=f) + P(C=t, S=t, R=f)$

$$P(C = t | S=t, R= f) = 0.01 / 0.21$$

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

+c	0.5
-c	0.5

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



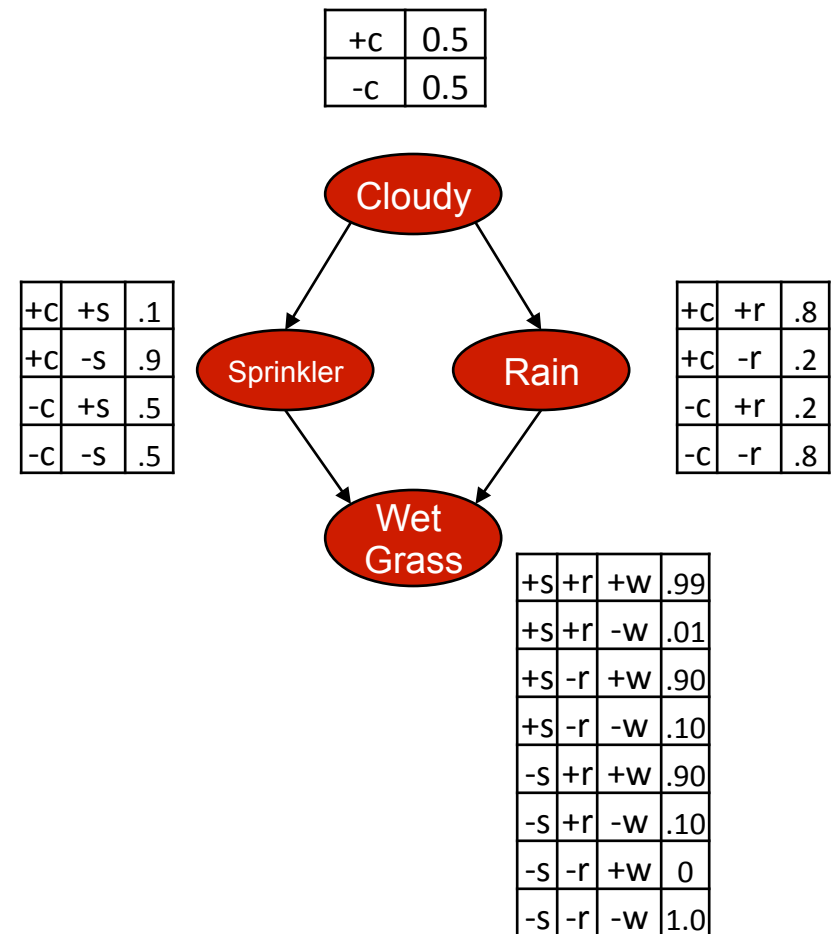
+c	+r	.8
+c	-r	.2
-c	+r	.2
-c	-r	.8

+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



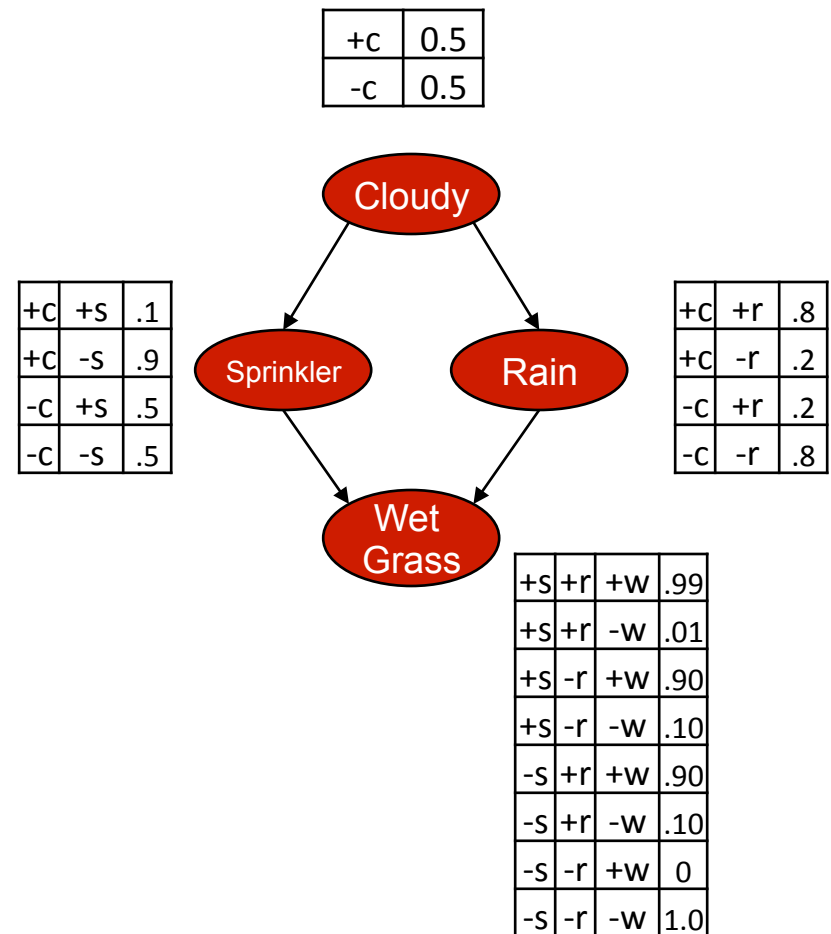
GIBBS SAMPLING EXAMPLE

- 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)



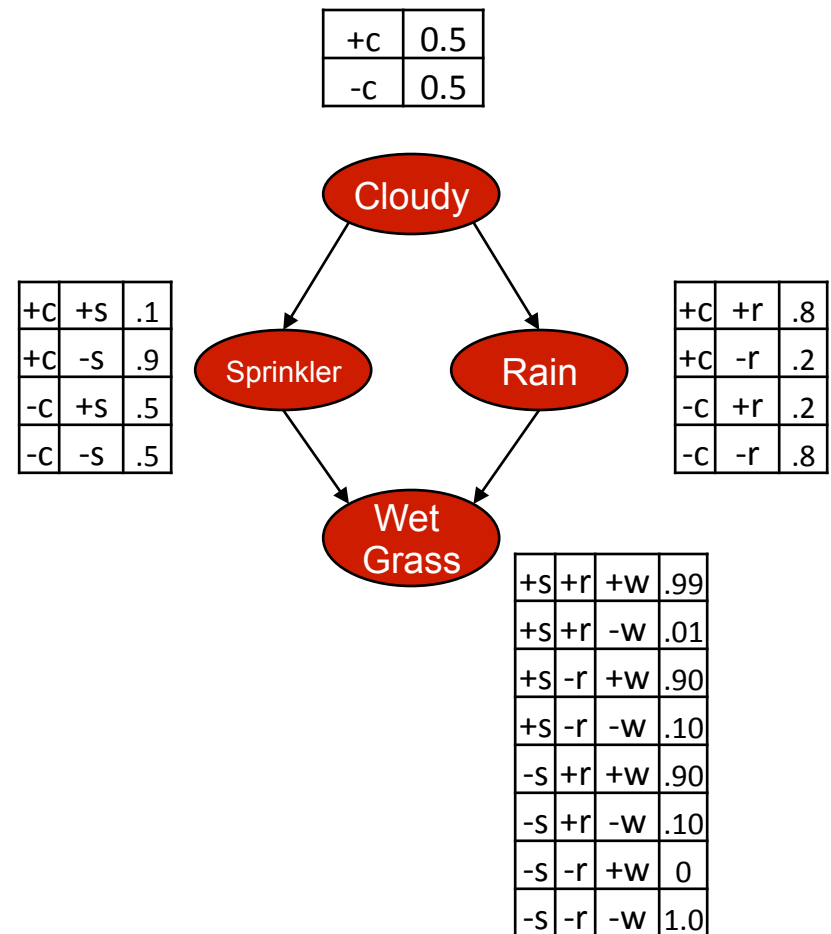
GIBBS SAMPLING EXAMPLE

- 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?**



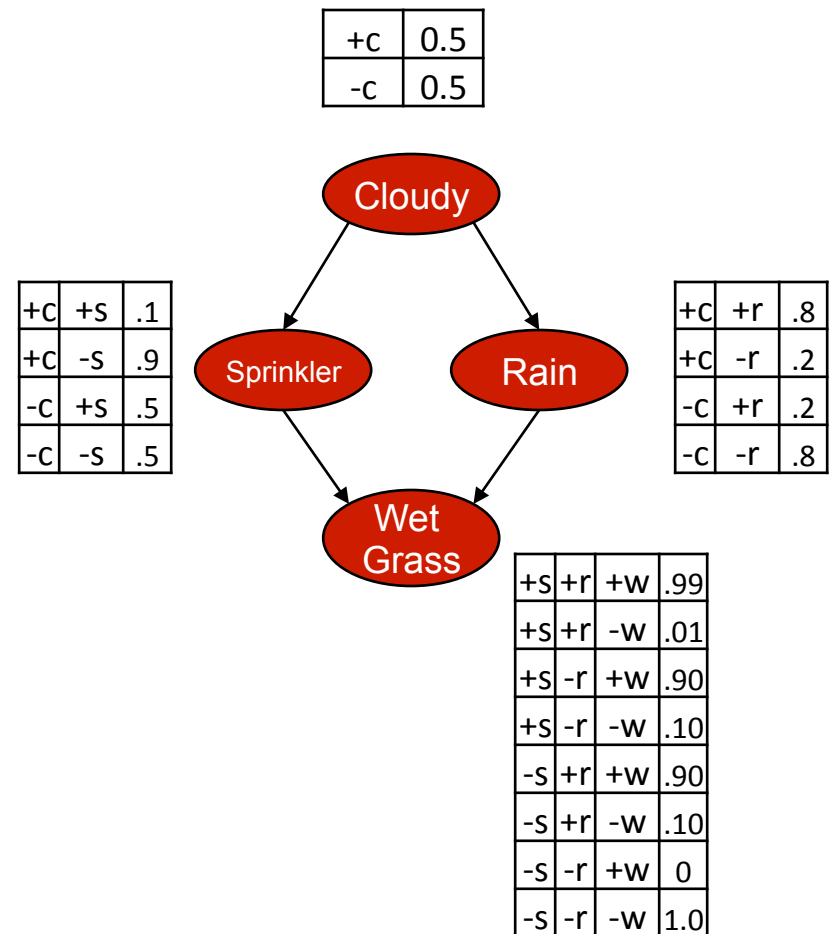
GIBBS SAMPLING EXAMPLE

- 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) (inconsistent w/evid)**
 5. Not sure



GIBBS SAMPLING

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})

This involve
inference!



$mb(Z_i)$ = Markov Blanket of Z_i



from Russell & Norvig

Carnegie Mellon University 36

POLL

ARE GIBBS SAMPLES INDEPENDENT?

1. YES 2. NO 3. NOT SURE

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})

$mb(Z_i) = \text{Markov Blanket of } Z_i$



MARKOV BLANKET SAMPLING

- Want to show $P(Z_i | \text{mb}(Z_i))$ is same as $P(Z_i | \text{all other variables})$
 - Implies conditional independence of Z_i from rest of network given its Markov Blanket
- Derive equation for computing $P(Z_i | \text{mb}(Z_i))$



PROBABILITY GIVEN MARKOV BLANKET

$$P(x'_i | mb(X_i)) = \alpha P(x'_i | parents(X_i)) \times \prod_{Y_j \in Children(X_i)} P(y_j | parents(Y_j))$$



WHY IS GIBBS CONSISTENT?

- Sampling process settles into a stationary distribution where long-term fraction of time spent in each state is exactly equal to posterior probability
 - → Implies that if draw enough samples from this stationary distribution, will get consistent estimate because sampling from true posterior



MARKOV CHAIN

- Let $P(\mathbf{x} \rightarrow \mathbf{x}')$ be probability the sampling process makes a transition from \mathbf{x} (some state) to \mathbf{x}' (some other state)
 - E.g. $(t,t,f,t) \rightarrow (t,f,f,t)$
- Run sampling for t steps
- $P_t(\mathbf{x})$ is probability system is in state \mathbf{x} at time t
- Next state $P_{t+1}(\mathbf{x}') = \text{Sum}_{\mathbf{x}} P_t(\mathbf{x}) P(\mathbf{x} \rightarrow \mathbf{x}')$



STATIONARY DISTRIBUTION

- Let $P(\mathbf{x} \rightarrow \mathbf{x}')$ be probability the process makes a transition from \mathbf{x} to \mathbf{x}'
- $P_t(\mathbf{x})$ is probability system is in state \mathbf{x} at time t
- $P_{t+1}(\mathbf{x}') = \text{Sum}_{\mathbf{x}} P_t(\mathbf{x}) P(\mathbf{x} \rightarrow \mathbf{x}')$
- Reached stationary distribution if $P_{t+1}(\mathbf{x}') = P_t(\mathbf{x})$
- Call stationary distribution π
 - Must satisfy $\pi(\mathbf{x}') = \sum_{\mathbf{x}} \pi(\mathbf{x}) P(\mathbf{x} \rightarrow \mathbf{x}')$ for all \mathbf{x}'
- If $P(\mathbf{x} \rightarrow \mathbf{x}')$ is ergodic, exactly one such π for any given $P(\mathbf{x} \rightarrow \mathbf{x}')$



DETAILED BALANCE

- Let $P(\mathbf{x} \rightarrow \mathbf{x}')$ be probability the process makes a transition from \mathbf{x} to \mathbf{x}'
- $P_t(\mathbf{x})$ is probability system is in state \mathbf{x} at time t
- Stationary distribution π
 - Satisfies $\pi(\mathbf{x}') = \sum_{\mathbf{x}} \pi(\mathbf{x}) P(\mathbf{x} \rightarrow \mathbf{x}')$ for all \mathbf{x}'
- Detailed balance: inflow = outflow
- $\pi(\mathbf{x}) P(\mathbf{x} \rightarrow \mathbf{x}') = \pi(\mathbf{x}') P(\mathbf{x}' \rightarrow \mathbf{x})$ for all \mathbf{x}, \mathbf{x}'

$$\sum_x \pi(x') P(x \rightarrow x') = \sum_x \pi(x') P(x' \rightarrow x) = \pi(x')$$



- Proof on board



PROVING GIBBS SAMPLES FROM TRUE POSTERIOR

- General Gibbs: sample the value of a new variable conditioned on all the other variables
- Can prove this version of Gibbs satisfies detailed balance equation with stationary distribution of $P(X|e)$
- Then use prior result that sampling conditioned on all variables is equivalent to sampling given Markov Blanket for Bayes Nets
- See text for recap



GIBBS SAMPLING

- Samples are valid once reach stationary distribution
- When do we reach stationary distribution?
 - Unclear...



WHAT YOU SHOULD KNOW

- Define probabilistic inference
- How to define a Bayes Net given a real example
- How joint can be used to answer any query
- Complexity of exact inference
- Approximation inference (direct, likelihood, Gibbs)
 - Be able to implement and run algorithm
 - Compare benefits and limitations of each

