

# Intelligent Autonomous Agents and Cognitive Robotics

## Topic 5: Bayesian Networks

Ralf Möller, Rainer Marrone  
Hamburg University of Technology

## Uncertainty in prior knowledge

- Diagnosis:
  - ♦  $\text{Toothache} \Rightarrow \text{Cavity} \vee \text{GumProblem} \vee \text{Abscess} \vee \dots$
  - $\text{RootInfection} \vee \dots \vee \text{Cavity} \Rightarrow \text{Toothache}$
- The connection between toothaches and cavity is just not a logical consequence. For medical diagnosis logic does not seem to be appropriate.

## Probability

Probabilistic assertions **summarize** effects of

- **laziness:**  
It is too much work to list the complete set of antecedents or consequents needed to ensure an exceptionless rule and too hard to use such rules
- **theoretical ignorance:**  
no complete theory, e.g., medical science has no complete theory for the domain.
- **practical ignorance:**  
lack of relevant facts, initial conditions, tests, etc.

3

## Making decisions under uncertainty

Suppose I believe the following:

$$\begin{aligned} P(A_{25} \text{ gets me there on time} \mid \dots) &= 0.04 \\ P(A_{90} \text{ gets me there on time} \mid \dots) &= 0.70 \\ P(A_{120} \text{ gets me there on time} \mid \dots) &= 0.95 \\ P(A_{1440=24h} \text{ gets me there on time} \mid \dots) &= 0.9999 \end{aligned}$$

- Which action to choose?

Depends on my **preferences** for missing flight vs. time spent waiting, etc.

- ♦ **Utility theory** is used to represent and use preferences
- ♦ **Decision theory** = **probability theory** + **utility theory**

↓  
Later in this lecture

4

## Example world

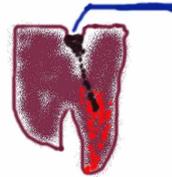
**Example:** *Dentist problem with four variables:*

*Toothache* (I have a toothache)

*Cavity* (I have a cavity)

*Catch* (steel probe catches in my tooth)

*Weather* (*sunny,rainy,cloudy,snow*)



5

## Prior probability

- **Prior** or **unconditional probabilities** of propositions

e.g.,  $P(\text{Cavity} = \text{true}) = 0.1$  and  $P(\text{Weather} = \text{sunny}) = 0.72$   
correspond to belief prior to arrival of any (new) evidence

- **Probability distribution**

gives values for all possible assignments

(*sunny,rainy,cloudy,snow*):

$P(\text{Weather}) = \langle 0.72, 0.1, 0.08, 0.1 \rangle$

(**normalized**, i.e., sums to 1 because one must be the case)

6

## Full joint probability distribution

- **Joint probability distribution** for a set of random variables gives the probability of every atomic event on those random variables

$P(\text{Weather}, \text{Cavity}) = a 4 \times 2$  matrix of values:

<i>Weather</i> =	sunny	rainy	cloudy	snow	
<i>Cavity</i> = true	0.144	0.02	0.016	0.02	= 0.2
<i>Cavity</i> = false	0.576	0.08	0.064	0.08	= 0.8
					= 1.0

- Full joint probability distribution: all random variables involved
  - ♦  $P(\text{Toothache}, \text{Catch}, \text{Cavity}, \text{Weather})$
- **Every question about a domain can be answered by the full joint distribution**

7

## Conditional probability

- **Conditional or posterior probabilities** (after having received some information)  
e.g.,  $P(\text{cavity} \mid \text{toothache}) = 0.8$
- Definition of **conditional probability** (in terms of uncond. prob.):  
 $P(a \mid b) = P(a \wedge b) / P(b)$  if  $P(b) > 0$
- **Product rule** gives an alternative formulation ( $\wedge$  is commutative):  
 $P(a \wedge b) = P(a \mid b) P(b) = P(b \mid a) P(a)$
- **Chain rule** is derived by successive application of product rule:

$$\begin{aligned}
 P(X_1, \dots, X_n) &= P(X_1, \dots, X_{n-1}) P(X_n \mid X_1, \dots, X_{n-1}) \\
 &= P(X_1, \dots, X_{n-2}) P(X_{n-1} \mid X_1, \dots, X_{n-2}) P(X_n \mid X_1, \dots, X_{n-1}) \\
 &= \prod_{i=1}^n P(X_i \mid X_1, \dots, X_{i-1})
 \end{aligned}$$

8

## Bayes rule

**Product rule:**  $P(x, y) = P(x|y)P(y) = P(y|x)P(x)$



$$P(\text{cause}|\text{effect}) = \frac{P(\text{effect}|\text{cause}) * P(\text{cause})}{P(\text{effect})} = \frac{\text{likelihood} * \text{prior}}{\text{evidence}}$$

$$P(X|Y) = \alpha P(Y|X)P(X)$$

9

## Inference by enumeration

- Start with the joint probability distribution:

	<i>toothache</i>		$\neg$ <i>toothache</i>	
	<i>catch</i>	$\neg$ <i>catch</i>	<i>catch</i>	$\neg$ <i>catch</i>
<i>cavity</i>	<b>.108</b>	<b>.012</b>	<b>.072</b>	<b>.008</b>
$\neg$ <i>cavity</i>	<b>.016</b>	<b>.064</b>	<b>.144</b>	<b>.576</b>

- For any proposition  $\phi$ , sum the atomic events where it is true:  $P(\phi) = \sum_{\omega \in \phi} P(\omega)$

10

## Inference by enumeration

- Start with the joint probability distribution:

	<i>toothache</i>		$\neg$ <i>toothache</i>	
	<i>catch</i>	$\neg$ <i>catch</i>	<i>catch</i>	$\neg$ <i>catch</i>
<i>cavity</i>	<b>.108</b>	<b>.012</b>	<b>.072</b>	<b>.008</b>
$\neg$ <i>cavity</i>	<b>.016</b>	<b>.064</b>	<b>.144</b>	<b>.576</b>

$$P(\text{cavity} \vee \text{toothache}) = 0.108 + 0.012 + 0.072 + 0.008 + 0.016 + 0.064 = 0.28$$

11

## Inference by enumeration

- Start with the joint probability distribution:

	<i>toothache</i>		$\neg$ <i>toothache</i>	
	<i>catch</i>	$\neg$ <i>catch</i>	<i>catch</i>	$\neg$ <i>catch</i>
<i>cavity</i>	<b>.108</b>	<b>.012</b>	<b>.072</b>	<b>.008</b>
$\neg$ <i>cavity</i>	<b>.016</b>	<b>.064</b>	<b>.144</b>	<b>.576</b>

- Can also compute conditional probabilities:

$$\begin{aligned}
 P(\neg \text{cavity} \mid \text{toothache}) &= \frac{P(\neg \text{cavity} \wedge \text{toothache})}{P(\text{toothache})} \\
 &\stackrel{\text{Product rule}}{=} \frac{0.016 + 0.064}{0.108 + 0.012 + 0.016 + 0.064} \\
 &= 0.4
 \end{aligned}$$

12

## Normalization

	toothache		$\neg$ toothache	
	catch	$\neg$ catch	catch	$\neg$ catch
cavity	.108	.012	.072	.008
$\neg$ cavity	.016	.064	.144	.576

- Denominator  $P(\mathbf{z})$  (or  $P(\text{toothache})$  in the example before) can be viewed as a **normalization constant**  $\alpha$

$$\begin{aligned}
 P(\text{Cavity} \mid \text{toothache}) &= P(\text{Cavity}, \text{toothache}) / P(\text{toothache}) \\
 &= \alpha P(\text{Cavity}, \text{toothache}) \\
 &= \alpha [P(\text{Cavity}, \text{toothache}, \text{catch}) + P(\text{Cavity}, \text{toothache}, \neg \text{catch})] \\
 &= \alpha [\langle 0.108, 0.016 \rangle + \langle 0.012, 0.064 \rangle] \\
 &= \alpha \langle 0.12, 0.08 \rangle = \langle 0.6, 0.4 \rangle
 \end{aligned}$$

$$\begin{aligned}
 \alpha * (0,12 + 0,08) &= 1 \\
 \alpha &= 1 / 0,2 = 5 \\
 5 * 0,12 &= 0,6 \\
 5 * 0,08 &= 0,4
 \end{aligned}$$

General idea: compute distribution on query variable by fixing **evidence variables** (toothache) and summing over **hidden variables** (Catch)

13

## General inference procedure

Typically, we are interested in the posterior joint distribution of the **query variables**  $Y$  given specific values  $\mathbf{e}$  for the **evidence variables**  $E$ .  $X$  are all variables of the modeled world.

Let the **hidden variables** be  $H = X - Y - E$  then the required summation of joint entries is done by summing out the hidden variables:

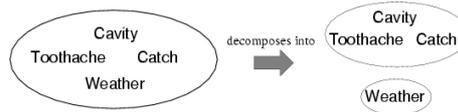
$$P(Y \mid E = \mathbf{e}) = \alpha P(Y, E = \mathbf{e}) = \alpha \sum_{\mathbf{h}} P(Y, E = \mathbf{e}, H = \mathbf{h})$$

- The terms in the summation are joint entries because  $Y$ ,  $E$  and  $H$  together exhaust the set of random variables ( $X$ )
- Obvious problems:
  - Space complexity  $O(d^n)$  to store the joint distribution where  $d$  is the largest arity and  $n$  denotes the number of random variables
  - Worst-case time complexity  $O(d^n)$
  - How to find the numbers for  $O(d^n)$  entries?

14

# Independence

- $A$  and  $B$  are independent iff  
 $P(A|B) = P(A)$  or  $P(B|A) = P(B)$  or  $P(A, B) = P(A) P(B)$



$$P(\text{Toothache}, \text{Catch}, \text{Cavity}, \text{Weather}) = P(\text{Toothache}, \text{Catch}, \text{Cavity}) P(\text{Weather})$$

- 32 entries table can be constructed from 8 and 4 entries;
- Absolute independence powerful but rare
- How can we check whether we have independent variables in the full joint?

15

# Example #1

Bread	Bagels	Butter	$p(r,a,u)$
0	0	0	0.24
0	0	1	0.06
0	1	0	0.12
0	1	1	0.08
1	0	0	0.12
1	0	1	0.18
1	1	0	0.04
1	1	1	0.16

Bread	$p(r)$
0	
1	

$$P(a,u)=P(a)P(u)?$$

$$P(r,a)=P(r)P(a)?$$

16

### Example #1

Butter	p(u)
0	0.52
1	0.48

Bagels	p(a)
0	0.6
1	0.4

Bread	p(r)
0	0.5
1	0.5

Bread	Bagels	Butter	p(r,a,u)
0	0	0	0.24
0	0	1	0.06
0	1	0	0.12
0	1	1	0.08
1	0	0	0.12
1	0	1	0.18
1	1	0	0.04
1	1	1	0.16

Bagels	Butter	p(a,u)
0	0	?
0	1	
1	0	
1	1	

$P(a,u)=P(a)P(u)?$   $P(r,a)=P(r)P(a)?$

17

### Example #1

Butter	p(u)
0	0.52
1	0.48

Bagels	p(a)
0	0.6
1	0.4

Bread	p(r)
0	0.5
1	0.5

Bread	Bagels	Butter	p(r,a,u)
0	0	0	0.24
0	0	1	0.06
0	1	0	0.12
0	1	1	0.08
1	0	0	0.12
1	0	1	0.18
1	1	0	0.04
1	1	1	0.16

Bagels	Butter	p(a,u)
0	0	0.36
0	1	0.24
1	0	0.16
1	1	0.24

$\neq 0.52 \cdot 0.6 = 0.312$

Bread	Bagels	p(r,a)
0	0	0.3
0	1	0.2
1	0	0.3
1	1	0.2

$P(a,u)=P(a)P(u)?$  **NO**  $P(r,a)=P(r)P(a)?$  **YES**

18

## Conditional independence

- $P(\text{Toothache}, \text{Cavity}, \text{Catch})$  has  $2^3 - 1 = 7$  independent entries
- If I have a cavity, the probability that the probe catches in doesn't depend on whether I have a toothache:
  - (1)  $P(\text{catch} \mid \text{toothache}, \text{cavity}) = P(\text{catch} \mid \text{cavity})$
  - (2)  $P(\text{catch} \mid \text{toothache}, \neg\text{cavity}) = P(\text{catch} \mid \neg\text{cavity})$
- Catch is **conditionally independent** of Toothache given Cavity:  
 $P(\text{Catch} \mid \text{Toothache}, \text{Cavity}) = P(\text{Catch} \mid \text{Cavity})$
- Equivalent statements:  
 $P(\text{Toothache} \mid \text{Catch}, \text{Cavity}) = P(\text{Toothache} \mid \text{Cavity})$   
 $P(\text{Toothache}, \text{Catch} \mid \text{Cavity}) = P(\text{Toothache} \mid \text{Cavity}) P(\text{Catch} \mid \text{Cavity})$

19

## Conditional independence contd.

- Write out full joint distribution using chain rule:  

$$P(\text{Toothache}, \text{Catch}, \text{Cavity})$$

$$= P(\text{Toothache} \mid \text{Catch}, \text{Cavity}) P(\text{Catch}, \text{Cavity})$$

$$= P(\text{Toothache} \mid \text{Catch}, \text{Cavity}) P(\text{Catch} \mid \text{Cavity}) P(\text{Cavity})$$
 conditional independence  

$$= P(\text{Toothache} \mid \text{Cavity}) P(\text{Catch} \mid \text{Cavity}) P(\text{Cavity})$$
- i.e.,  $2 + 2 + 1 = 5$  independent numbers
- In most cases, the use of conditional independence reduces the size of the representation of the joint distribution from exponential in  $n$  to linear in  $n$ .
- Conditional independence is our most basic and robust form of knowledge about uncertain environments.

20

## Car Example

- Three variables:
  - ♦ Gas, Battery, Starts
- $P(\text{Battery}|\text{Gas}) = P(\text{Battery})$   
Gas and Battery are independent
- $P(\text{Battery}|\text{Gas}, \text{Starts}) \neq P(\text{Battery}|\text{Starts})$   
*Gas and Battery are not independent given Starts*
- Independence does not imply conditional independence.
- Conditional independence does not imply independence

21

## Question

- How can we make use of
    - ♦ independence
    - ♦ and conditional independence
- Need a model that can express this

22

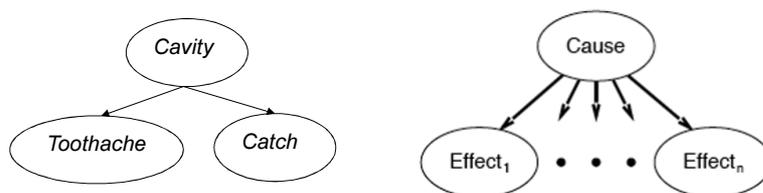
## Bayesian networks

- A simple, graphical notation for **conditional independence assertions** and hence for compact specification of the full joint distributions
- Syntax:
  - ♦ a set of nodes, one per variable
  - ♦ a directed, acyclic graph (link  $\approx$  "directly influences")
  - ♦ a conditional distribution for each node given its parents:  

$$P(X_i | \text{Parents}(X_i))$$
- In the simplest case, conditional distribution represented as a **conditional probability table (CPT)** giving the distribution over  $X_i$  for each combination of parent values

23

## Simplest Bayesian Network



$$P(\text{Cause} | \text{Effect}_1, \text{Effect}_2, \dots) = \alpha P(\text{Cause}) \prod_i P(\text{Effect}_i | \text{Cause})$$

- also called Naïve Bayesian networks
- conditional independence of all effect variables

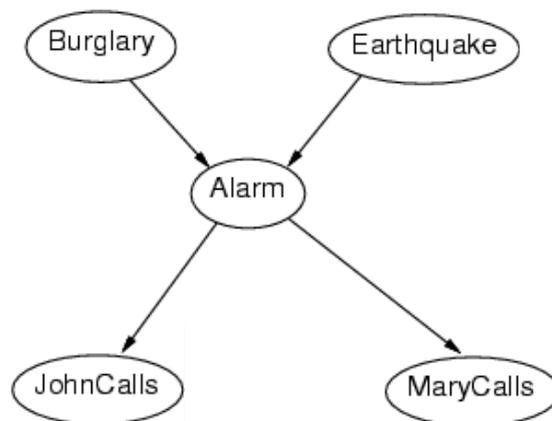
24

## More complex example

- I'm at work, neighbor John calls to say my alarm is ringing, but neighbor Mary calls but not as often as John. Sometimes it's set off by minor earthquakes but also on burglary. Is there a burglar?
- Variables: *Burglary*, *Earthquake*, *Alarm*, *JohnCalls*, *MaryCalls*
- Network topology reflects "causal" knowledge:
  - ♦ A burglar can set the alarm off
  - ♦ An earthquake can set the alarm off
  - ♦ The alarm can cause Mary to call
  - ♦ The alarm can cause John to call

25

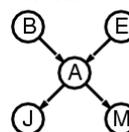
## Example contd.



26

## Compactness

- A CPT for Boolean  $X_i$  with  $k$  Boolean parents has  $2^k$  rows for the combinations of parent values
- Each row requires one number  $p$  for  $X_i = \text{true}$  (the number for  $X_i = \text{false}$  is just  $1-p$ )
- If each of  $n$  Boolean variables has no more than  $k$  parents, the complete network requires  $O(n \cdot 2^k)$  numbers i.e., grows linearly with  $n$ , vs.  $O(2^n)$  for the full joint distribution
- For burglary net?  $1 + 1 + 4 + 2 + 2 = 10$  numbers (vs.  $2^5 - 1 = 31$ )



$k$  parents with  $n$  values each and  $m$  values for the child node of the parents?  
 Number of independent values =  $n^k \cdot (m-1)$

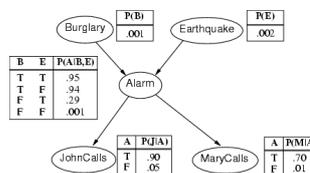
27

## Semantics

The full joint distribution can be rewritten using the *chain rule*:

$$P(X_1, \dots, X_n) = \prod_{i=1}^n P(X_i | X_1, \dots, X_{i-1})$$

$$P(X_1, \dots, X_n) = \prod_{i=1}^n P(X_i | \text{parent}(X_i))$$



Assumption: Independence and Conditional independence assertions are correctly modeled

28

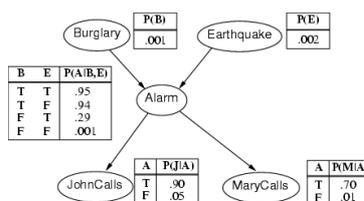
## Semantics

The full joint distribution is defined as the product of the local conditional distributions:

$$P(X_1, \dots, X_n) = \prod_{i=1}^n P(X_i | \text{parent}(X_i))$$

e.g.,  $P(j \wedge m \wedge a \wedge \neg b \wedge \neg e)$

$$\begin{aligned} &= P(j | a) P(m | a) P(a | \neg b, \neg e) P(\neg b) P(\neg e) \\ &= 0.90 \times 0.7 \times 0.001 \times 0.999 \times 0.998 \\ &\approx 0.00063 \end{aligned}$$



29

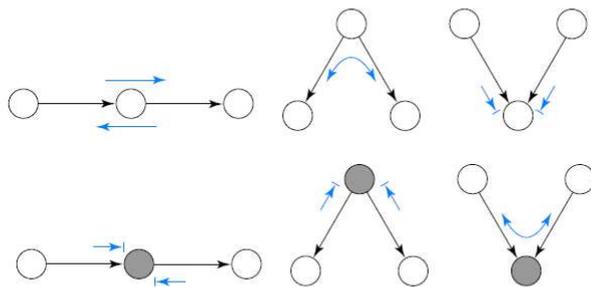
## Encoding conditional independence via d-separation

- We can determine if conditional independence holds by a graph separation criterion called **d-separation** (*direction dependent separation*)
- X and Y are **d-separated** if there is no active path between them.
- The formal definition of **active** is somewhat involved. The Bayes Ball Algorithm gives a nice graphical definition.

30

## The six rules of Bayes Ball

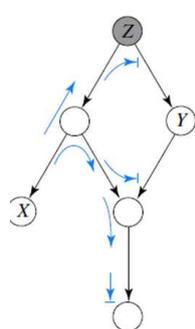
An undirected path is active if a Bayes ball travelling along it never encounters the “stop” symbol:  $\rightarrow \perp$



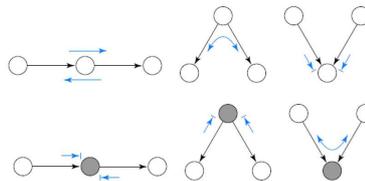
If there are no active paths from  $X$  to  $Y$  when  $\{Z_1, \dots, Z_k\}$  are shaded, then  $X \perp\!\!\!\perp Y \mid \{Z_1, \dots, Z_k\}$ .

31

## A double-header: two games of Bayes Ball

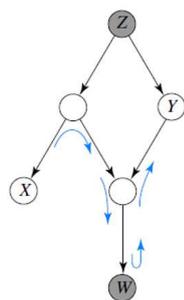


no active paths  
 $X \perp\!\!\!\perp Y \mid Z$

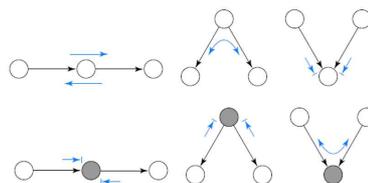


32

## A double-header: two games of Bayes Ball



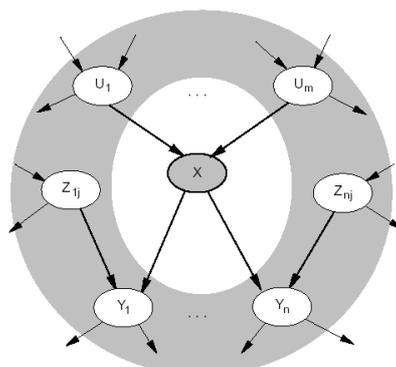
one active path  
 $X \perp\!\!\!\perp Y \mid \{W, Z\}$



33

## Markov Blanket

- Markov blanket: Parents + children + children's parents
- Node is conditionally independent of all other nodes in network, given its Markov Blanket -> simplifies computation -> gather information on the nodes of the Markov Blanket?



34

## Constructing Bayesian networks

- 1. Choose an ordering of variables  $X_1, \dots, X_n$ .  
Cause should precede effects.
- 2. For  $i = 1$  to  $n$ 
  - ♦ add  $X_i$  to the network
  - ♦ select parents from  $X_1, \dots, X_{i-1}$  such that
 
$$P(X_i | \text{Parents}(X_i)) = P(X_i | X_1, \dots, X_{i-1})$$

This choice of parents guarantees:

$$\begin{aligned}
 P(X_1, \dots, X_n) &= \prod_{i=1}^n P(X_i | X_1, \dots, X_{i-1}) \\
 \text{(chain rule)} & \\
 &= \prod_{i=1}^n P(X_i | \text{Parents}(X_i)) \\
 \text{(by construction)} &
 \end{aligned}$$

35

## Example

- Suppose we choose the ordering  $M, J, A, B, E$

$P(J | M) = P(J)$ ? **No**

MaryCalls

JohnCalls

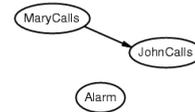
36

## Example

- Suppose we choose the ordering  $M, J, A, B, E$

$P(A | J, M) = P(A)$ ? **No**

$P(A | J, M) = P(A | J)$ ? **No**



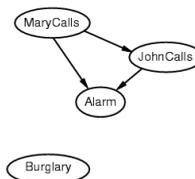
37

## Example

- Suppose we choose the ordering  $M, J, A, B, E$

$P(B | A, J, M) = P(B)$ ? **No**

$P(B | A, J, M) = P(B | A)$ ? **Yes**



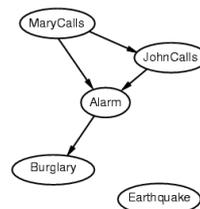
38

## Example

- Suppose we choose the ordering  $M, J, A, B, E$

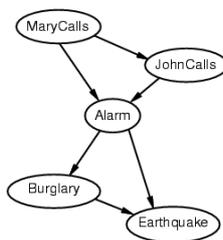
$P(E | B, A, J, M) = P(E | A)$ ? **No**

$P(E | B, A, J, M) = P(E | A, B)$ ? **Yes**



39

## Example contd.



- Deciding conditional independence is hard in noncausal directions
- (Causal models and conditional independence seem hardwired for humans!)
- Network is less compact:  $1 + 2 + 4 + 2 + 4 = 13$  numbers needed instead of 10.

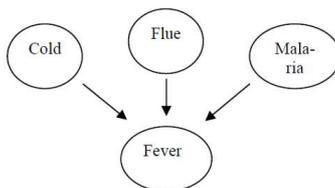
40

## Efficient implementation of CPTs

- The number of independent entries grow exponentially with the number of parents.
- Two ways to overcome this
  - ♦ Restrict the number of parents if possible
  - ♦ Instead of free distributions, often canonical (parameterized) distributions are suggested. One popular example of such a pattern is the noisy OR for discrete cases.

41

## Example



The noisy OR is a generalization of the logical OR. Three assumptions:

1. All possible causes  $U_i$  for an event  $X$  are listed (you can add a *leak* node)
2. Negated causes  $\neg U_i$  do not have any influence on  $X$
3. Independent failure probability  $q_i$  for each cause alone.

$$q_{\text{cold}} = P(\neg \text{fever} \mid \text{cold}, \neg \text{flu}, \neg \text{malaria}) = 0.6 ,$$

$$q_{\text{flu}} = P(\neg \text{fever} \mid \neg \text{cold}, \text{flu}, \neg \text{malaria}) = 0.2 ,$$

$$q_{\text{malaria}} = P(\neg \text{fever} \mid \neg \text{cold}, \neg \text{flu}, \text{malaria}) = 0.1 .$$

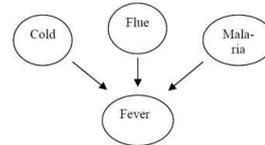
42

# Example

$$q_{\text{cold}} = P(\neg \text{fever} \mid \text{cold}, \neg \text{flu}, \neg \text{malaria}) = 0.6,$$

$$q_{\text{flu}} = P(\neg \text{fever} \mid \neg \text{cold}, \text{flu}, \neg \text{malaria}) = 0.2,$$

$$q_{\text{malaria}} = P(\neg \text{fever} \mid \neg \text{cold}, \neg \text{flu}, \text{malaria}) = 0.1.$$



$$P(\neg x \mid o_1, o_2, \dots, o_r, \neg o_{r+1}, \dots, \neg o_n) = \prod_{i=1}^r q_i$$

Cold	Flu	Malaria	$P(\text{Fever})$	$P(\neg \text{Fever})$
F	F	F	---	---
F	F	T	---	<b>0.1</b>
F	T	F	---	<b>0.2</b>
F	T	T	---	---
T	F	F	---	<b>0.6</b>
T	F	T	---	---
T	T	F	---	---
T	T	T	---	---

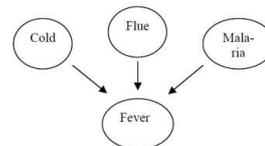
43

# Example

$$q_{\text{cold}} = P(\neg \text{fever} \mid \text{cold}, \neg \text{flu}, \neg \text{malaria}) = 0.6,$$

$$q_{\text{flu}} = P(\neg \text{fever} \mid \neg \text{cold}, \text{flu}, \neg \text{malaria}) = 0.2,$$

$$q_{\text{malaria}} = P(\neg \text{fever} \mid \neg \text{cold}, \neg \text{flu}, \text{malaria}) = 0.1.$$



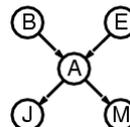
$$P(x \mid o_1, o_2, \dots, o_r, \neg o_{r+1}, \dots, \neg o_n) = 1 - \prod_{i=1}^r q_i$$

Cold	Flu	Malaria	$P(\text{Fever})$	$P(\neg \text{Fever})$
F	F	F	0.0	1.0
F	F	T	0.9	<b>0.1</b>
F	T	F	0.8	<b>0.2</b>
F	T	T	0.98	$0.02 = 0.2 \times 0.1$
T	F	F	0.4	<b>0.6</b>
T	F	T	0.94	$0.06 = 0.6 \times 0.1$
T	T	F	0.88	$0.12 = 0.6 \times 0.2$
T	T	T	0.988	$0.012 = 0.6 \times 0.2 \times 0.1$

44

## Last Time

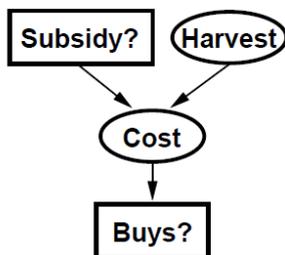
- Structure and semantic of BN
  - ♦ Modelling of independence and conditional independence
  - ♦ Causal and non-causal networks
  - ♦ d-separation, Markov blanket
  - ♦ Efficient CPTs, e.g., noisy OR, trees, Min, Max, ...



45

## Hybrid (discrete+continuous) networks

Discrete (*Subsidy?* and *Buys?*); continuous (*Harvest* and *Cost*)



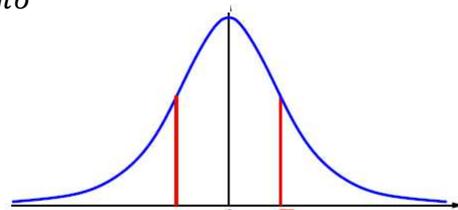
Option 1: discretization—possibly large errors, large CPTs  
 Option 2: finitely parameterized canonical families

- 1) Continuous variable, discrete+continuous parents (e.g., *Cost*)
- 2) Discrete variable, continuous parents (e.g., *Buys?*)

46

## Continuous variables: Gaussian density

$$P(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/2\sigma^2}$$



$$\text{Mean } \mu = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\text{Variance } \sigma^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \mu)^2$$

$$\text{Standard deviation } \sigma = \sqrt{\sigma^2}$$

47

## Continuous child variables

- Need one *conditional density* function for child variable given
  - continuous parents
  - for each discrete value of parents



$$P(c|h, \text{subsidy}) = N(a_t h + b_t, \sigma_t^2)(c) = \frac{1}{\sigma_t \sqrt{2\pi}} e^{-\frac{1}{2} \left( \frac{c - (a_t h + b_t)}{\sigma_t} \right)^2}$$

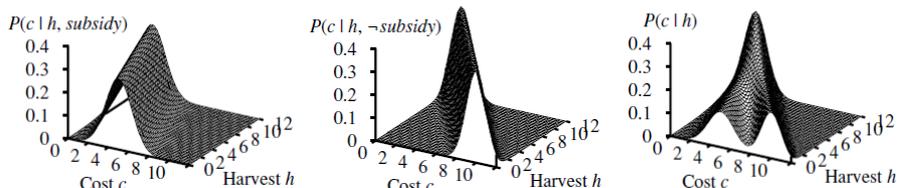
$$P(c|h, \neg \text{subsidy}) = N(a_f h + b_f, \sigma_f^2)(c) = \frac{1}{\sigma_f \sqrt{2\pi}} e^{-\frac{1}{2} \left( \frac{c - (a_f h + b_f)}{\sigma_f} \right)^2}$$

- Mean **Cost** varies *linearly* with **Harvest**, variance fixed
- Linear variation is unreasonable over the full range but works if the likely range of **Harvest** is narrow

48

## Continuous child variables

- Determine a Gaussian for *subsidy* and  $\neg$ *subsidy*
- What happens if subsidy is not given  $P(c|h)$ ?



All-continuous network with LG distributions  
 $\Rightarrow$  full joint distribution is a multivariate Gaussian

Discrete+continuous LG network is a conditional Gaussian network i.e., a multivariate Gaussian over all continuous variables for each combination of discrete variable values

49

## Discrete variabel cont. parents

- Probability of Buys given Cost should be a soft threshold

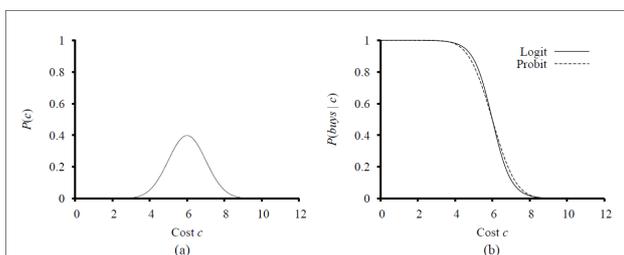


Figure 14.7 FILES: . (a) A normal (Gaussian) distribution for the cost threshold, centered on  $\mu = 6.0$  with standard deviation  $\sigma = 1.0$ . (b) Logit and probit distributions for the probability of *buys* given *cost*, for the parameters  $\mu = 6.0$  and  $\sigma = 1.0$ .

Use integral  $\Phi(x) = \int^x N(0, 1)(x)dx$   
 Leads to  $P(\text{buys} \mid \text{Cost} = c) = \Phi((-c + \mu)/\sigma)$  Probit  
 Alternativ  $P(\text{buys} \mid \text{Cost} = c) = \frac{1}{1 + \exp(-2\frac{-c+\mu}{\sigma})}$  Logit

50

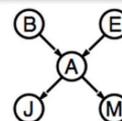
## Inference tasks

- **Simple queries:**  $P(X_1, \dots, X_n | e_2, e_4, e_5)$
- **Optimal decisions:** decision networks include utility information; inference must handle utility nodes.
- **Value of information:** which evidence to seek next?
- **Sensitivity analysis:** which probability values are more critical?
- **Explanation:** why do I need a new engine?

51

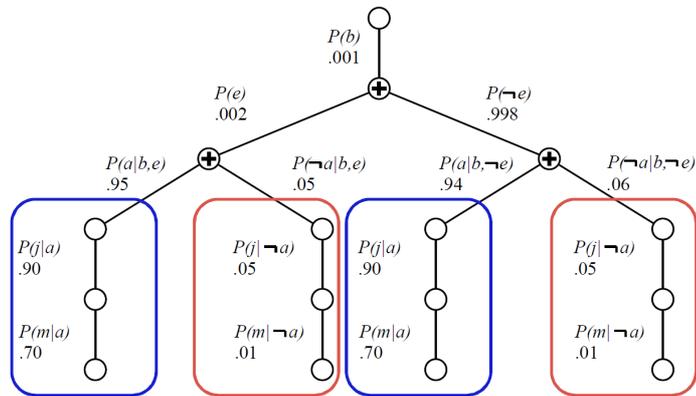
## Inference by enumeration

$$\begin{aligned}
 P(b|j,m) &= \alpha P(b,j,m) \\
 &= \alpha \sum_a \sum_e P(b \wedge j \wedge m \wedge a \wedge e) \text{ [marginalization]} \\
 &= \alpha \sum_a \sum_e P(b)P(e)P(a|b,e)P(j|a)P(m|a) \text{ [BN]} \\
 &= \alpha P(b) \sum_e P(e) \sum_a P(a|b,e)P(j|a)P(m|a) \text{ [re-ordering]}
 \end{aligned}$$



52

## Evaluation Tree



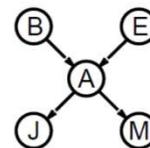
Enumeration is inefficient: repeated computation  
 e.g., computes  $P(j|a)P(m|a)$  for each value of  $e$

53

## Irrelevant variables

Consider the query  $P(\text{JohnCalls} | \text{Burglary} = \text{true})$

$$P(J|b) = \alpha P(b) \sum_e P(e) \sum_a P(a|b,e) P(J|a) \sum_m P(m|a)$$



**What about M?**

**We sum over all possible values of m**

**For each row it means that the value is 1**

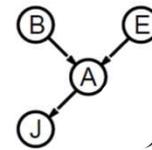
62

## Irrelevant variables

Consider the query  $P(\text{JohnCalls} | \text{Burglary} = \text{true})$

$$P(J|b) = \alpha P(b) \sum_e P(e) \sum_a P(a|b, e) P(J|a) \sum_m P(m|a)$$

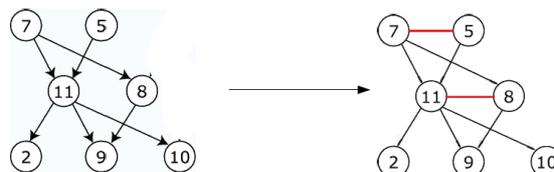
For each row it means that the value is 1



63

## Moral Graph: Markov Blanket

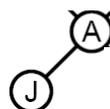
- The moral graph is an undirected graph that is obtained as follows:
  - ♦ connect all parents of all nodes
  - ♦ make all directed links undirected
- Note:
  - ♦ the moral graph connects each node to all nodes of its **Markov blanket**
    - it is already connected to parents and children
    - now it is also connected to the parents of its children



64

## Irrelevant variables continued:

- m-separation:
  - ♦ A is m-separated from B by C iff it is separated by C in the moral graph
- Example:
  - ♦ J is m-separated from E by A



**Theorem 2:** Y is irrelevant if it is m-separated from X by E

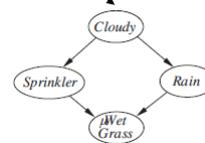
- Example:
 

*For  $P(\text{JohnCalls} | \text{Alarm}=\text{true})$ ,  
Burglary, Earthquake and MarryCalls are irrelevant.*

65

## Approximate Inference In Bayesian Network

- Singly connected networks (or polytrees):
  - ♦ any two nodes are connected by at most one (undirected) path
  - ♦ time and space cost of variable elimination linear in the size of the network (number of CPT entries; number of parents  $O(d^n)$ ).
- Multiply connected networks: NP-hard
- We need **approximate inference techniques!!!!!!**
- Monte Carlo algorithm
  - ♦ Widely used to estimate quantities that are difficult to calculate exactly
  - ♦ Randomized sampling algorithm
  - ♦ Accuracy depends on the number of samples
  - ♦ Two families
    - Direct sampling
    - Markov chain sampling



66

## Inference by stochastic simulation

Basic idea:

- 1) Draw  $N$  samples from a sampling distribution  $S$
- 2) Compute an approximate posterior probability  $\hat{P}$
- 3) Show this converges to the true probability  $P$

0.5



Outline:

- Sampling from an empty network
- Rejection sampling: reject samples disagreeing with evidence
- Likelihood weighting: use evidence to weight samples
- Markov chain Monte Carlo (MCMC): sample from a stochastic process whose stationary distribution is the true posterior

67

## Sampling from empty network

- Generating samples from a network that has no evidence associated with it (*empty network*)
- Basic idea
  - ♦ sample a value for each variable in topological order
  - ♦ using the specified conditional probabilities

```

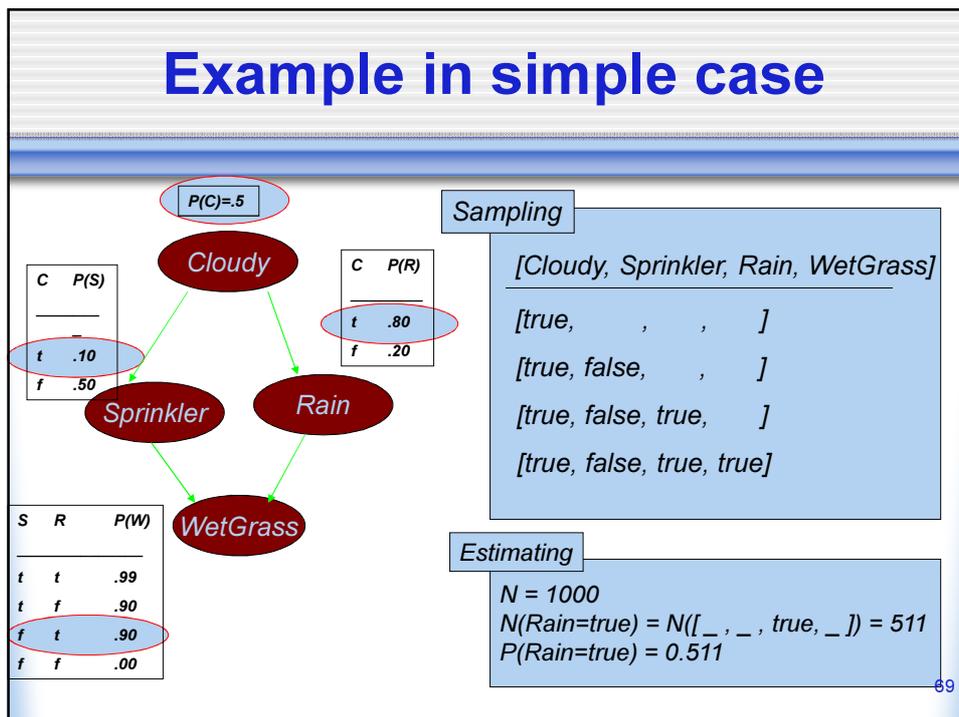
function PRIOR-SAMPLE( $bn$ ) returns an event sampled from  $bn$ 
inputs:  $bn$ , a belief network specifying joint distribution  $\mathbf{P}(X_1, \dots, X_n)$ 

 $x \leftarrow$  an event with  $n$  elements
for  $i = 1$  to  $n$  do
   $x_i \leftarrow$  a random sample from  $\mathbf{P}(X_i \mid \text{parents}(X_i))$ 
  given the values of  $\text{Parents}(X_i)$  in  $x$ 
return  $x$ 

```

68

## Example in simple case



## Properties

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

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

- Used to compute conditional probabilities
- Procedure
  - ◆ Generating sample from prior distribution specified by the Bayesian Network
  - ◆ Rejecting all that do not match the evidence
  - ◆ Estimating probability

71

## Rejection Sampling

$\hat{P}(X|e)$  estimated from samples agreeing with  $e$

```

function REJECTION-SAMPLING( $X, e, bn, N$ ) returns an estimate of  $P(X|e)$ 
  local variables:  $N$ , a vector of counts over  $X$ , initially zero
  for  $j = 1$  to  $N$  do
     $x \leftarrow$  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( $N[X]$ )
  
```

72

## Rejection Sampling Example

- Let us assume we want to estimate  $P(\text{Rain}|\text{Sprinkler} = \text{true})$  with 100 samples
- 100 samples
  - ♦ 73 samples => Sprinkler = false
  - ♦ 27 samples => Sprinkler = true
    - 8 samples => Rain = true
    - 19 samples => Rain = false
- $P(\text{Rain}|\text{Sprinkler} = \text{true}) = \text{NORMALIZE}(\{8, 19\}) = \{0.296, 0.704\}$
- The true answer ist  $\langle 0.3, 0.7 \rangle$
- Problem
  - ♦ It rejects too many samples

73

## Analysis of rejection sampling

$$\begin{aligned}
 \hat{P}(X|\mathbf{e}) &= \alpha N_{PS}(X, \mathbf{e}) && \text{(algorithm defn.)} \\
 &= N_{PS}(X, \mathbf{e}) / N_{PS}(\mathbf{e}) && \text{(normalized by } N_{PS}(\mathbf{e})\text{)} \\
 &\approx P(X, \mathbf{e}) / P(\mathbf{e}) && \text{(property of PRIORSAMPLE)} \\
 &= P(X|\mathbf{e}) && \text{(defn. of conditional probability)}
 \end{aligned}$$

Hence rejection sampling returns consistent posterior estimates

Problem: hopelessly expensive if  $P(\mathbf{e})$  is small

$P(\mathbf{e})$  drops off exponentially with number of evidence variables!

75

## Likelihood Weighting

- Goal
  - ◆ Avoiding inefficiency of rejection sampling
- Idea
  - ◆ Generating only events consistent with evidence
  - ◆ Each event is weighted by likelihood that the event accords to the evidence

76

## Likelihood weighting

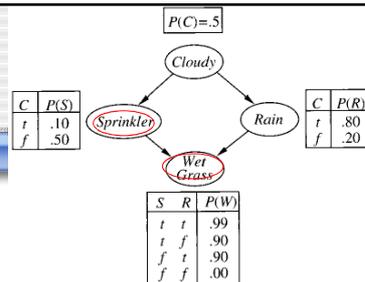
Idea: fix evidence variables, sample only nonevidence variables, and weight each sample by the likelihood it accords the evidence

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

```
function WEIGHTED-SAMPLE( $bn, e$ ) returns an event and a weight
   $x \leftarrow$  an event with  $n$  elements;  $w \leftarrow 1$ 
  for  $i = 1$  to  $n$  do
    if  $X_i$  has a value  $x_i$  in  $e$ 
      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  $x, w$ 
```

77

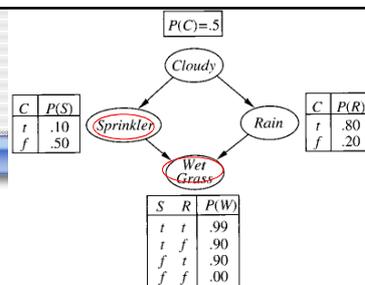
## Likelihood Weighting Example



- $P(\text{Rain} | \text{Sprinkler}=\text{true}, \text{WetGrass}=\text{true})?$
- Sampling, start with weight=1
  - ♦ Sample from  $P(\text{Cloudy}) = \{0.5, 0.5\} \Rightarrow \text{true}$
  - ♦ Sprinkler is an evidence variable with value **true**  
 $w \leftarrow w * P(\text{Sprinkler}=\text{true} | \text{Cloudy}=\text{true}) = 0.1$
  - ♦ Sample from  $P(\text{Rain} | \text{Cloudy}=\text{true}) = \{0.8, 0.2\} \Rightarrow \text{true}$
  - ♦ WetGrass is an evidence variable with value **true**  
 $w \leftarrow w * P(\text{WetGrass}=\text{true} | \text{Sprinkler}=\text{true}, \text{Rain}=\text{true}) = 0.099$
  - ♦ **[true, true, true, true]** with weight 0.099

78

## Likelihood Weighting Example



- $P(\text{Rain} | \text{Sprinkler}=\text{true}, \text{WetGrass}=\text{true})?$
- Sampling, start with weight=1
  - ♦ Sample from  $P(\text{Cloudy}) = \{0.5, 0.5\} \Rightarrow \text{false}$
  - ♦ Sprinkler is an evidence variable with value **true**  
 $w \leftarrow w * P(\text{Sprinkler}=\text{true} | \text{Cloudy}=\text{false}) = 0.5$
  - ♦ Sample from  $P(\text{Rain} | \text{Cloudy}=\text{false}) = \{0.2, 0.8\} \Rightarrow \text{false}$
  - ♦ WetGrass is an evidence variable with value **true**  
 $w \leftarrow w * P(\text{WetGrass}=\text{true} | \text{Sprinkler}=\text{true}, \text{Rain}=\text{false}) = 0.45$
  - ♦ **[true, true, true, true]** with weight 0.45
- Estimating
  - ♦ Accumulating weights to either Rain=true or Rain=false
  - ♦ Normalize

79

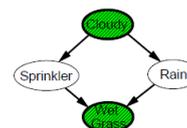
## Likelihood analysis

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

81

## Markov Chain Monte Carlo

- Let's think of the network as being in a particular current state specifying a value for every variable
- MCMC generates each event by making a random change to the preceding event
- The next state is generated by randomly sampling a value for one of the non evidence variables  $X_i$ , **conditioned on the current values of the variables in the MarkovBlanket of  $X_i$**
- Likelihood Weighting only takes into account the evidences of the parents.

82

## Gibbs sampling

- Gibbs sampling is a MCMC method
  - ♦ State of the network => current assignment
  - ♦ Generate next state by sampling one non-evidence variable given Markov blanket
  - ♦ Sample each variable in turn ( can choose it random)

```

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

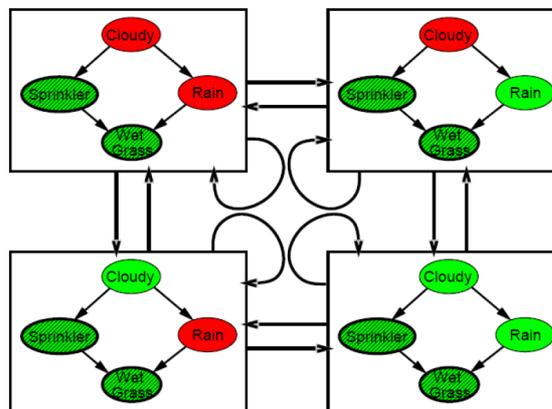
  initialize  $x$  with random values for the variables in  $Z$ 
  for  $j = 1$  to  $N$  do
    for each  $Z_i$  in  $Z$  do
      set the value of  $Z_i$  in  $x$  by sampling from  $P(Z_i|mb(Z_i))$ 
       $N[x] \leftarrow N[x] + 1$  where  $x$  is the value of  $X$  in  $x$ 
  return NORMALIZE( $N$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

**Figure 14.16** The Gibbs sampling algorithm for approximate inference in Bayesian networks; this version cycles through the variables, but choosing variables at random also works.

83

## Example

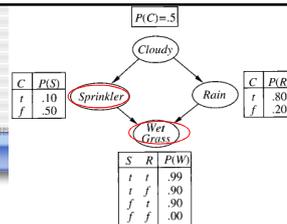
With *Sprinkler = true*, *WetGrass = true*, there are four states:



Wander about for a while, average what you see

84

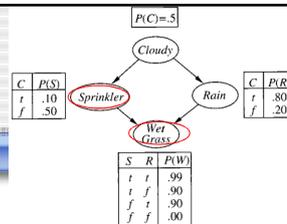
## Gibbs Example



- Query  $P(\text{Rain} \mid \text{Sprinkler} = \text{true}, \text{WetGrass} = \text{true})$
- Initial state is  $[\text{true}, \text{true}, \text{false}, \text{true}]$   $[\text{Cloudy}, \text{Sprinkler}, \text{Rain}, \text{WetGrass}]$
- The following steps are executed repeatedly:
  - $\text{Cloudy}$  is sampled, given the current values of its Markov Blanket variables  
So, we sample from  $P(\text{Cloudy} \mid \text{Sprinkler} = \text{true}, \text{Rain} = \text{false})$   
The result is  $\text{Cloudy} = \text{false}$  (???????)
  - Now current state is  $[\text{false}, \text{true}, \text{false}, \text{true}]$  and counts are updated
  - $\text{Rain}$  is sampled, given the current values of its Markov Blanket variables  
Sample from  $P(\text{Rain} \mid \text{Cloudy} = \text{false}, \text{Sprinkler} = \text{true}, \text{WetGrass} = \text{true})$   
First create the distribution we want to sample from.  
 $\rightarrow \text{Rain} = \text{true}$ .
  - Current state is  $[\text{false}, \text{true}, \text{true}, \text{true}]$
- After all the iterations, let's say the process visited 20 states where rain is true and 60 states where rain is false then the answer of the query is  $\text{NORMALIZE}(\{20, 60\}) = \{0.25, 0.75\}$

85

## Sample distribution



Want to sample  $\text{Cloudy}$ .

The current state is  $[\text{Cloudy?}, \text{true}, \text{false}, \text{true}]$

What is the Markov blanket, the sampling distribution?

*evidence*      *sampled*

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

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

$$\propto \langle 0.5, 0.5 \rangle * \langle 0.1, 0.5 \rangle * \langle 0.2, 0.8 \rangle =$$

$$\propto \langle 0.5, 0.5 \rangle * \langle 0.1 * 0.2, 0.5 * 0.8 \rangle =$$

$$\propto \langle 0.5, 0.5 \rangle * \langle 0.02, 0.4 \rangle =$$

$$\propto \langle 0.01, 0.2 \rangle \approx \langle 0.05, 0.95 \rangle$$

$[\text{false}, \text{true}, \text{false}, \text{true}]$  with probability 0,95

$[\text{true}, \text{true}, \text{false}, \text{true}]$  with probability 0,05

86

## Summary

- Bayesian networks provide a natural representation for (causally induced) conditional independence
- Topology + CPTs = compact representation of joint distribution
- Generally easy for domain experts to construct (if not too big)
- Exact inference by variable elimination
  - ♦ polytime on polytrees, NP-hard on general graphs
  - ♦ space can be exponential as well
- Approximate inference based on sampling and counting help to overcome complexity of exact inference

87