

# UNCERTAIN KNOWLEDGE AND REASONING

8

## 8 UNCERTAIN KNOWLEDGE AND REASONING

---

8.1 Uncertainty

8.2 Probability

- Syntax and semantics ● Inference
- Independence ● Bayes' rule

8.3 Bayesian networks

8.4 Probabilistic reasoning<sup>+</sup>

8.5 Dynamic Bayesian networks<sup>+</sup>

8.6 Causal Inference<sup>\*</sup>

8.7 Probabilistic programming<sup>\*</sup>

8.8 Probabilistic logic<sup>\*</sup>

# Uncertainty

---

Let action  $A_t$  = leave for airport  $t$  minutes before flight  
Will  $A_t$  get me there on time?

Problems

- 1) partial observability (road state, other drivers' plans, etc.)
- 2) noisy sensors (traffic radio)
- 3) uncertainty in action outcomes (flat tire, etc.), etc.

Hence a purely logical approach either

- 1) risks falsehood: " $A_{25}$  will get me there on time"
- or 2) leads to conclusions that are too weak for decision making  
" $A_{25}$  will get me there on time if there's no accident on the bridge and it doesn't rain and my tires remain intact etc."

( $A_{1440}$  might reasonably be said to get me there on time  
but I'd have to stay overnight in the airport ...)

# Uncertainty knowledge representation and reasoning

## Nonmonotonic logic

Assume  $A_{25}$  works unless contradicted by evidence

Issues: How to handle quantitation? Reasonable assumptions?

## Rules with fudge factors:

$A_{25} \mapsto_{0.3} \textit{AtAirportOnTime}$

$\textit{Sprinkler} \mapsto_{0.99} \textit{WetGrass}$

$\textit{WetGrass} \mapsto_{0.7} \textit{Rain}$

Issues: problems with combination, e.g., *Sprinkler* causes *Rain*?

Fuzzy logic handles **degree of truth** NOT uncertainty e.g.,

*WetGrass* is true to degree 0.2

## Probability

Given the available evidence,

$A_{25}$  will get me there on time with probability 0.04

Qualitative vs. quantitative  $\Rightarrow$  Logic vs. probability  $\Leftarrow$  Prob. logics

# Probability

---

Probabilistic assertions **summarize** effects of

**laziness**: failure to enumerate exceptions, qualifications, etc.

**ignorance**: lack of relevant facts, initial conditions, etc.

Subjective (posterior, conditional, Bayesian) probability

Probabilities relate propositions to one's own state of knowledge

e.g.,  $P(A_{25} | \text{no reported accidents}) = 0.06$

These are **not** claims of a “probabilistic tendency” in the current situation

(but might be learned from past experience of similar situations)

Probabilities of propositions change with new evidence

e.g.,  $P(A_{25} | \text{no reported accidents, 5 a.m.}) = 0.15$

(Analogous to logical entailment  $KB \models \alpha$ , not truth but nonmonotonic in nature)

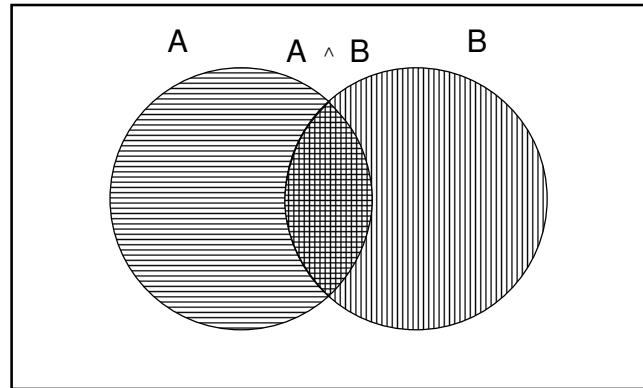
# Why use probability?

---

The definitions imply that certain logically related events must have related probabilities

E.g.,  $P(a \vee b) = P(a) + P(b) - P(a \wedge b)$

True



de Finetti (1931): an agent who bets according to probabilities that violate these axioms can be forced to bet so as to lose money regardless of the outcome

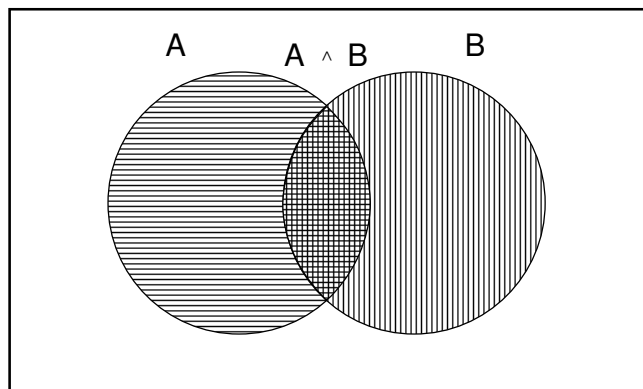
# Axioms of probability

---

For any propositions  $A, B$

1.  $0 \leq P(A) \leq 1$
2.  $P(\text{True}) = 1$  and  $P(\text{False}) = 0$
3.  $P(A \vee B) = P(A) + P(B) - P(A \wedge B)$

True



A probability is a measure over a set of events that satisfies three axioms  $\Rightarrow$  probability theory is analogous to logical theory (axioms)

e.g.,  $P(\neg a) = 1 - P(a)$  is derived from the axioms

$P(a \vee b) = P(a) + P(b) - P(a \wedge b)$  (inclusion-exclusion principle)

# Syntax and semantics

---

Traditional probability theory has informal language  
needs to be formalized for agents

Begin with a set  $\Omega$  — sample space

e.g., 6 possible rolls of a die

$\omega \in \Omega$  is a sample point (outcome/possible world/atomic event/data)

A probability space or probability model is a sample space  
with an assignment  $P(\omega)$  for every  $\omega \in \Omega$  s.t.

$$(1) 0 \leq P(\omega) \leq 1$$

$$(2) \sum_{\omega} P(\omega) = 1$$

e.g.,  $P(1) = P(2) = P(3) = P(4) = P(5) = P(6) = 1/6$

An event  $A$  is any subset of  $\Omega$

$$P(A) = \sum_{\{\omega \in A\}} P(\omega)$$

e.g.,  $P(\text{die roll} < 4) = P(1) + P(2) + P(3) = 1/6 + 1/6 + 1/6 = 1/2$



# Random variables

---

A **random variable** is a function from sample points to some **range**

- **Booleans** (propositions)

e.g., *Cavity* (do I have a cavity?)

*Cavity = true* is a proposition, also written *Cavity*

- **Discrete** (**finite** or **infinite**)

e.g., *Weather* is one of  $\langle \textit{sunny}, \textit{rain}, \textit{cloudy}, \textit{snow} \rangle$

*Weather = rain* is a proposition

Values must be exhaustive and mutually exclusive

- **Continuous** or **real** (**bounded** or **unbounded**)

e.g., *Temp = 21.6*; also allow, e.g., *Temp < 22.0*

Arbitrary Boolean combinations of basic propositions

# Probability distribution

---

$P$  induces a (prob.) distribution for any r.v. (random variable)  $X$

$$P(X = x_i) = \sum_{\{\omega: X(\omega) = x_i\}} P(\omega)$$

gives values for all possible assignments

E.g.,  $P(\text{Odd} = \text{true}) = P(1) + P(3) + P(5) = 1/6 + 1/6 + 1/6 = 1/2$

The probability of a proposition  $\text{Odd} = \text{true}$  as the sum of the probabilities of worlds in which it holds

# Propositions

---

Think of a proposition as the event (set of sample points)  
where the proposition is true

Given Boolean r.v.s  $A$  and  $B$

event  $a$  = set of sample points where  $A(\omega) = \text{true}$

event  $\neg a$  = set of sample points where  $A(\omega) = \text{false}$

event  $a \wedge b$  = points where  $A(\omega) = \text{true}$  and  $B(\omega) = \text{true}$

The sample points are defined by the values of a set of r.v.s  
i.e., the sample space is Cartesian product of the ranges of the r.v.s

# Propositions

---

For **Boolean** r.v.s

sample point (possible world) = propositional logic model

e.g.,  $A = true$ ,  $B = false$ , or  $a \wedge \neg b$

A possible world is defined to be an assignment of values to all of the r.v.s under consideration

– possible worlds are mutually exclusive and exhaustive, why??

Proposition = disjunction of atomic events (clausal form)

e.g.,  $(a \vee b) \equiv (\neg a \wedge b) \vee (a \wedge \neg b) \vee (a \wedge b)$

$\Rightarrow P(a \vee b) = P(\neg a \wedge b) + P(a \wedge \neg b) + P(a \wedge b)$

For any proposition  $\phi$ , the possible world (model)  $\omega$  where it is true

$\omega \models \phi$

**Hint:** (propositional) logic + probability  $\Rightarrow$  probabilistic logic

# Prior probability

---

Prior (unconditional probabilities) of propositions

$$\text{e.g., } P(\textit{Cavity} = \textit{true}) = 0.1$$

$$P(\textit{Weather} = \textit{sunny}) = 0.72$$

correspond to belief prior to the arrival of any (new) evidence

Probability distribution gives values for all possible assignments

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

(normalized, i.e., sums to 1)

# Joint probability distribution

---

Joint probability distribution for a set of r.v.s gives the probability of every atomic event on those r.v.s (i.e., every sample point)

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

<i>Weather =</i>	<i>sunny</i>	<i>rain</i>	<i>cloudy</i>	<i>snow</i>
<i>Cavity = true</i>	0.144	0.02	0.016	0.02
<i>Cavity = false</i>	0.576	0.08	0.064	0.08

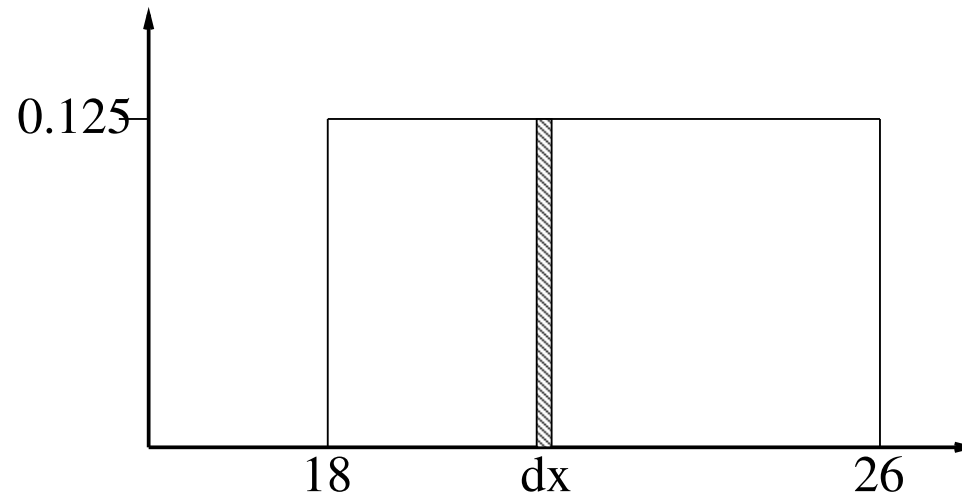
Every question about a domain can be answered by the joint distribution because every event is a sum of sample points

# Probability for continuous variables<sup>#</sup>

---

Express distribution as a parameterized function of value

$$P(X = x) = U[18, 26](x) = \text{uniform density between 18 and 26}$$



Here  $P$  is a **density**; integrates to 1

$P(X = 20.5) = 0.125$  really means

$$\lim_{dx \rightarrow 0} P(20.5 \leq X \leq 20.5 + dx) / dx = 0.125$$

# Conditional probability

---

Conditional (posterior) probabilities

e.g.,  $P(\text{cavity}|\text{toothache}) = 0.8$

i.e., given evidence that *toothache* is all I know

NOT “if *toothache* then 80% chance of *cavity*”

Full joint (conditional probability) distribution for all of the r.v.s

$\mathbf{P}(\text{Cavity}|\text{Toothache}) = 2\text{-element vector of } 2\text{-element vectors}$

If we know more, e.g., *cavity* is also given, then we have

$P(\text{cavity}|\text{toothache}, \text{cavity}) = 1$

Note: the less specific belief remains valid after more evidence arrives, but is not always useful

New evidence may be irrelevant, allowing simplification, e.g.,

$P(\text{cavity}|\text{toothache}, \text{49ersWin}) = P(\text{cavity}|\text{toothache}) = 0.8$

This kind of inference, sanctioned by domain knowledge, is crucial



# Conditional probability

---

Defn. of conditional probability by unconditional probabilities

$$P(a|b) = \frac{P(a \wedge b)}{P(b)} \text{ if } P(b) \neq 0$$

**Product rule** gives an alternative formulation

$$P(a \wedge b) = P(a|b)P(b) = P(b|a)P(a)$$

A general version holds for whole distributions, e.g.,

$$\mathbf{P}(\textit{Weather}, \textit{Cavity}) = \mathbf{P}(\textit{Weather}|\textit{Cavity})\mathbf{P}(\textit{Cavity})$$

(View as a  $4 \times 2$  set of equations, **not** matrix mult.)

**Chain rule** is derived by successive application of product rule

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

# Inference

Probabilistic **inference** is the computation of posterior probabilities for query propositions given observed evidence

where the full joint distribution can be viewed as the KB  
from which answers to all questions may be derived

Start with the joint 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

$$P(\phi) = \sum_{\omega:\omega\models\phi} P(\omega)$$

# Inference by enumeration

---

Start with the joint 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

$$P(\phi) = \sum_{\omega:\omega\models\phi} P(\omega)$$

E.g.,  $P(\textit{toothache}) = 0.108 + 0.012 + 0.016 + 0.064 = 0.2$

# Inference by enumeration

---

Start with the joint 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

$$P(\phi) = \sum_{\omega:\omega\models\phi} P(\omega)$$

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

# Inference by enumeration

Start with the joint 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

$$\begin{aligned} P(\neg cavity | toothache) &= \frac{P(\neg cavity \wedge toothache)}{P(toothache)} \\ &= \frac{0.016 + 0.064}{0.108 + 0.012 + 0.016 + 0.064} = 0.4 \end{aligned}$$

# Normalization

	<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>

Denominator can be viewed as a normalization constant  $\alpha$

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

Idea: compute distribution on query variable  
by fixing **evidence variables** and summing over **hidden variables**

## Inference by enumeration contd.

---

Let  $\mathbf{X}$  be all the variables. *ASK*

the posterior joint distribution of the query variables  $\mathbf{Y}$   
given specific values  $\mathbf{e}$  for the evidence variables  $\mathbf{E}$

Let the hidden variables be  $\mathbf{H} = \mathbf{X} - \mathbf{Y} - \mathbf{E}$

$\Rightarrow$  the required summation of joint entries is done by *summing out*  
the hidden variables:

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

The terms in the summation are joint entries because  $\mathbf{Y}$ ,  $\mathbf{E}$ , and  $\mathbf{H}$   
together exhaust the set of random variables

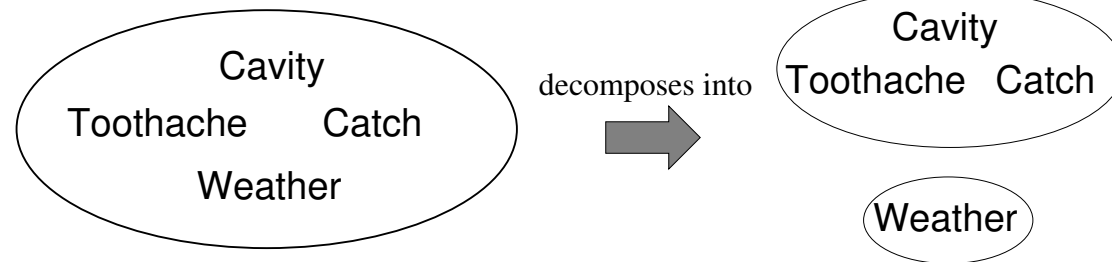
Problems

- 1) Worst-case time complexity  $O(d^n)$  where  $d$  is the largest arity
- 2) Space complexity  $O(d^n)$  to store the joint distribution
- 3) How to find the numbers for  $O(d^n)$  entries?

# Independence

$A$  and  $B$  are independent iff

$$\mathbf{P}(A|B) = \mathbf{P}(A) \quad \text{or} \quad \mathbf{P}(B|A) = \mathbf{P}(B) \quad \text{or} \quad \mathbf{P}(A, B) = \mathbf{P}(A)\mathbf{P}(B)$$



$$\begin{aligned} &\mathbf{P}(\textit{Toothache}, \textit{Catch}, \textit{Cavity}, \textit{Weather}) \\ &= \mathbf{P}(\textit{Toothache}, \textit{Catch}, \textit{Cavity})\mathbf{P}(\textit{Weather}) \end{aligned}$$

32 entries reduced to 12; for  $n$  independent biased coins,  $2^n \rightarrow n$

Absolute independence is powerful but rare

Dentistry is a large field with hundreds of variables, none of which are independent. What to do?



# Conditional independence

---

$\mathbf{P}(\textit{Toothache}, \textit{Cavity}, \textit{Catch})$  has  $2^3 - 1 = 7$  independent entries

If I have a cavity, the probability that the probe catches in it doesn't depend on whether I have a toothache

$$(1) P(\textit{catch}|\textit{toothache}, \textit{cavity}) = P(\textit{catch}|\textit{cavity})$$

The same independence holds if I haven't got a cavity

$$(2) P(\textit{catch}|\textit{toothache}, \neg\textit{cavity}) = P(\textit{catch}|\neg\textit{cavity})$$

*Catch* is **conditionally independent** of *Toothache* given *Cavity*

$$\mathbf{P}(\textit{Catch}|\textit{Toothache}, \textit{Cavity}) = \mathbf{P}(\textit{Catch}|\textit{Cavity})$$

Equivalent statements

$$\mathbf{P}(\textit{Toothache}|\textit{Catch}, \textit{Cavity}) = \mathbf{P}(\textit{Toothache}|\textit{Cavity})$$

$$\mathbf{P}(\textit{Toothache}, \textit{Catch}|\textit{Cavity}) = \mathbf{P}(\textit{Toothache}|\textit{Cavity})\mathbf{P}(\textit{Catch}|\textit{Cavity})$$

# Conditional independence

---

Write out full joint distribution using chain rule

$$\begin{aligned} & \mathbf{P}(\textit{Toothache}, \textit{Catch}, \textit{Cavity}) \\ &= \mathbf{P}(\textit{Toothache} | \textit{Catch}, \textit{Cavity}) \mathbf{P}(\textit{Catch}, \textit{Cavity}) \\ &= \mathbf{P}(\textit{Toothache} | \textit{Catch}, \textit{Cavity}) \mathbf{P}(\textit{Catch} | \textit{Cavity}) \mathbf{P}(\textit{Cavity}) \\ &= \mathbf{P}(\textit{Toothache} | \textit{Cavity}) \mathbf{P}(\textit{Catch} | \textit{Cavity}) \mathbf{P}(\textit{Cavity}) \end{aligned}$$

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 uncertainty

# Bayes' rule

---

Product rule  $P(a \wedge b) = P(a|b)P(b) = P(b|a)P(a)$

$$\Rightarrow \text{Bayes' rule } P(a|b) = \frac{P(b|a)P(a)}{P(b)}$$

or in distribution form

$$\mathbf{P}(Y|X) = \frac{\mathbf{P}(X|Y)\mathbf{P}(Y)}{\mathbf{P}(X)} = \alpha\mathbf{P}(X|Y)\mathbf{P}(Y)$$

Useful for assessing **diagnostic** probability from **causal** probability

$$P(Cause|Effect) = \frac{P(Effect|Cause)P(Cause)}{P(Effect)}$$

E.g., let  $M$  be meningitis,  $S$  be stiff neck

$$P(m|s) = \frac{P(s|m)P(m)}{P(s)} = \frac{0.8 \times 0.0001}{0.1} = 0.0008$$

# Naive Bayes

---

Bayes' rule and conditional independence

$$\begin{aligned} & \mathbf{P}(Cavity|toothache \wedge catch) \\ &= \alpha \mathbf{P}(toothache \wedge catch|Cavity)\mathbf{P}(Cavity) \\ &= \alpha \mathbf{P}(toothache|Cavity)\mathbf{P}(catch|Cavity)\mathbf{P}(Cavity) \end{aligned}$$

This is an example of a **naive Bayes** model (Bayesian classifier)

$$\mathbf{P}(Cause, Effect_1, \dots, Effect_n) = \mathbf{P}(Cause) \prod_i \mathbf{P}(Effect_i|Cause)$$



Total number of parameters is **linear** in  $n$

## Example: Wumpus World

---

1,4	2,4	3,4	4,4
1,3	2,3	3,3	4,3
1,2 <b>B</b> <b>OK</b>	2,2	3,2	4,2
1,1 <b>OK</b>	2,1 <b>B</b> <b>OK</b>	3,1	4,1

$P_{ij} = true$  iff  $[i, j]$  contains a pit

$B_{ij} = true$  iff  $[i, j]$  is breezy

Include only  $B_{1,1}, B_{1,2}, B_{2,1}$  in the probability model

## Specifying the probability model

---

The full joint distribution is  $\mathbf{P}(P_{1,1}, \dots, P_{4,4}, B_{1,1}, B_{1,2}, B_{2,1})$

Apply product rule:  $\mathbf{P}(B_{1,1}, B_{1,2}, B_{2,1} \mid P_{1,1}, \dots, P_{4,4})\mathbf{P}(P_{1,1}, \dots, P_{4,4})$

(Do it this way to get  $P(\textit{Effect} \mid \textit{Cause})$ )

First term: 1 if pits are adjacent to breezes, 0 otherwise

Second term: pits are placed randomly, probability 0.2 per square:

$$\mathbf{P}(P_{1,1}, \dots, P_{4,4}) = \prod_{i,j=1,1}^{4,4} \mathbf{P}(P_{i,j}) = 0.2^n \times 0.8^{16-n}$$

for  $n$  pits

# Observations and query

---

We know the following facts:

$$b = \neg b_{1,1} \wedge b_{1,2} \wedge b_{2,1}$$

$$known = \neg p_{1,1} \wedge \neg p_{1,2} \wedge \neg p_{2,1}$$

Query is  $\mathbf{P}(P_{1,3}|known, b)$

Define  $Unknown = P_{ij}$ s other than  $P_{1,3}$  and  $Known$

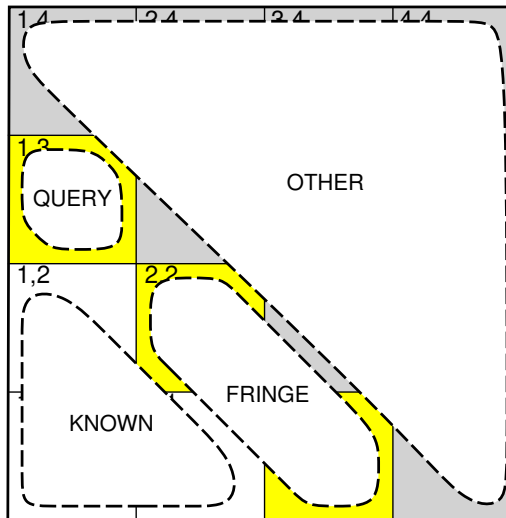
For inference by enumeration, we have

$$\mathbf{P}(P_{1,3}|known, b) = \alpha \sum_{unknown} \mathbf{P}(P_{1,3}, unknown, known, b)$$

Grows exponentially with number of squares

# Using conditional independence

Basic insight: observations are conditionally independent of other hidden squares, given neighboring hidden squares



Define  $Unknown = Fringe \cup Other$

$$\mathbf{P}(b|P_{1,3}, Known, Unknown) = \mathbf{P}(b|P_{1,3}, Known, Fringe)$$

Manipulate the query into a form where we can use this

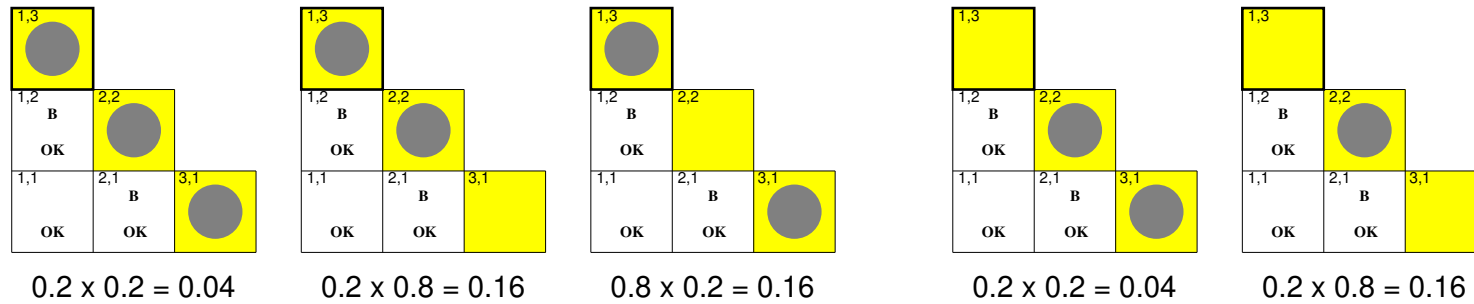


# Using conditional independence<sup>#</sup>

---

$$\begin{aligned}\mathbf{P}(P_{1,3}|\textit{known}, b) &= \alpha \sum_{\textit{unknown}} \mathbf{P}(P_{1,3}, \textit{unknown}, \textit{known}, b) \\ &= \alpha \sum_{\textit{unknown}} \mathbf{P}(b|P_{1,3}, \textit{known}, \textit{unknown})\mathbf{P}(P_{1,3}, \textit{known}, \textit{unknown}) \\ &= \alpha \sum_{\textit{fringe}} \sum_{\textit{other}} \mathbf{P}(b|\textit{known}, P_{1,3}, \textit{fringe}, \textit{other})\mathbf{P}(P_{1,3}, \textit{known}, \textit{fringe}, \textit{other}) \\ &= \alpha \sum_{\textit{fringe}} \sum_{\textit{other}} \mathbf{P}(b|\textit{known}, P_{1,3}, \textit{fringe})\mathbf{P}(P_{1,3}, \textit{known}, \textit{fringe}, \textit{other}) \\ &= \alpha \sum_{\textit{fringe}} \mathbf{P}(b|\textit{known}, P_{1,3}, \textit{fringe}) \sum_{\textit{other}} \mathbf{P}(P_{1,3}, \textit{known}, \textit{fringe}, \textit{other}) \\ &= \alpha \sum_{\textit{fringe}} \mathbf{P}(b|\textit{known}, P_{1,3}, \textit{fringe}) \sum_{\textit{other}} \mathbf{P}(P_{1,3})P(\textit{known})P(\textit{fringe})P(\textit{other}) \\ &= \alpha P(\textit{known})\mathbf{P}(P_{1,3}) \sum_{\textit{fringe}} \mathbf{P}(b|\textit{known}, P_{1,3}, \textit{fringe})P(\textit{fringe}) \sum_{\textit{other}} P(\textit{other}) \\ &= \alpha' \mathbf{P}(P_{1,3}) \sum_{\textit{fringe}} \mathbf{P}(b|\textit{known}, P_{1,3}, \textit{fringe})P(\textit{fringe})\end{aligned}$$

# Using conditional independence



$$\mathbf{P}(P_{1,3}|known, b) = \alpha' \langle 0.2(0.04 + 0.16 + 0.16), 0.8(0.04 + 0.16) \rangle$$

$$\approx \langle 0.31, 0.69 \rangle$$

$$\mathbf{P}(P_{2,2}|known, b) \approx \langle 0.86, 0.14 \rangle$$

# Bayesian networks

---

BNs: a graphical notation for conditional independence assertions and hence for compact specification of full joint distributions  
alias **Probabilistic Graphical Models (PGMs)**

## Syntax

a set of nodes, one per variable

a **directed acyclic graph** (DAG, link  $\rightarrow$  “directly influences”)

a conditional distribution for each node given its parents

$$\mathbf{P}(X_i | \text{Parents}(X_i))$$

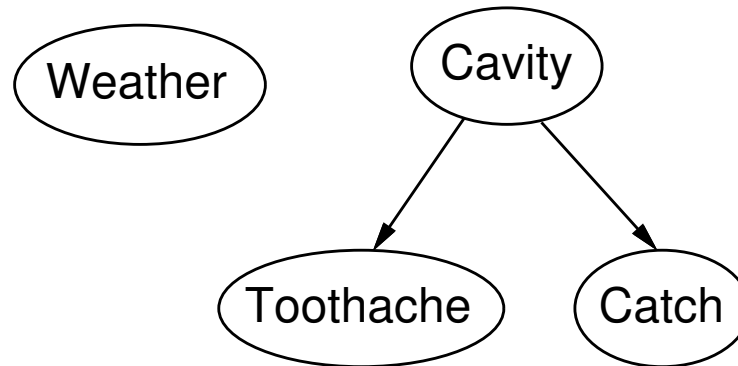
In the simplest case, conditional distribution is represented as a **conditional probability table (CPT)**

giving the distribution over  $X_i$  for each combination of parent values

## Example: Bayesian networks

---

Topology of the network encodes conditional independence assertions



*Weather* is independent of the other variables

*Toothache* and *Catch* are conditionally independent given *Cavity*

## Example: burglary network

---

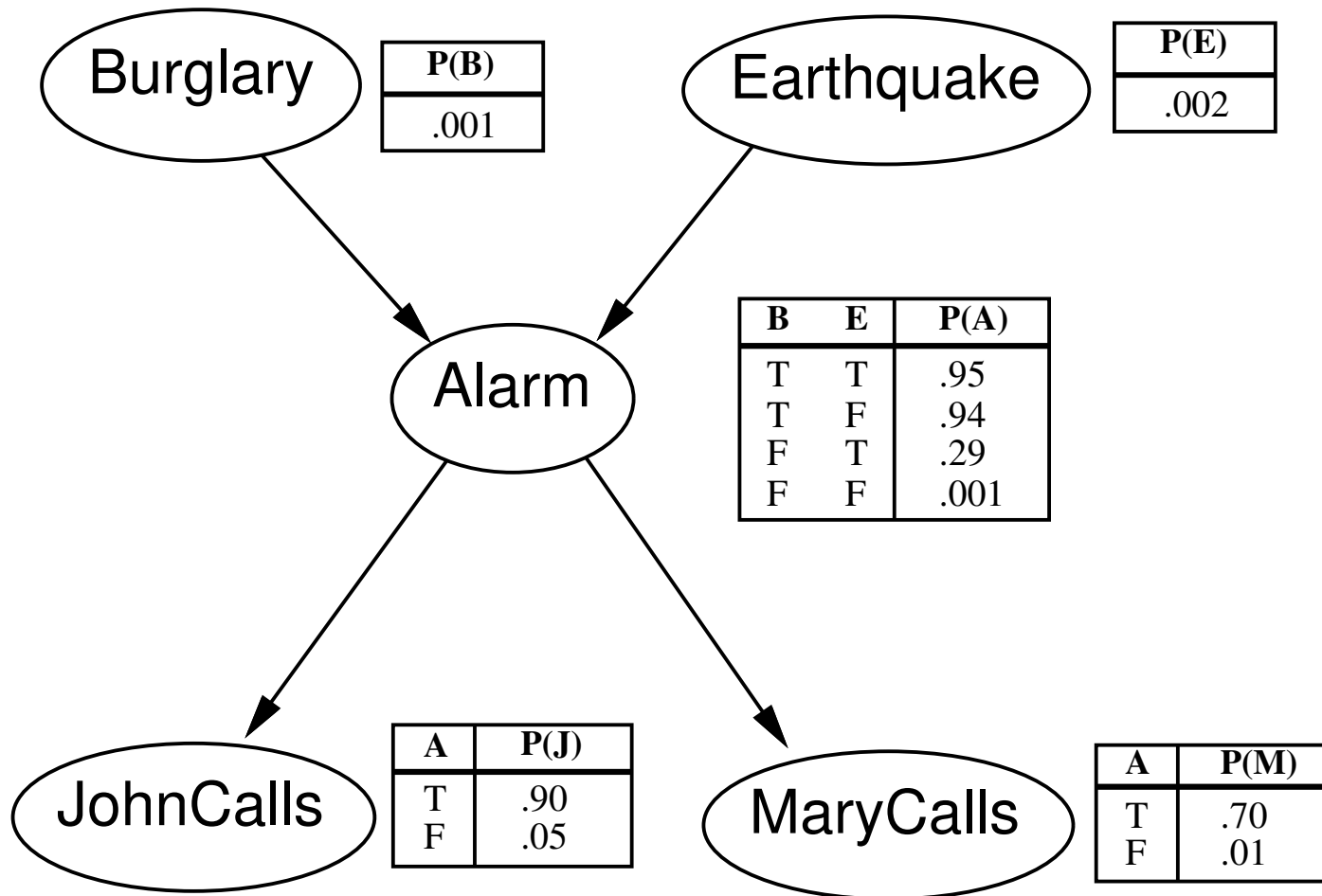
I'm at work, neighbor John calls to say my alarm is ringing, but neighbor Mary doesn't call. Sometimes it's set off by minor earthquakes. Is there a burglar?

Variables: *Burglar*, *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

# Example: burglary network



# Compactness

---

A CPT for Boolean  $X_i$  with  $k$  Boolean parents has  $2^k$  rows for the combinations of parent values

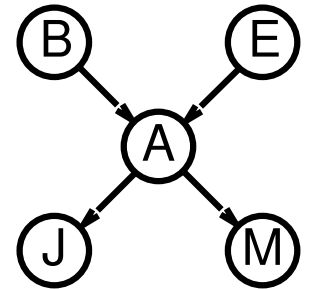
Each row requires one number  $p$  for  $X_i = true$  (the number for  $X_i = false$  is just  $1 - p$ )

If each variable 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.  $2^5 - 1 = 31$ )

In certain cases (assumptions of conditional independency), BNs make  $O(2^n) \Rightarrow O(kn)$  (NP  $\Rightarrow$  P !)



# Global semantics<sup>+</sup>

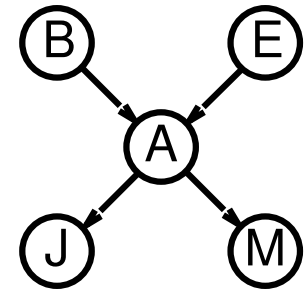
---

**Global** semantics defines the full joint distribution as the product of the local conditional distributions

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

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

=





# Global semantics<sup>+</sup>

---

Global semantics defines the full joint distribution as the product of the local conditional distributions

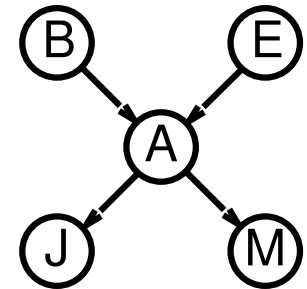
$$P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i | \text{parents}(X_i))$$

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

$$= P(j|a)P(m|a)P(a|\neg b, \neg e)P(\neg b)P(\neg e)$$

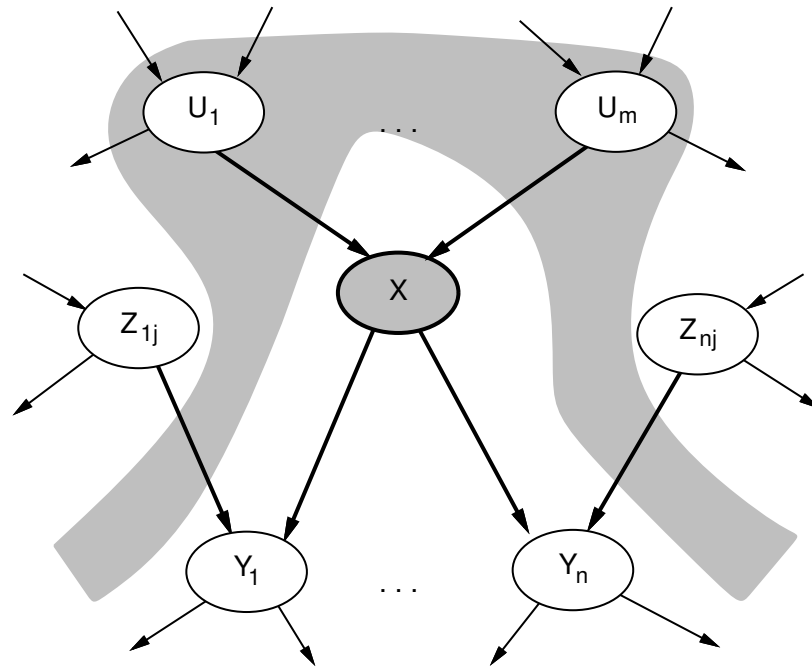
$$= 0.9 \times 0.7 \times 0.001 \times 0.999 \times 0.998$$

$$\approx 0.00063$$



# Local semantics<sup>+</sup>

**Local** semantics: each node is conditionally independent of its nondescendants ( $Z_{i,j}$ ) given its parents ( $U_i$  in the gray area)

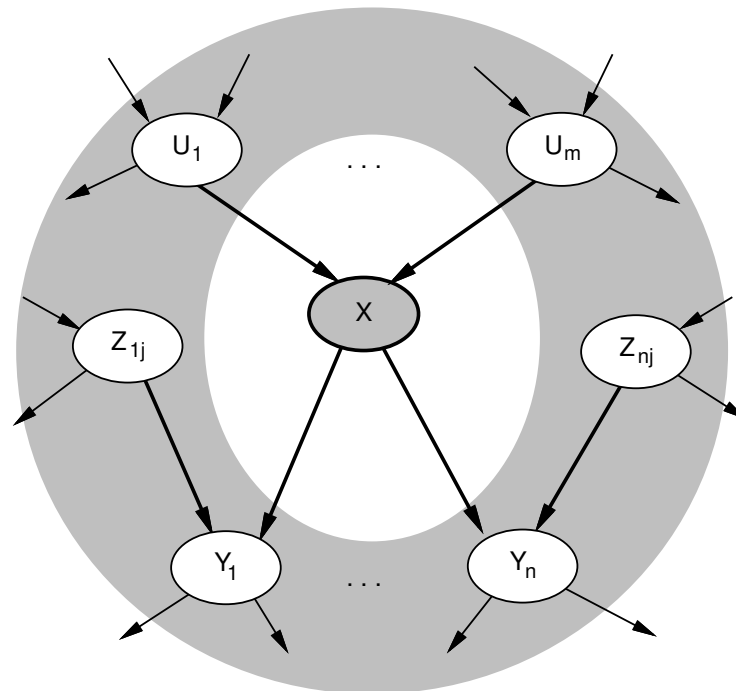


Theorem: Local semantics  $\Leftrightarrow$  global semantics

# Markov blanket<sup>+</sup>

---

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



# Constructing Bayesian networks

---

**Algorithm:** a series of locally testable assertions of conditional independence guarantees the required global semantics

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

This choice of parents guarantees the global semantics:

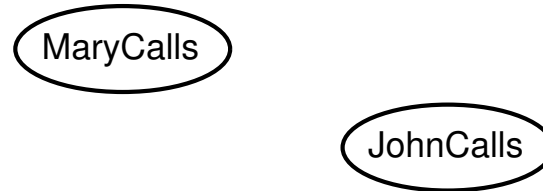
$$\begin{aligned}\mathbf{P}(X_1, \dots, X_n) &= \prod_{i=1}^n \mathbf{P}(X_i | X_1, \dots, X_{i-1}) \quad (\text{chain rule}) \\ &= \prod_{i=1}^n \mathbf{P}(X_i | Parents(X_i)) \quad (\text{by construction})\end{aligned}$$

Each node is conditionally independent of its other predecessors in the node (partial) ordering, given its parents

# Example: burglary network

---

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

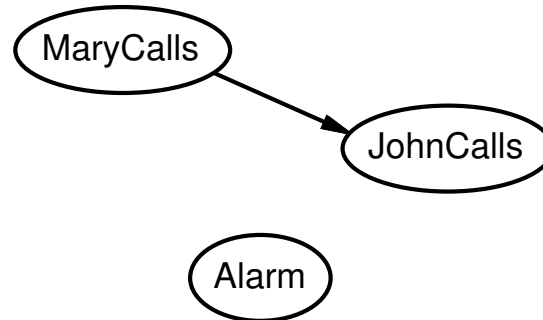


$$P(J|M) = P(J)?$$

# Example: burglary network

---

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



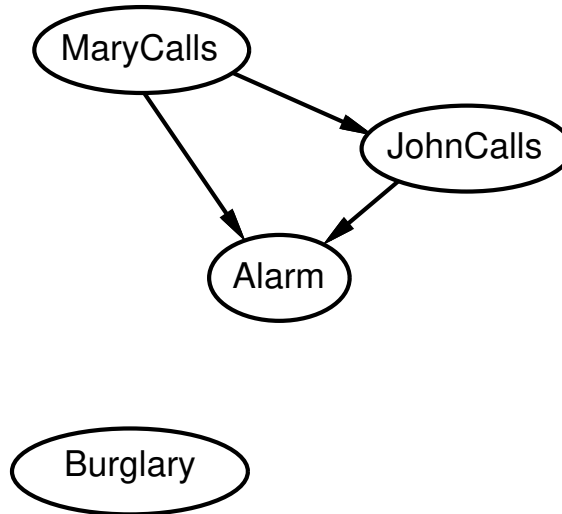
$P(J|M) = P(J)$ ? No

$P(A|J, M) = P(A|J)$ ?  $P(A|J, M) = P(A)$ ?

# Example: burglary network

---

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



$$P(J|M) = P(J)? \quad \text{No}$$

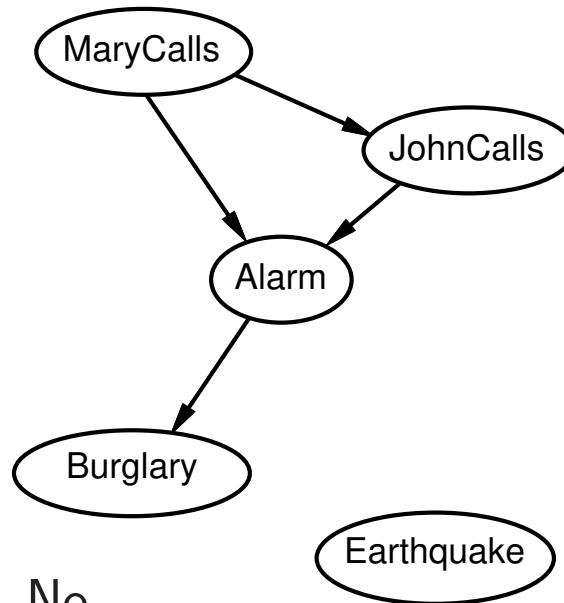
$$P(A|J, M) = P(A|J)? \quad P(A|J, M) = P(A)? \quad \text{No}$$

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

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

# Example: burglary network

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



$$P(J|M) = P(J)? \quad \text{No}$$

$$P(A|J, M) = P(A|J)? \quad P(A|J, M) = P(A)? \quad \text{No}$$

$$P(B|A, J, M) = P(B|A)? \quad \text{Yes}$$

$$P(B|A, J, M) = P(B)? \quad \text{No}$$

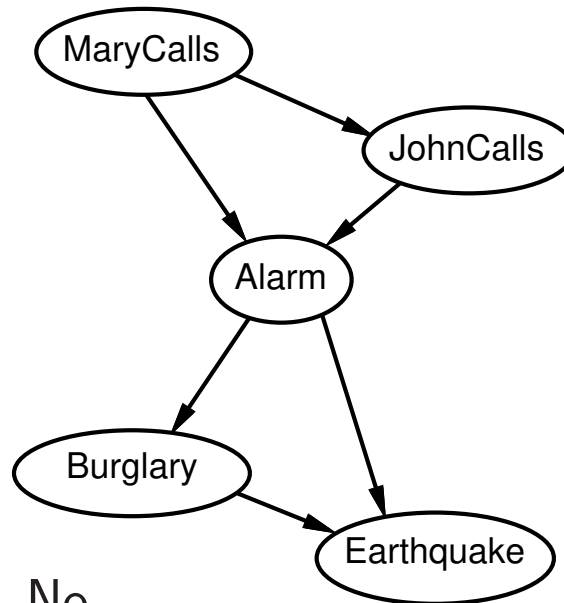
$$P(E|B, A, J, M) = P(E|A)?$$

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



# Example: burglary network

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



$$P(J|M) = P(J)? \quad \text{No}$$

$$P(A|J, M) = P(A|J)? \quad P(A|J, M) = P(A)? \quad \text{No}$$

$$P(B|A, J, M) = P(B|A)? \quad \text{Yes}$$

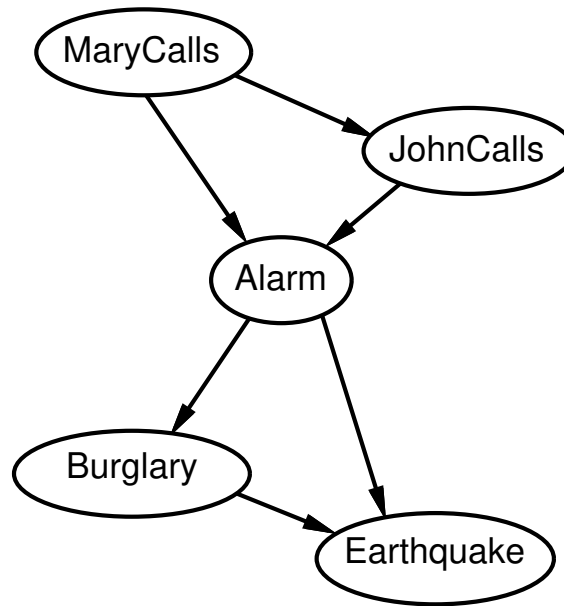
$$P(B|A, J, M) = P(B)? \quad \text{No}$$

$$P(E|B, A, J, M) = P(E|A)? \quad \text{No}$$

$$P(E|B, A, J, M) = P(E|A, B)? \quad \text{Yes}$$

## Example: burglary network

---



Assessing conditional probabilities is hard in noncausal directions  
The network can be far more compact than the full joint distribution  
But, this network is less compact:  $1 + 2 + 4 + 2 + 4 = 13$   
(due to the ordering of the variables)

# Probabilistic reasoning<sup>+</sup>

---

- Exact inference
  - enumeration
  - variable elimination
- Approximate inference\*
  - stochastic simulation
  - Markov chain Monte Carlo

## Reasoning tasks in BNs (PGMs)#

---

Simple queries: compute posterior marginal  $\mathbf{P}(X_i|\mathbf{E} = \mathbf{e})$

e.g.,  $P(\text{NoGas}|\text{Gauge} = \text{empty}, \text{Lights} = \text{on}, \text{Starts} = \text{false})$

Conjunctive queries:  $\mathbf{P}(X_i, X_j|\mathbf{E} = \mathbf{e}) = \mathbf{P}(X_i|\mathbf{E} = \mathbf{e})\mathbf{P}(X_j|X_i, \mathbf{E} = \mathbf{e})$

Optimal decisions: decision networks include utility information  
probabilistic inference required for  $P(\text{outcome}|\text{action}, \text{evidence})$

Value of information: which evidence to seek next?

Sensitivity analysis: which probability values are most critical?

Explanation/Causal inference: why do I need a nucleic acid detection (for coronavirus)?

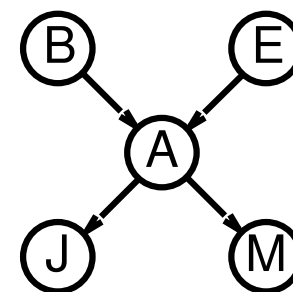
## Inference by enumeration

---

A slightly intelligent way to sum out variables from the joint without actually constructing its explicit representation

Simple query on the burglary network

$$\begin{aligned} & \mathbf{P}(B|j, m) \\ &= \mathbf{P}(B, j, m) / P(j, m) \\ &= \alpha \mathbf{P}(B, j, m) \\ &= \alpha \sum_e \sum_a \mathbf{P}(B, e, a, j, m) \end{aligned}$$



Rewrite full joint entries using product of CPT entries

$$\begin{aligned} & \mathbf{P}(B|j, m) \\ &= \alpha \sum_e \sum_a \mathbf{P}(B)P(e)\mathbf{P}(a|B, e)P(j|a)P(m|a) \\ &= \alpha \mathbf{P}(B) \sum_e P(e) \sum_a \mathbf{P}(a|B, e)P(j|a)P(m|a) \end{aligned}$$

Recursive depth-first enumeration:  $O(n)$  space,  $O(d^n)$  time

# Enumeration algorithm<sup>#</sup>

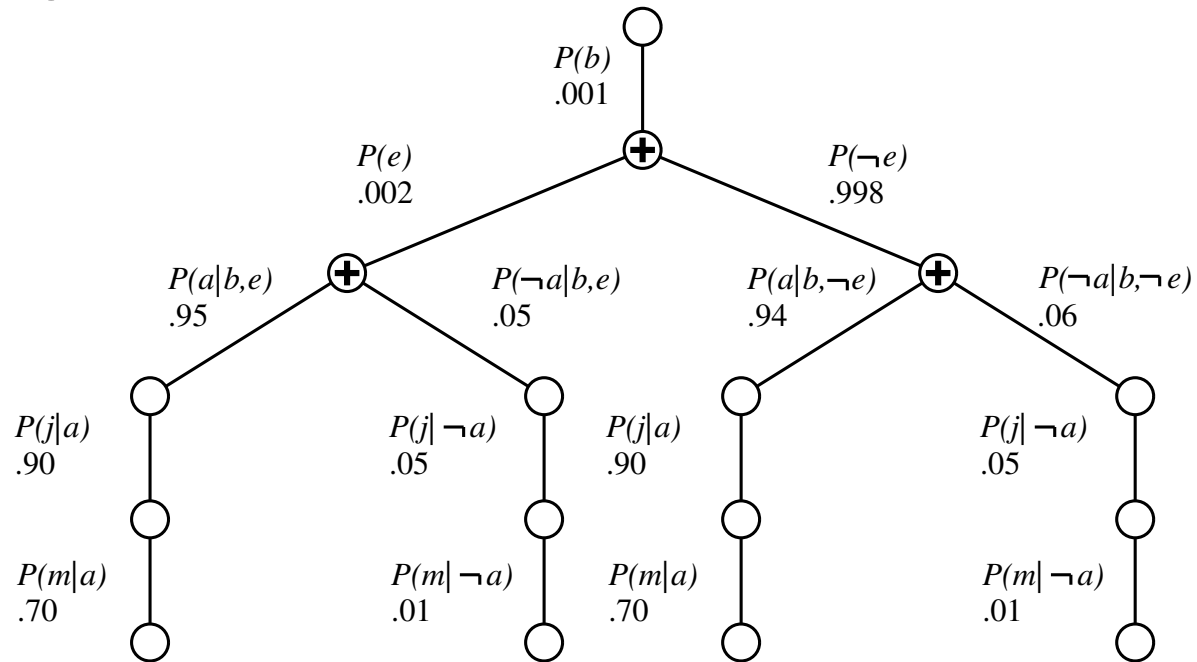
```
def ENUMERATION-ASK( $X, \mathbf{e}, bn$ )
inputs:  $X$ , the query variable
         $\mathbf{e}$ , observed values for variables  $\mathbf{E}$ 
         $bn$ , a Bayes net with variables  $vars$ 

 $Q(X) \leftarrow$  a distribution over  $X$ , initially empty
for each value  $x_i$  of  $X$  do
     $Q(x_i) \leftarrow$  ENUMERATE-ALL( $vars, \mathbf{e}_{x_i}$ )
        where  $\mathbf{e}_{x_i}$  is  $\mathbf{e}$  extended with  $X = x_i$ 
return NORMALIZE( $Q(X)$ ) // a distribution over  $X$ 

def ENUMERATE-ALL( $vars, \mathbf{e}$ )
if EMPTY?( $vars$ ) then return 1.0
 $V \leftarrow$  FIRST( $vars$ )
if  $V$  is an evidence variable with value  $v$  in  $\mathbf{e}$ 
then return  $P(v \mid parents(V)) \times$  ENUMERATE-ALL(REST( $vars$ ),  $\mathbf{e}$ )
else return  $\sum_v P(v \mid parents(V)) \times$  ENUMERATE-ALL(REST( $vars$ ),  $\mathbf{e}_v$ )
    where  $\mathbf{e}_v$  is  $\mathbf{e}$  extended with  $V = v$ 
```

# Evaluation tree

Summing at the “+” nodes



Enumeration is inefficient: repeated computation

e.g., computes  $P(j|a)P(m|a)$  for each value of  $e$   
 improved by eliminating repeated variables

# Inference by variable elimination

---

**Variable elimination:** carry out summations right-to-left, storing intermediate results (**factors**) to avoid recomputation

$$\begin{aligned} \mathbf{P}(B|j, m) &= \alpha \underbrace{\mathbf{P}(B)}_B \sum_e \underbrace{P(e)}_E \sum_a \underbrace{\mathbf{P}(a|B, e)}_A \underbrace{P(j|a)}_J \underbrace{P(m|a)}_M \\ &= \alpha \mathbf{P}(B) \sum_e P(e) \sum_a \mathbf{P}(a|B, e) P(j|a) f_M(a) \\ &= \alpha \mathbf{P}(B) \sum_e P(e) \sum_a \mathbf{P}(a|B, e) f_J(a) f_M(a) \\ &= \alpha \mathbf{P}(B) \sum_e P(e) \sum_a f_A(a, b, e) f_J(a) f_M(a) \\ &= \alpha \mathbf{P}(B) \sum_e P(e) f_{\bar{A}JM}(b, e) \text{ (sum out } A) \\ &= \alpha \mathbf{P}(B) f_{\bar{E}\bar{A}JM}(b) \text{ (sum out } E) \\ &= \alpha f_B(b) \times f_{\bar{E}\bar{A}JM}(b) \end{aligned}$$



## Variable elimination: Basic operations

---

Summing out a variable from a product of factors

move any constant factors outside the summation

add up submatrices in the pointwise product of remaining factors

$$\sum_x f_1 \times \cdots \times f_k = f_1 \times \cdots \times f_i \sum_x f_{i+1} \times \cdots \times f_k = f_1 \times \cdots \times f_i \times f_{\bar{X}}$$

assuming  $f_1, \dots, f_i$  do not depend on  $X$

Pointwise product of factors  $f_1$  and  $f_2$

$$\begin{aligned} f_1(x_1, \dots, x_j, y_1, \dots, y_k) \times f_2(y_1, \dots, y_k, z_1, \dots, z_l) \\ = f(x_1, \dots, x_j, y_1, \dots, y_k, z_1, \dots, z_l) \end{aligned}$$

$$\text{e.g., } f_1(a, b) \times f_2(b, c) = f(a, b, c)$$

# Variable elimination algorithm<sup>#</sup>

---

```
def ELIMINATION-ASK( $X, \mathbf{e}, bn$ )
  inputs:  $X$ , the query variable
          $\mathbf{e}$ , observed values for variables  $\mathbf{E}$ 
          $bn$ , a Bayes net with variables  $vars$ 

  factors  $\leftarrow []$ 
  for each  $var$  in ORDER( $vars$ ) do
    factors  $\leftarrow$  [MAKE-FACTOR( $V, \mathbf{e}$ )] + |factors
    if  $V$  is a hidden variable then factors  $\leftarrow$  SUM-OUT( $V, factors$ )
  return NORMALIZE(POINTWISEPRODUCT( $factors$ ))
```

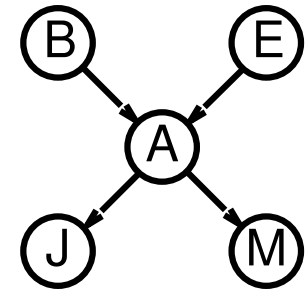
## Irrelevant variables\*

---

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

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

Sum over  $m$  is identically 1;  $M$  is **irrelevant** to the query



**Theorem:**  $Y$  is irrelevant unless  $Y \in \text{Ancestors}(\{X\} \cup \mathbf{E})$

Here,  $X = \text{JohnCalls}$ ,  $\mathbf{E} = \{\text{Burglary}\}$ , and  
 $\text{Ancestors}(\{X\} \cup \mathbf{E}) = \{\text{Alarm}, \text{Earthquake}\}$   
so  $\text{MaryCalls}$  is irrelevant

(Compare this to backward chaining from the query in Horn clause KBs)

## Irrelevant variables\*

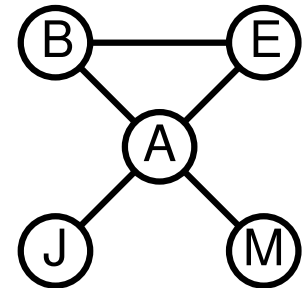
---

Defn: **moral graph** of BN: marry all parents and drop arrows

Defn: **A** is ***m-separated*** from **B** by **C** iff separated by **C** in the moral graph

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

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



# Complexity of exact inference

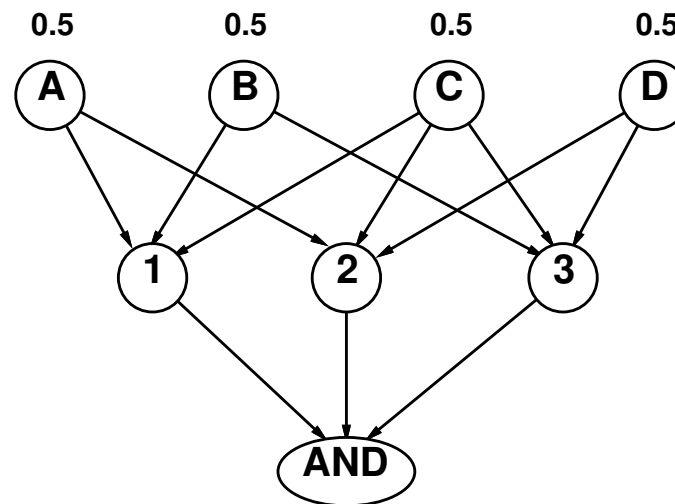
## Singly connected networks (or polytrees)

- any two nodes are connected by at most one (undirected) path
- time and space cost of variable elimination are  $O(d^k n)$

## Multiply connected networks

- can reduce 3SAT to exact inference  $\Rightarrow$  NP-hard
- equivalent to **counting** 3SAT models  $\Rightarrow$  #P-complete

1.  $A \vee B \vee C$
2.  $C \vee D \vee \sim A$
3.  $B \vee C \vee \sim D$

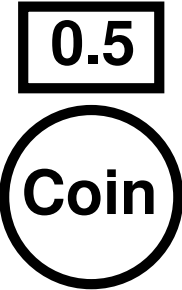


# Inference by stochastic simulation\*

---

## 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  $P$



## Methods

- 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

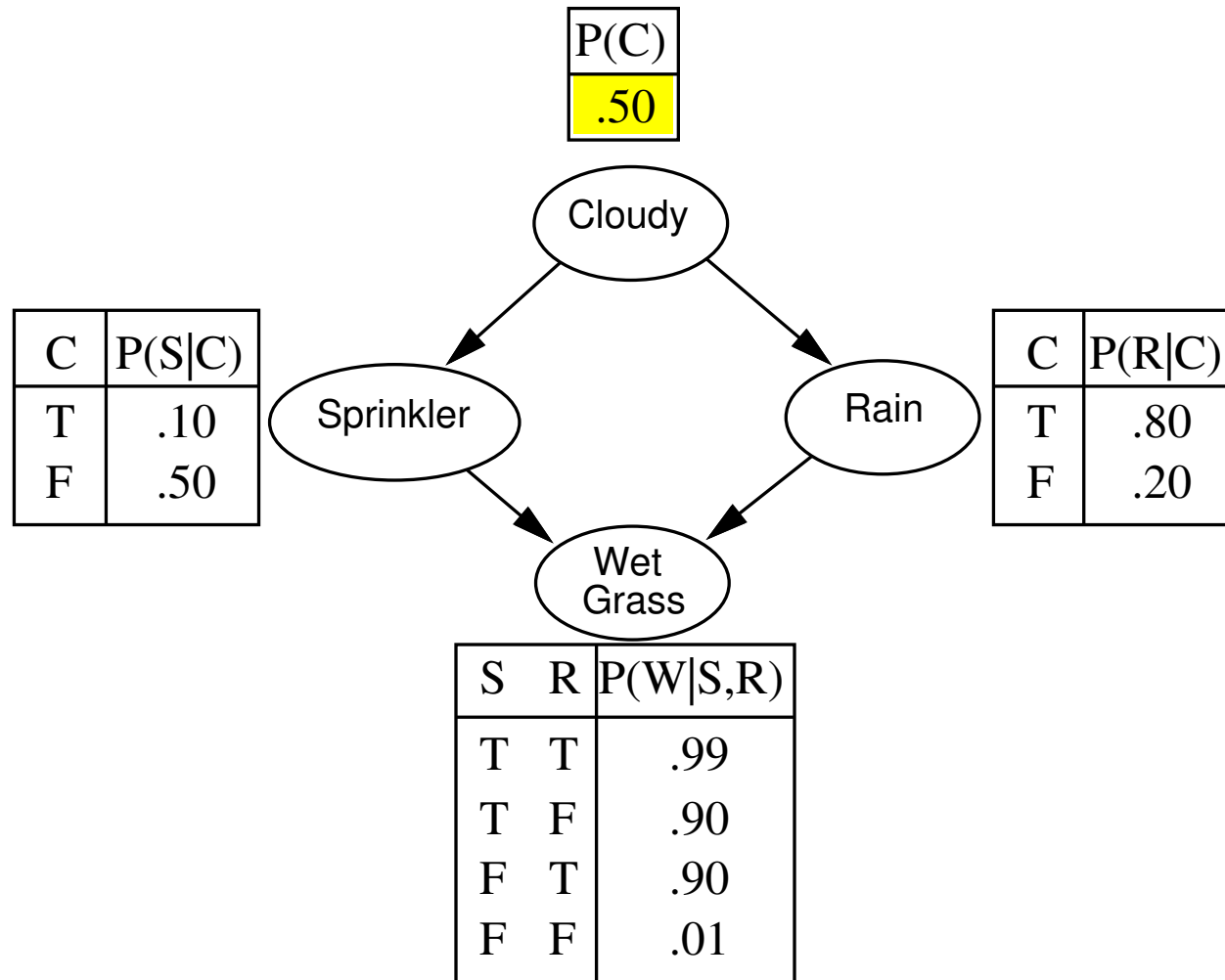
# Sampling from an empty network

---

Direct sampling from a network that has no evidence associated  
(sampling each variable in turn, in topological order)

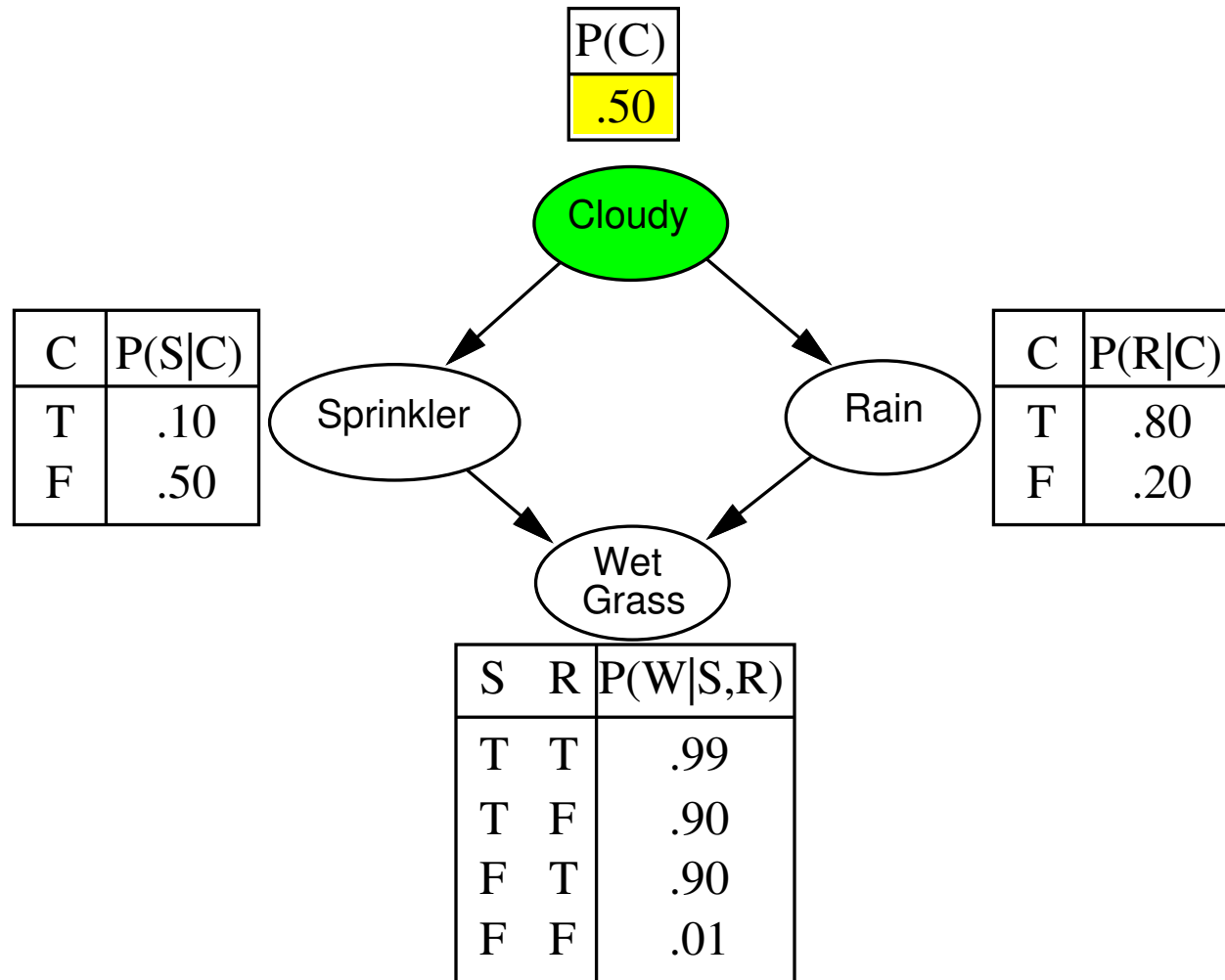
```
def PRIOR-SAMPLE(bn)
  inputs: bn, a BN specifying joint distribution  $\mathbf{P}(X_1, \dots, X_n)$ 
   $\mathbf{x} \leftarrow$  an event with  $n$  elements
  for each variable  $X_i$  in  $X_1, \dots, X_n$  do
     $\mathbf{x}[i] \leftarrow$  a random sample from  $\mathbf{P}(X_i \mid \text{Parents}(X_i))$ 
  return  $\mathbf{x}$  // an event sampled from the prior specified by bn
```

# Example: prior sampling

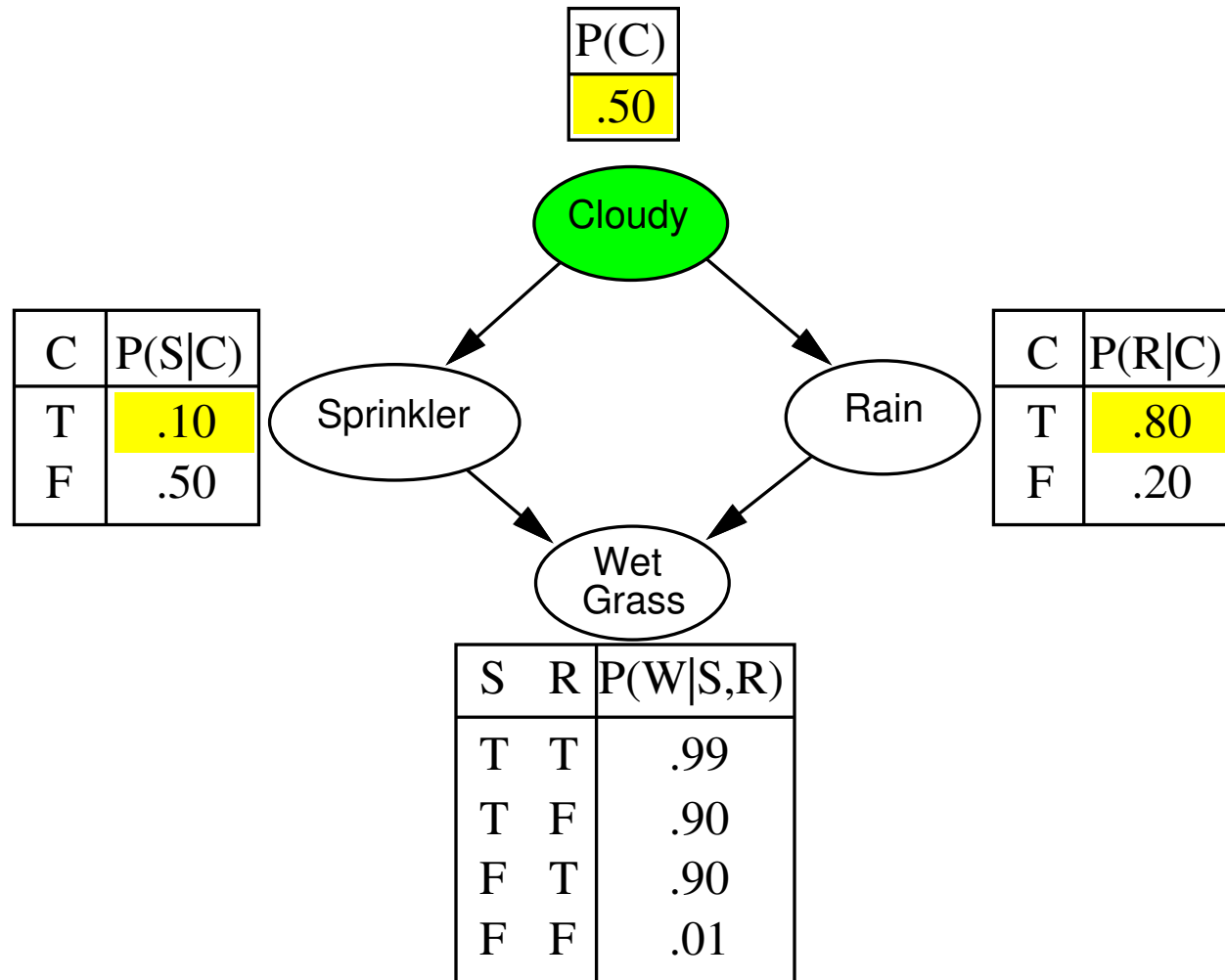




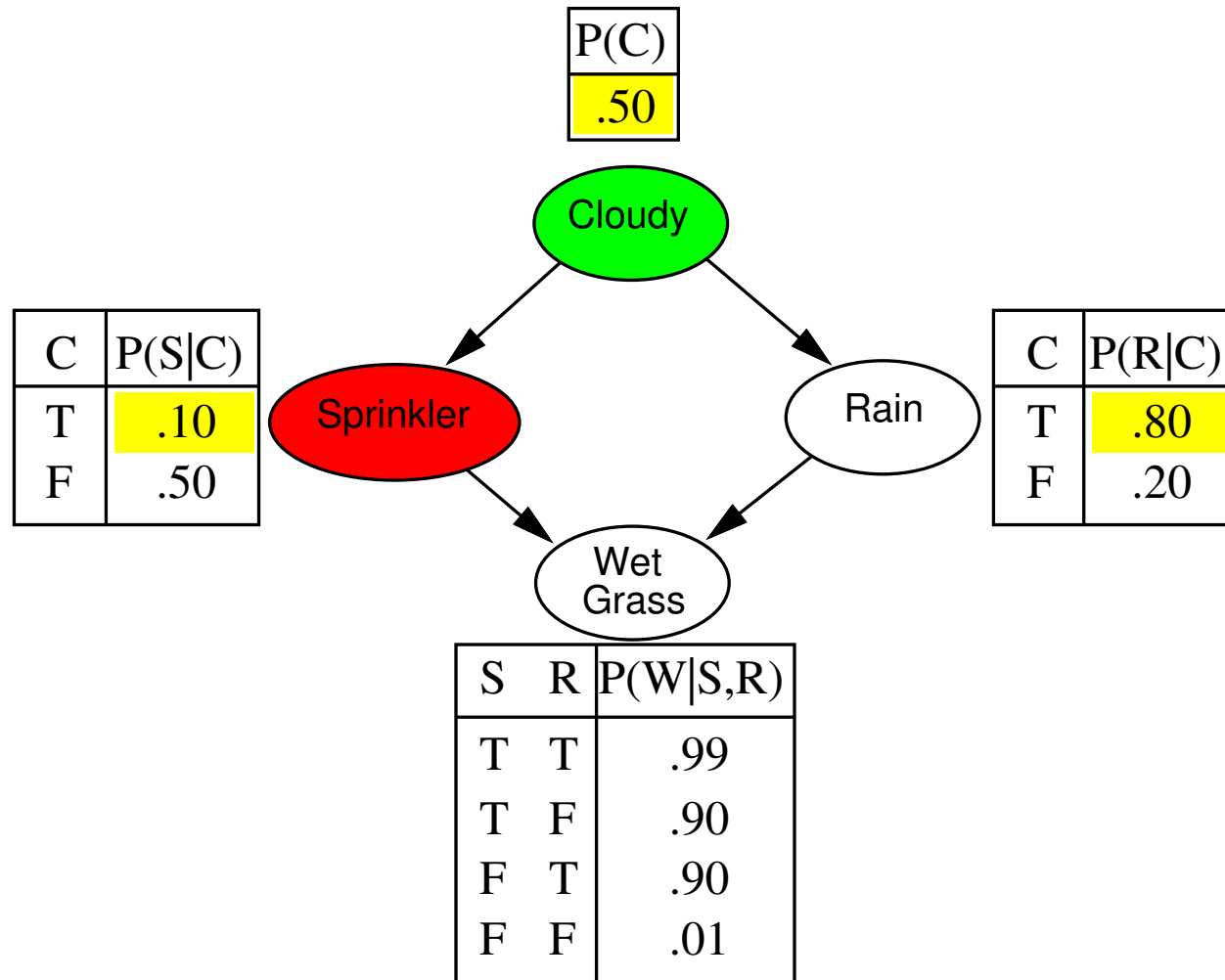
# Example: prior sampling



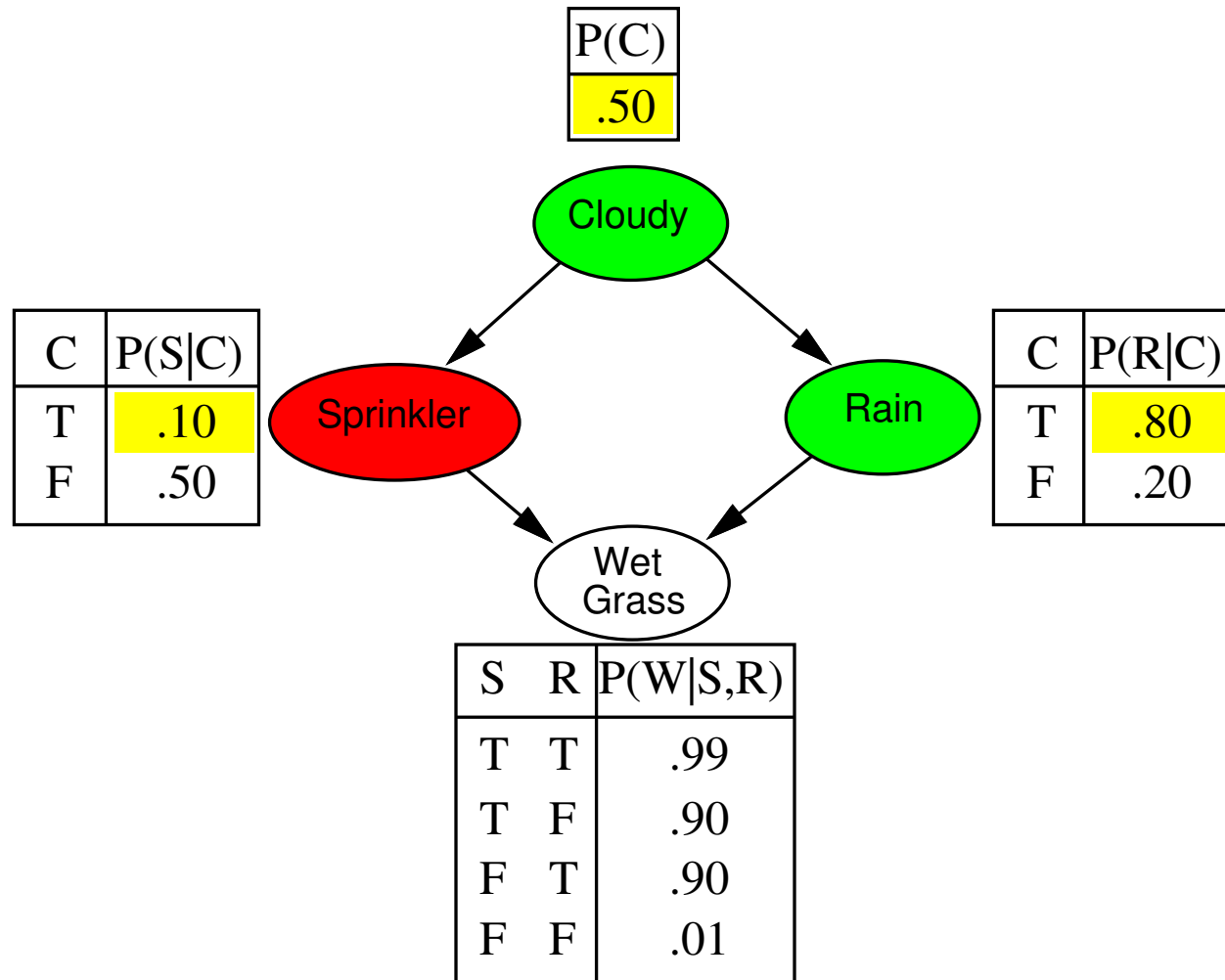
# Example: prior sampling



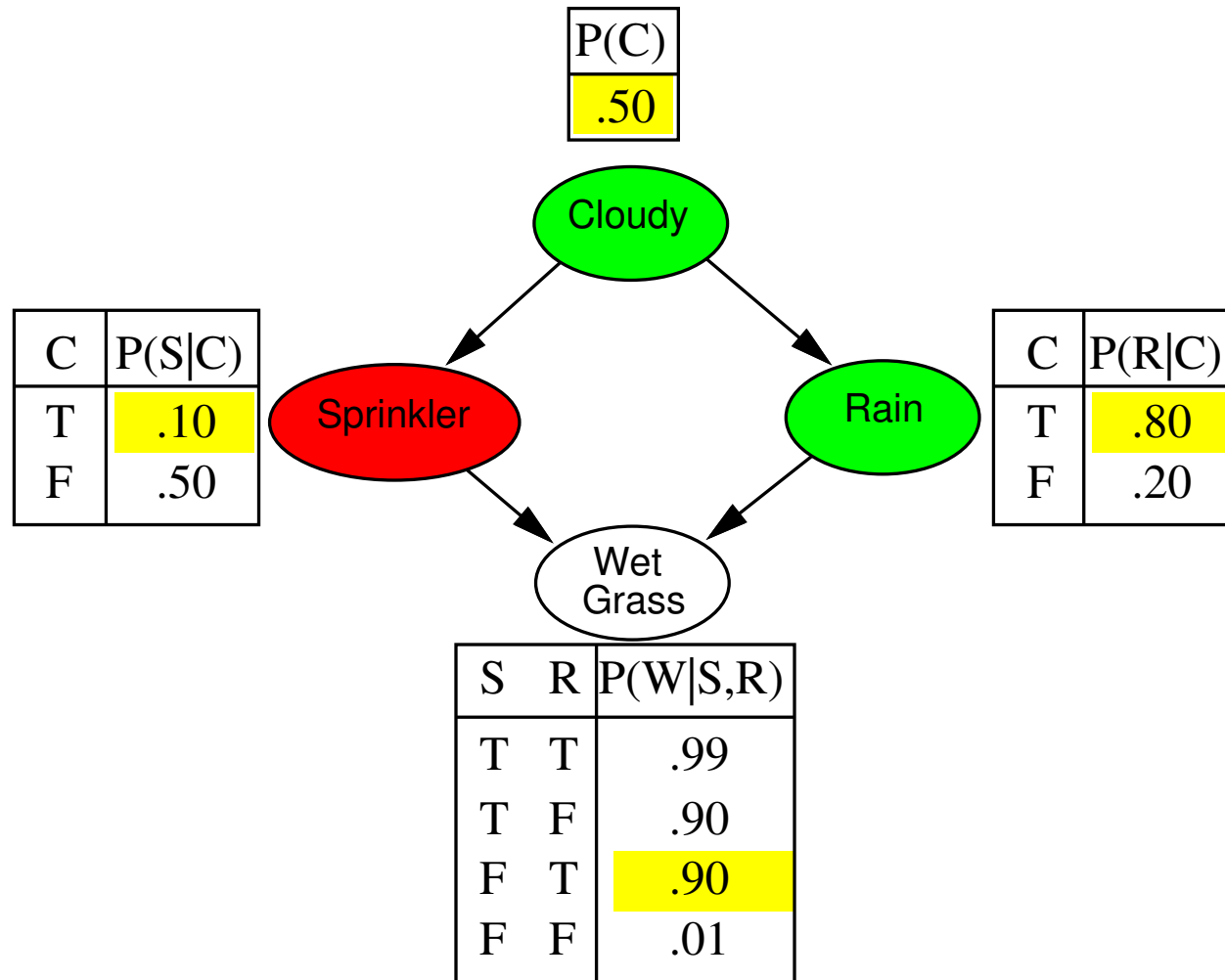
# Example: prior sampling



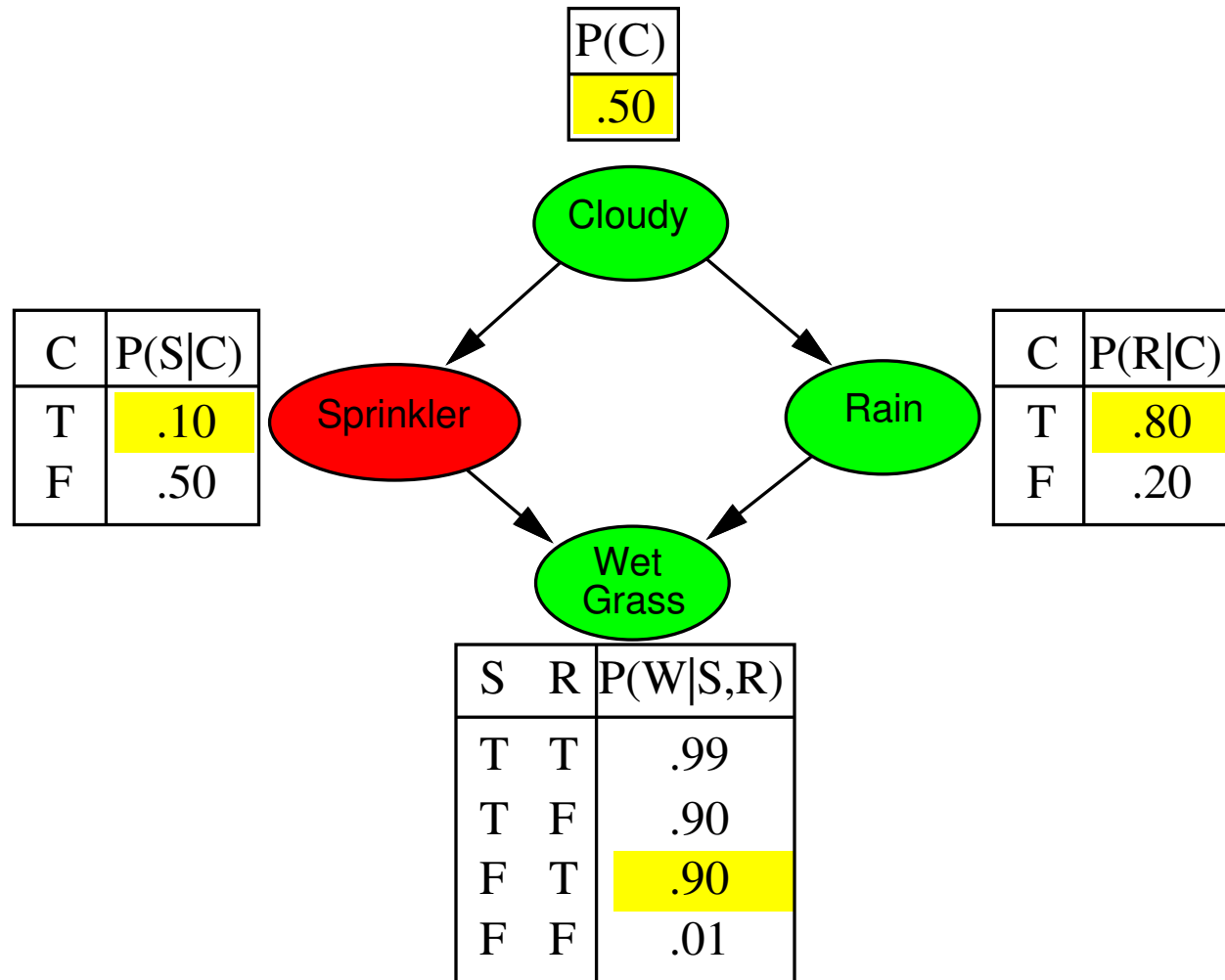
# Example: prior sampling



# Example: prior sampling



# Example: prior sampling



## Sampling from an empty network contd.

---

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

---

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

```
def REJECTION-SAMPLING( $X, \mathbf{e}, bn, N$ )  
  inputs:  $X$ , the query variable  
            $\mathbf{e}$ , observed values for variables  $E$   
            $bn$ , a BN  
            $N$ , the total number of samples to be generated  
  local variables:  $C$ , a vector of counts for each value of  $X$ , initially zero  
  for  $j = 1$  to  $N$  do  
     $\mathbf{x} \leftarrow$  PRIOR-SAMPLE( $bn$ )  
    if  $\mathbf{x}$  is consistent with  $\mathbf{e}$  then // do not match the evidence  
       $C[j] \leftarrow C[j] + 1$  where  $x_j$  is the value of  $X$  in  $\mathbf{x}$   
  return NORMALIZE( $C$ ) // an estimate of  $P(X|\mathbf{e})$ 
```



## Example: rejection sampling

---

Estimate  $\mathbf{P}(Rain|Sprinkler = true)$  using 100 samples

27 samples have  $Sprinkler = true$

Of these, 8 have  $Rain = true$  and 19 have  $Rain = false$ .

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

Similar to a basic real-world empirical estimation procedure

## Rejection sampling contd.

---

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

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

# Likelihood weighting

---

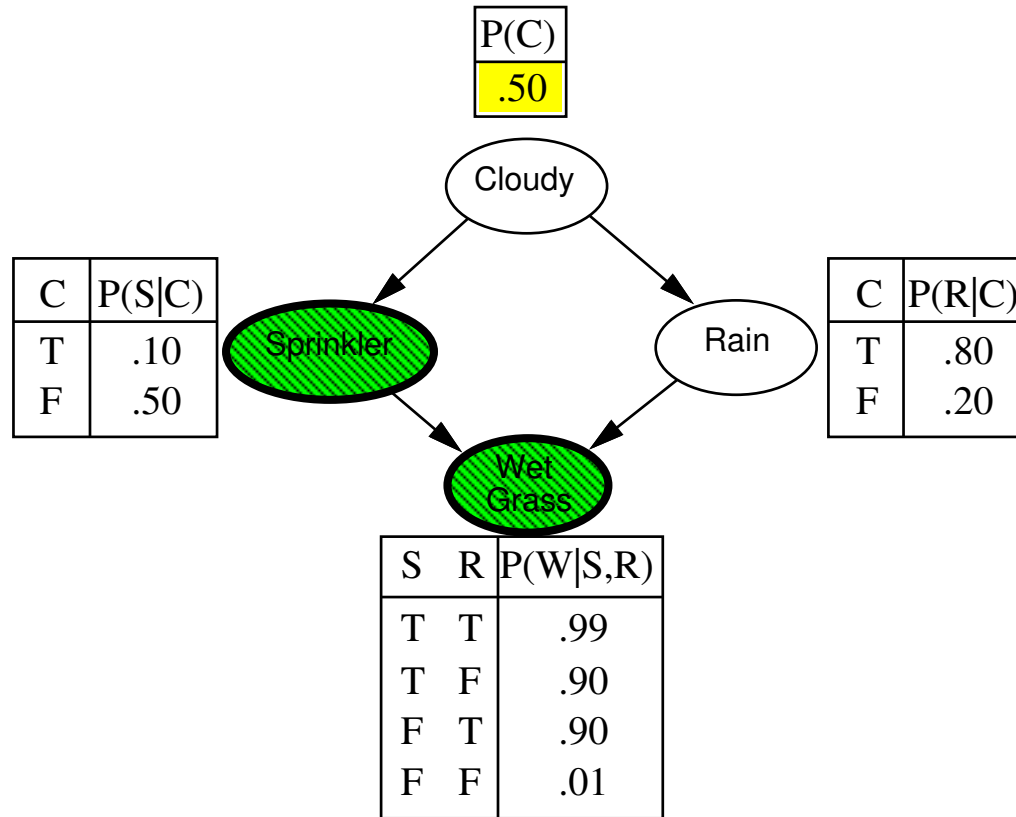
## Idea

- fix evidence variables
- sample only nonevidence variables
- weight each sample by the likelihood it accords the evidence

# Likelihood weighting

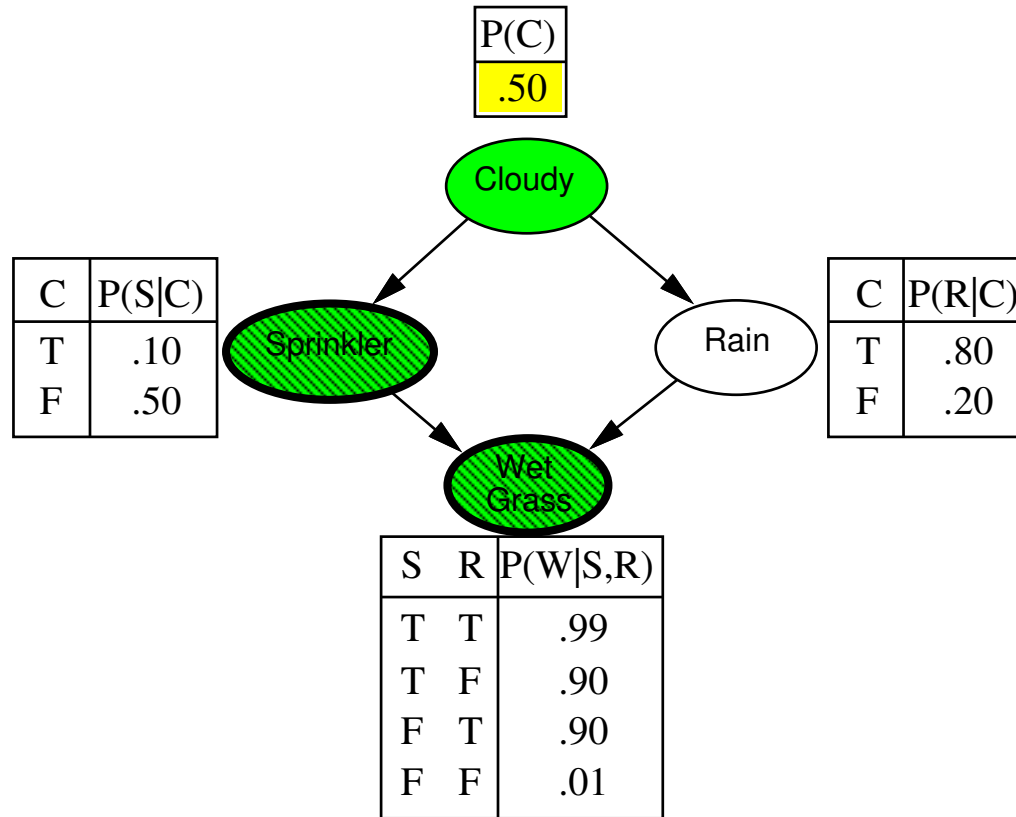
```
def LIKELIHOOD-WEIGHTING( $X, \mathbf{e}, bn, N$ )
  inputs:  $X, \mathbf{e}$ , the query variable, observed values for variables  $E$ 
           $bn, N$ , a BN, the total number of samples to be generated
  local variables:  $W$ , a vector of weighted counts for each value of  $X$ , initially 0
  for  $j = 1$  to  $N$  do
     $\mathbf{x}, w \leftarrow$  WEIGHTED-SAMPLE( $bn, \mathbf{e}$ )
     $W[j] \leftarrow W[j] + w$  where  $x_j$  is the value of  $X$  in  $\mathbf{x}$ 
  return NORMALIZE( $W[X]$ )
def WEIGHTED-SAMPLE( $bn, \mathbf{e}$ )
   $\mathbf{x} \leftarrow$  an event with  $n$  elements from  $\mathbf{e}$ ;  $w \leftarrow 1$ 
  for  $i=1$  to  $n$  do
    if  $X_i$  is an evidence variable with value  $x_i$  in  $\mathbf{e}$ 
      then  $w \leftarrow w \times P(X_i = x_{ij} \mid Parents(X_i))$ 
      else  $\mathbf{x}[i] \leftarrow$  a random sample from  $\mathbf{P}(X_i \mid Parents(X_i))$ 
  return  $\mathbf{x}, w$ 
```

# Example: likelihood weighting



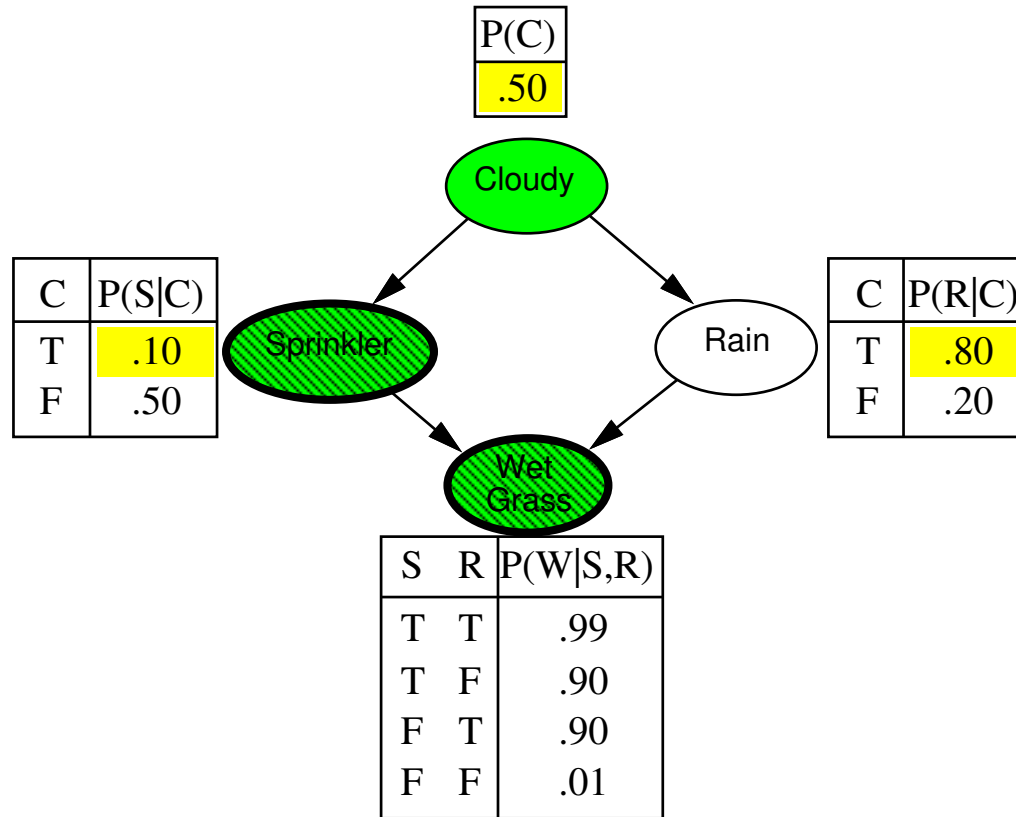
$$w = 1.0$$

# Example: likelihood weighting



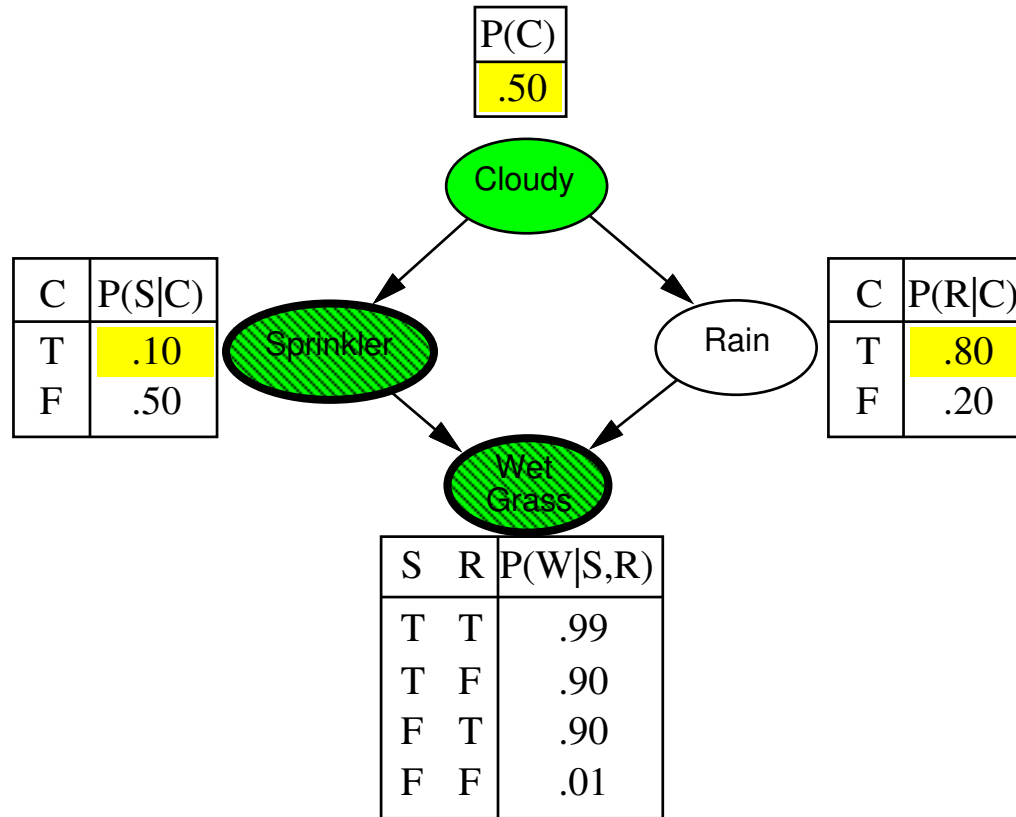
$$w = 1.0$$

# Example: likelihood weighting



$$w = 1.0$$

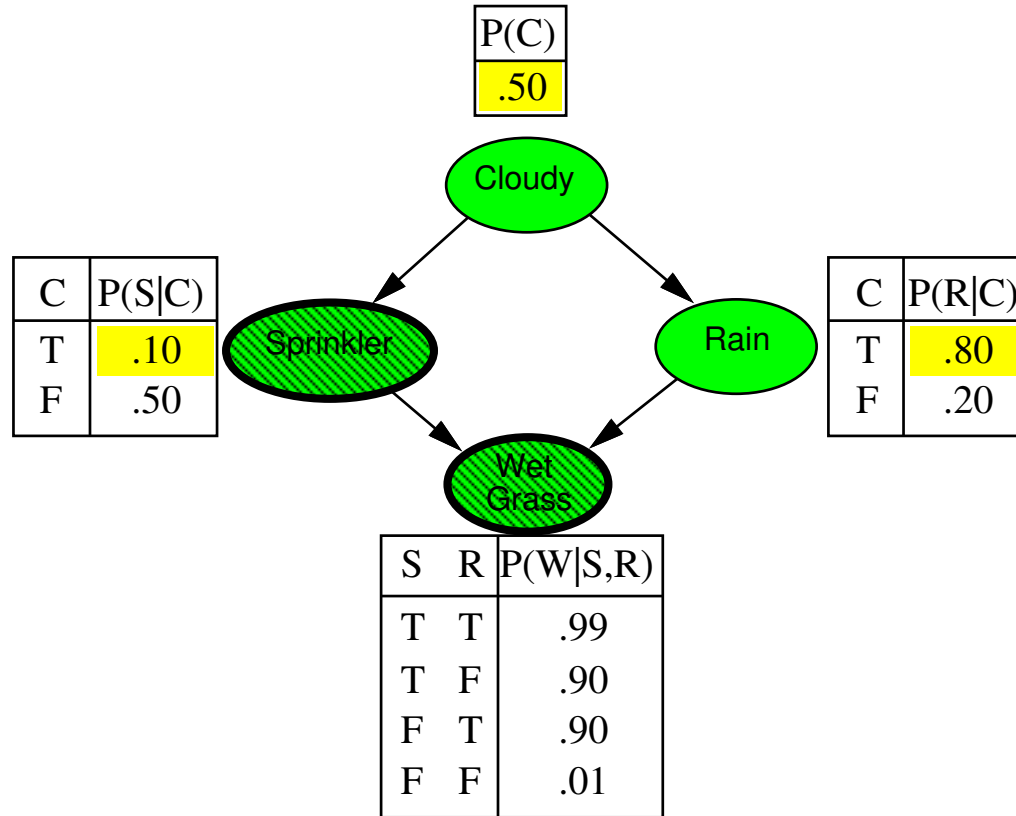
# Example: likelihood weighting



$$w = 1.0 \times 0.1$$

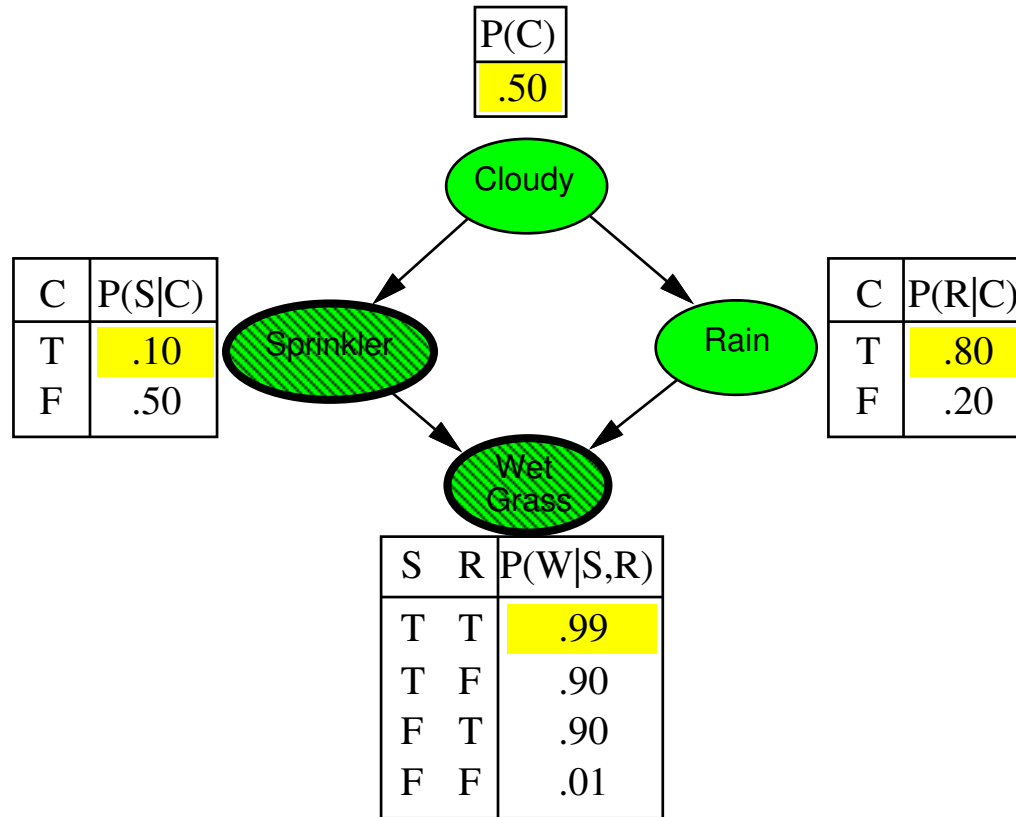


# Example: likelihood weighting



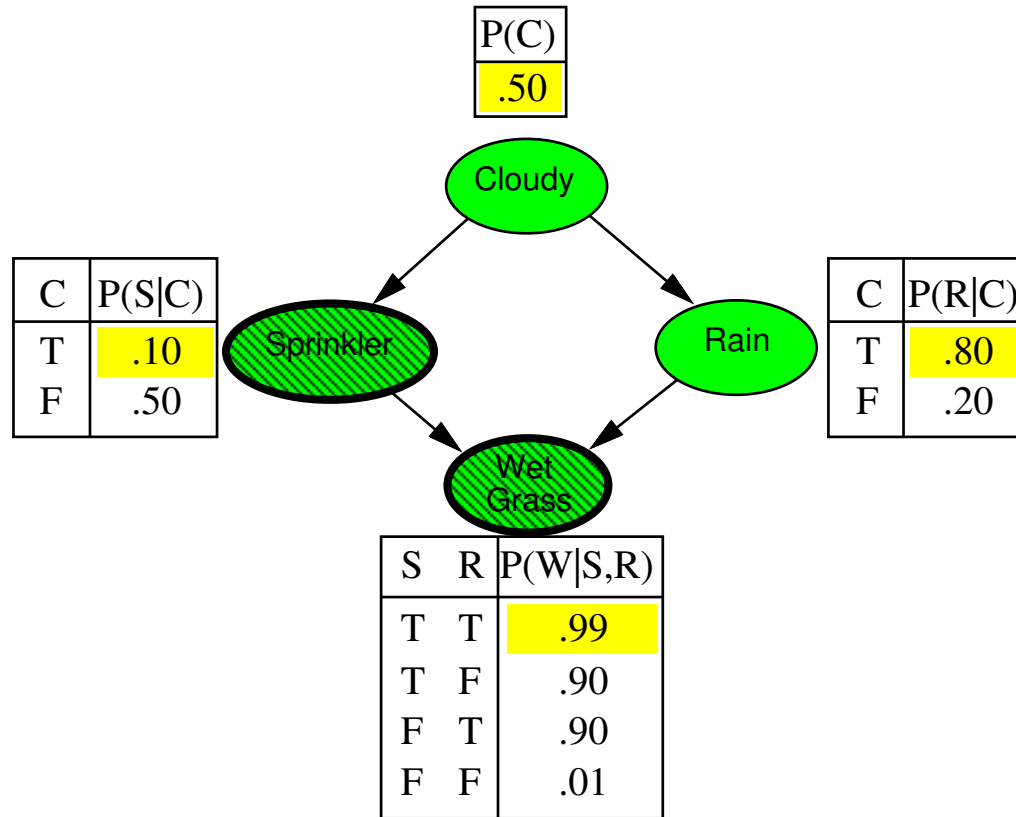
$$w = 1.0 \times 0.1$$

# Example: likelihood weighting



$$w = 1.0 \times 0.1$$

# Example: likelihood weighting



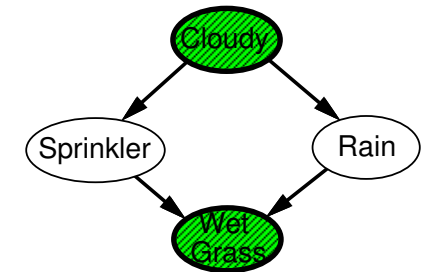
$$w = 1.0 \times 0.1 \times 0.99 = 0.099$$

## Likelihood weighting contd.

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  
 $\Rightarrow$  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

# Inference by Markov chain Monte Carlo\*

---

“State” of network = current assignment to all variables

⇒ the next state by making random changes to the current state

Generate the next state by sampling one variable given Markov blanket

recall Markov blanket: parents, children, and children’s parents

Sample each variable in turn, keeping evidence fixed

Specific **transition probability** with which the **stochastic process** moves from one state to another

defined by conditional distribution given Markov blanket of the variable being sampled

# MCMC Gibbs sampling

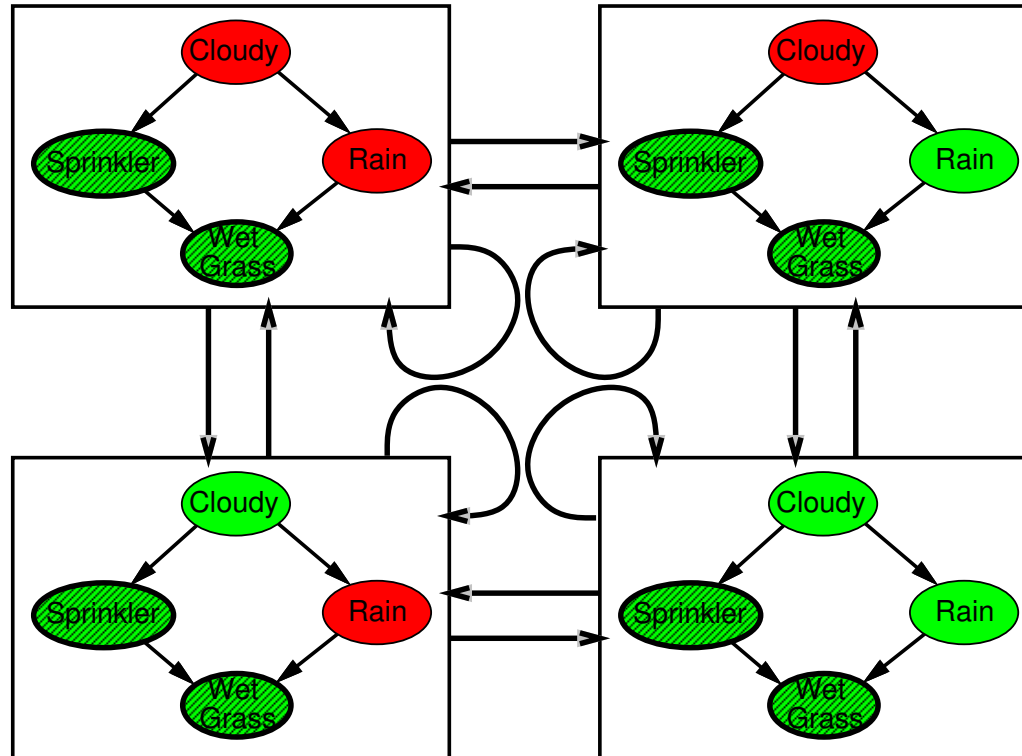
---

```
def MCMC-GIBBS-ASK( $X, \mathbf{e}, bn, N$ )
  local variables:  $C$ , a vector of counts for each value of  $X$ , initially zero
                   $Z$ , 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  $Z$ 
  for  $k = 1$  to  $N$  do // Can choose at random
    choose any variable  $Z_i$  from  $Z$  according to any distribution  $\rho(i)$ 
    set the value of  $Z_i$  in  $\mathbf{x}$  by sampling from  $\mathbf{P}(Z_i | mb(Z_i))$  // Markov blanket
     $C[j] \leftarrow C[j] + 1$  where  $x_j$  is the value of  $X$  in  $\mathbf{x}$ 
  return NORMALIZE( $C$ )
```

# Example: Markov chain

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



Wander about for a while

## Example: Gibbs sampling

---

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

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*

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

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



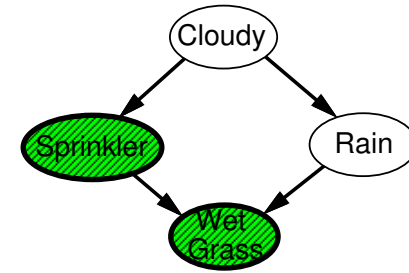
# Markov blanket sampling

---

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)

# Approximate inference

---

Exact inference by variable elimination:

- polytime on polytrees, NP-hard on general graphs
- space = time, very sensitive to topology

Approximate inference by LW (Likelihood Weighting), MCMC (Markov chain Monte Carlo):

- LW does poorly when there are lots of (downstream) evidence
- LW, MCMC generally insensitive to topology
- Convergence can be very slow with probabilities close to 1 or 0
- Can handle arbitrary combinations of discrete and continuous variables

# Dynamic Bayesian networks<sup>+</sup>

---

DBNs are BNs that represent **temporal** probability models

Basic idea: copy state and evidence variables for each time step

$\mathbf{X}_t$  = set of unobservable **state** variables at time  $t$   
e.g., *BloodSugar<sub>t</sub>*, *StomachContents<sub>t</sub>*, etc.

$\mathbf{E}_t$  = set of observable evidence variables at time  $t$   
e.g., *MeasuredBloodSugar<sub>t</sub>*, *PulseRate<sub>t</sub>*, *FoodEaten<sub>t</sub>*

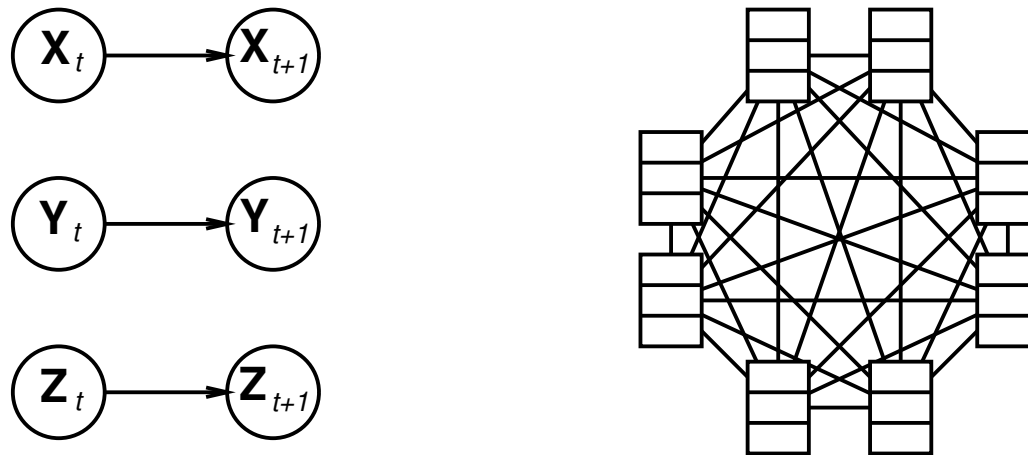
This assumes **discrete time**; step size depends on problem

Notation:  $\mathbf{X}_{a:b} = \mathbf{X}_a, \mathbf{X}_{a+1}, \dots, \mathbf{X}_{b-1}, \mathbf{X}_b$

$\mathbf{X}_t, \mathbf{E}_t$  contain arbitrarily many variables in a replicated Bayes net

# Hidden Markov models

HMMs: single-(state) variable DBNs  
every discrete DBN is an HMM  
(combine all the state variables in the DBN into a single one)



Sparse dependencies  $\Rightarrow$  exponentially fewer parameters;  
e.g., 20 state variables, three parents each  
DBN has  $20 \times 2^3 = 160$  parameters, HMM has  $2^{20} \times 2^{20} \approx 10^{12}$   
(analogous to BNs and full tabulated joint distributions)

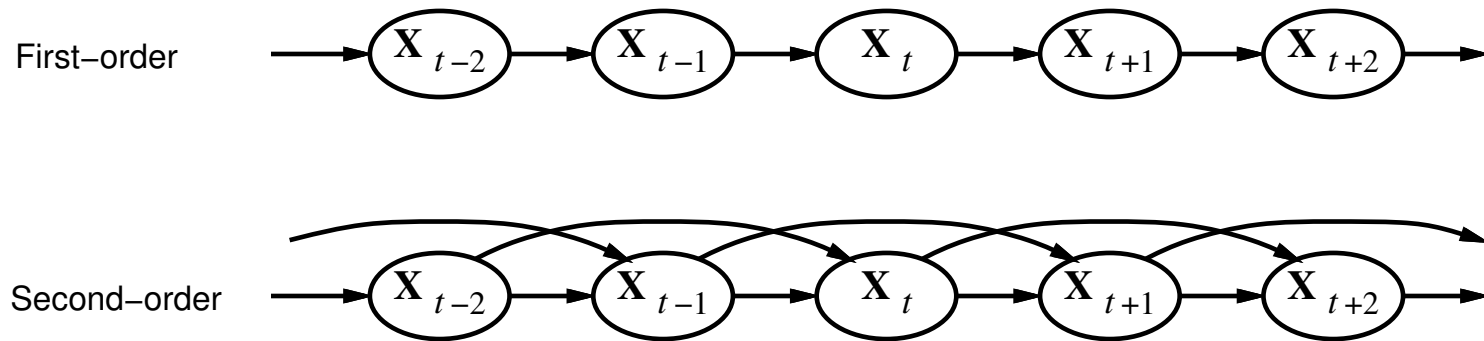
# Markov processes (Markov chains)

Construct a Bayes net from these variables: parents?

Markov assumption:  $\mathbf{X}_t$  depends on **bounded** subset of  $\mathbf{X}_{0:t-1}$

First-order Markov process:  $\mathbf{P}(\mathbf{X}_t | \mathbf{X}_{0:t-1}) = \mathbf{P}(\mathbf{X}_t | \mathbf{X}_{t-1})$

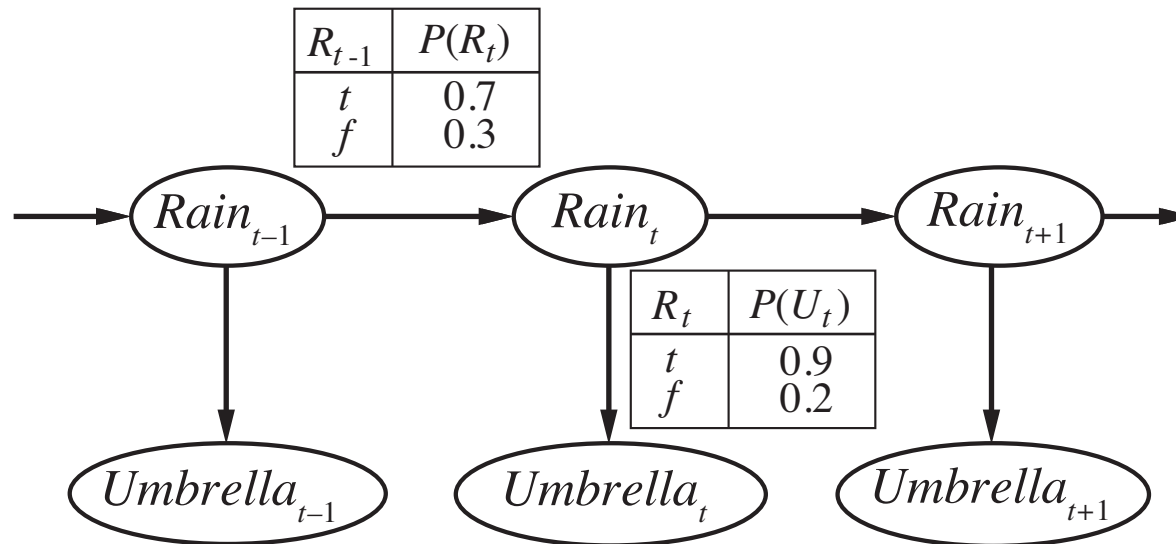
Second-order Markov process:  $\mathbf{P}(\mathbf{X}_t | \mathbf{X}_{0:t-1}) = \mathbf{P}(\mathbf{X}_t | \mathbf{X}_{t-2}, \mathbf{X}_{t-1})$



Sensor Markov assumption:  $\mathbf{P}(\mathbf{E}_t | \mathbf{X}_{0:t}, \mathbf{E}_{0:t-1}) = \mathbf{P}(\mathbf{E}_t | \mathbf{X}_t)$

Stationary process: transition model  $\mathbf{P}(\mathbf{X}_t | \mathbf{X}_{t-1})$  and sensor model  $\mathbf{P}(\mathbf{E}_t | \mathbf{X}_t)$  fixed for all  $t$

# Example: Markov processes



First-order Markov assumption not exactly true in real world

Possible fixes

1. **Increase order** of Markov process
2. **Augment state**, e.g., add  $Temp_t$ ,  $Pressure_t$

# HMMs

---

$X_t$  is a single, discrete variable (usually  $E_t$  is too)

Domain of  $X_t$  is  $\{1, \dots, S\}$

Transition matrix  $T_{ij} = P(X_t = j | X_{t-1} = i)$ , e.g.,  $\begin{pmatrix} 0.7 & 0.3 \\ 0.3 & 0.7 \end{pmatrix}$

Sensor matrix  $O_t$  for each time step, diagonal elements  $P(e_t | X_t = i)$

e.g., with  $U_1 = true$ ,  $O_1 = \begin{pmatrix} 0.9 & 0 \\ 0 & 0.2 \end{pmatrix}$

Forward and backward messages as column vectors

$$\mathbf{f}_{1:t+1} = \alpha \mathbf{O}_{t+1} \mathbf{T}^\top \mathbf{f}_{1:t}$$

$$\mathbf{b}_{k+1:t} = \mathbf{T} \mathbf{O}_{k+1} \mathbf{b}_{k+2:t}$$

Forward-backward algorithm needs time  $O(S^2t)$  and space  $O(St)$

# Inference tasks in HMMs

---

Filtering:  $\mathbf{P}(\mathbf{X}_t | \mathbf{e}_{1:t})$

belief state—input to the decision process of a rational agent

Prediction:  $\mathbf{P}(\mathbf{X}_{t+k} | \mathbf{e}_{1:t})$  for  $k > 0$

evaluation of possible action sequences;  
like filtering without the evidence

Smoothing:  $\mathbf{P}(\mathbf{X}_k | \mathbf{e}_{1:t})$  for  $0 \leq k < t$

better estimate of past states, essential for learning

Most likely explanation:  $\arg \max_{\mathbf{x}_{1:t}} P(\mathbf{x}_{1:t} | \mathbf{e}_{1:t})$

speech recognition, decoding with a noisy channel



# Filtering

---

Aim: devise a **recursive** state estimation algorithm

$$\mathbf{P}(\mathbf{X}_{t+1}|\mathbf{e}_{1:t+1}) = f(\mathbf{e}_{t+1}, \mathbf{P}(\mathbf{X}_t|\mathbf{e}_{1:t}))$$

$$\begin{aligned}\mathbf{P}(\mathbf{X}_{t+1}|\mathbf{e}_{1:t+1}) &= \mathbf{P}(\mathbf{X}_{t+1}|\mathbf{e}_{1:t}, \mathbf{e}_{t+1}) \\ &= \alpha \mathbf{P}(\mathbf{e}_{t+1}|\mathbf{X}_{t+1}, \mathbf{e}_{1:t}) \mathbf{P}(\mathbf{X}_{t+1}|\mathbf{e}_{1:t}) \\ &= \alpha \mathbf{P}(\mathbf{e}_{t+1}|\mathbf{X}_{t+1}) \mathbf{P}(\mathbf{X}_{t+1}|\mathbf{e}_{1:t})\end{aligned}$$

I.e., **prediction** + **estimation**. Prediction by summing out  $\mathbf{X}_t$

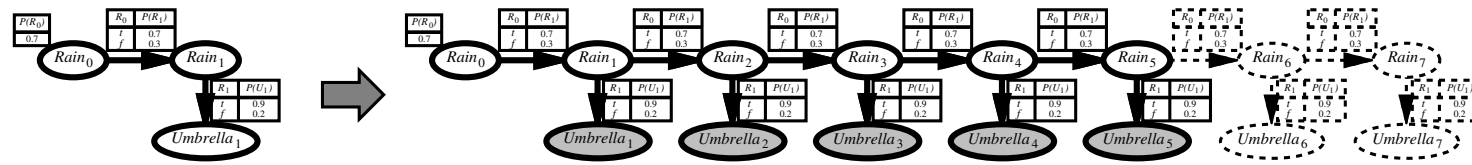
$$\begin{aligned}\mathbf{P}(\mathbf{X}_{t+1}|\mathbf{e}_{1:t+1}) &= \alpha \mathbf{P}(\mathbf{e}_{t+1}|\mathbf{X}_{t+1}) \sum_{\mathbf{x}_t} \mathbf{P}(\mathbf{X}_{t+1}|\mathbf{x}_t, \mathbf{e}_{1:t}) P(\mathbf{x}_t|\mathbf{e}_{1:t}) \\ &= \alpha \mathbf{P}(\mathbf{e}_{t+1}|\mathbf{X}_{t+1}) \sum_{\mathbf{x}_t} \mathbf{P}(\mathbf{X}_{t+1}|\mathbf{x}_t) P(\mathbf{x}_t|\mathbf{e}_{1:t})\end{aligned}$$

$\mathbf{f}_{1:t+1} = \text{FORWARD}(\mathbf{f}_{1:t}, \mathbf{e}_{t+1})$  where  $\mathbf{f}_{1:t} = \mathbf{P}(\mathbf{X}_t|\mathbf{e}_{1:t})$

Time and space **constant** (independent of  $t$ )

# Inference in DBNs

Naive method: **unroll** the network and run any exact algorithm



Problem: inference cost for each update grows with  $t$

**Rollup filtering**: add slice  $t + 1$ , “sum out” slice  $t$  using variable elimination

Largest factor is  $O(d^{n+1})$ , update cost  $O(d^{n+2})$   
(cf. HMM update cost  $O(d^{2n})$ )

Approximate inference by MCMC (Markov chain Monte Carlo) etc.

# Causal Inference\*

---

## Questions

- Observations: “What is we see  $A$ ?” (What is?)

$$P(y | A)$$

- Actions: “What if we do  $A$ ?” (What if?)

$$P(y | do(A))$$

- Counterfactuals: “What if we did things differerently?” (Why?)

$$P(y_{A'} | A)$$

E.g., recall  $C$ (limate)- $S$ (prinkler)- $R$ (rain)- $W$ (etness)

“Would the pavement be wet HAD the sprinkler been ON?”

$$(P(S | C) = 1)$$

Find if  $P(W_{S=1} = 1) = P(W = 1 | do(S = 1))$

Can drive counterfactuals from a model

# Graphical representations

---

- Observations → Bayesian networks
- Actions → Causal Bayesian networks
- Counterfactuals → Functional causal diagrams

## Hints

- Can reduce the action questions to symbolic calculus
- Can be estimated in polynomial time, complete algorithm (with the independence in the distribution)

# Probabilistic programming\*

---

Probability models: defined using executable code in any programming language that incorporates a source of randomness

⇒ Programs as probability models

⇐ probabilistic programming language (PPL)

– all of the expressive power of the underlying language

**Computationally universal:** they can represent any probability distribution that can be sampled from by a probabilistic Turing machine

# Generative program

---

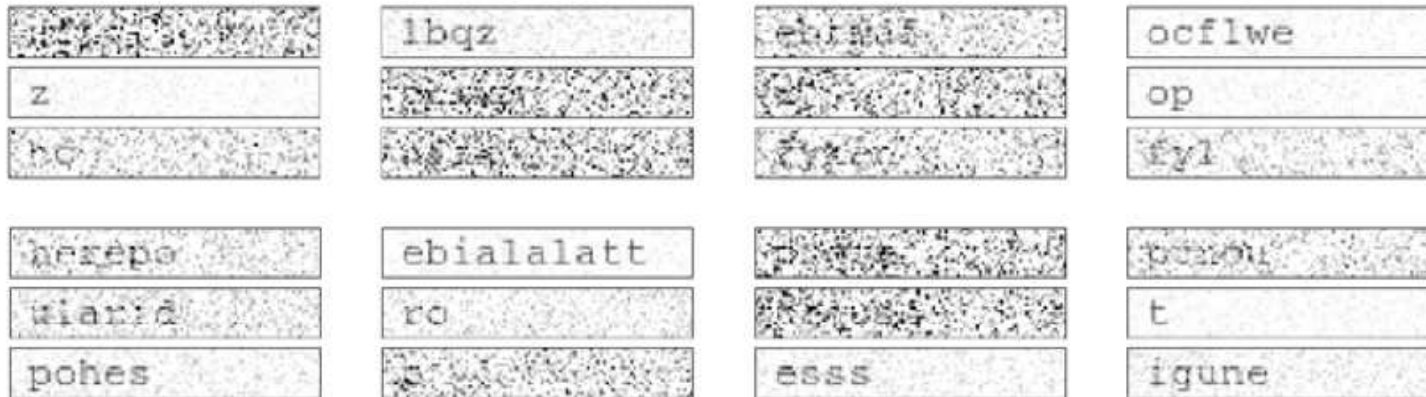
```
def GENERATE-IMAGE(())  
  letter ← GENERATE-LETTER(10)  
  return RENDER-NOISY-IMAGE(letter, 32, 128)  
def GENERATE-LETTER( $\lambda$ )  
   $n \sim \text{Poisson}(\lambda)$   
  for  $i=1$  to  $n$  do  
    letter  $\sim \text{UNIFORM-CHOICE}(\{a,b,c,\dots\})$   
  return letter  
def RENDER-NOISY-IMAGE(letter, width, height)  
  clean-image ← RENDER(letter, width, height, text-top=10, text-left=10)  
  noisy-image ← []  
  noise-variance  $\sim \text{UNIFORM-REAL}(0.1, 1)$   
  for row = 1 to width do  
    for col = 1 to height do  
      noisy-image[row, col]  $\sim \mathcal{N}(\text{clean-image}[\text{row}, \text{col}], \text{noise-variance})$   
  return noisy-image
```

## Example: reading text

---

The program that reads degraded (smudged or blurred) text (or CAPTCHAs), i.e., optical character recognition

– invoking GENERATE-IMAGE 9 times



# Probabilistic logic\*

---

Bayesian networks are essentially propositional:

- the set of random variables is fixed and finite
- each variable has a fixed domain of possible values

Probabilistic reasoning can be formalized as **probabilistic logic**

**First-order probabilistic logic** combines probability theory with the expressive power of first-order logic



# First-order probabilistic logic

---

Recall: Propositional probabilistic logic

- **Proposition** = disjunction of atomic events in which it is true
- **Possible world** (sample point)  $\omega$  = propositional logic model (an assignment of values to all of the r.v.s under consideration)
- $\omega \models \phi$ : for any proposition  $\phi$ , the  $\omega$  where it is true
- **probability model**: a set  $\Omega$  of possible worlds with a probability  $P(\omega)$  for each world  $\omega$

# First-order probabilistic logic

---

## FOPL

- Probability of any first-order logical sentence  $\phi$  as a sum over the possible worlds where it is true

$$P(\phi) = \sum_{\omega:\omega\models\phi} P(\omega)$$

- Conditional probabilities  $P(\phi|\mathbf{e})$  can be obtained similarly ask any question from the probability model  
 $\Rightarrow$  (first-order) belief networks

Problem: the set of first-order models is infinite

- the summation could be infeasible
- specifying a complete and consistent distribution over an infinite set of worlds could be very difficult

Analogous to the method of propositionalization for FOL

e.g. relational probability models (RPMs)

## Other approaches to uncertain reasoning

---

- Nonmonotonic reasoning
- Rule-based methods
- Dempster-Shafer theory
- Possibility theory
- Fuzzy logic
- Rough sets