

Probabilistic representation

Applied artificial intelligence (EDAI32)

Lecture “10”

2012-04-19

Elin A. Topp

Outline

- Uncertainty (chapter 13)
 - Uncertainty
 - Probability
 - Syntax and Semantics
 - Inference
 - Independence and Bayes' Rule
- Bayesian Networks (chapter 14.1-3)
 - Syntax
 - Semantics
 - Efficient representation of conditional distributions (parameterised distributions)

Outline

- Uncertainty (chapter 13)
 - Uncertainty
 - Probability
 - Syntax and Semantics
 - Inference
 - Independence and Bayes' Rule
- Bayesian Networks (chapter 14.1-3)
 - Syntax
 - Semantics
 - Efficient representation of conditional distributions (parameterised distributions)

Outline

- Uncertainty (chapter 13)
 - Uncertainty
 - Probability
 - Syntax and Semantics
 - Inference
 - Independence and Bayes' Rule
- Bayesian Networks (chapter 14.1-3)
 - Syntax
 - Semantics
 - Efficient representation of conditional distributions (parameterised distributions)

Uncertainty

Situation: Get to the airport in time for the flight (by car)

Action A_t := “Leave for airport t minutes before flight departs”

Question: will A_t get me there on time?

Deal with:

- 1) partial observability (road states, other drivers, ...)
- 2) noisy sensors (traffic reports)
- 3) uncertainty in action outcomes (flat tire, car failure, ...)
- 4) complexity of modeling and predicting traffic

Use pure logic? Well... :

1) risks falsehood: “ A_{25} will get me there on time”

or 2) leads to conclusions too weak for decision making:

“ A_{25} will get me there on time if there is no accident and it does not rain and my tires hold, and ...”

(A_{1440} would probably hold, but the waiting time would be intolerable, given the quality of airport food...)

Rational decision

*A*₂₅, *A*₉₀, *A*₁₈₀, *A*₁₄₄₀, ... what is “the right thing to do?”

Obviously dependent on relative importance of goals (being in time vs minimizing waiting time) AND on their respective likelihood of being achieved.

Uncertain reasoning: diagnosing a patient, i.e., find the CAUSE for the symptoms displayed.

“Diagnostic” rule: Toothache \Rightarrow Cavity ???

Complex rule: Toothache \Rightarrow Cavity \vee GumProblem \vee Abscess \vee ... ???

“Causal” rule: Cavity \Rightarrow Toothache ???

Rational decision

A₂₅, A₉₀, A₁₈₀, A₁₄₄₀, ... what is “the right thing to do?”

Obviously dependent on relative importance of goals (being in time vs minimizing waiting time) AND on their respective likelihood of being achieved.

Uncertain reasoning: diagnosing a patient, i.e., find the CAUSE for the symptoms displayed.

“Diagnostic” rule: Toothache \Rightarrow Cavity ??? No!

Complex rule: Toothache \Rightarrow Cavity \vee GumProblem \vee Abscess \vee ... ???

“Causal” rule: Cavity \Rightarrow Toothache ???

Rational decision

A_{25} , A_{90} , A_{180} , A_{440} , ... what is “the right thing to do?”

Obviously dependent on relative importance of goals (being in time vs minimizing waiting time) AND on their respective likelihood of being achieved.

Uncertain reasoning: diagnosing a patient, i.e., find the CAUSE for the symptoms displayed.

“Diagnostic” rule: Toothache \Rightarrow Cavity ??? No!

Complex rule: Toothache \Rightarrow Cavity \vee GumProblem \vee Abscess \vee ... ??? Too much!

“Causal” rule: Cavity \Rightarrow Toothache ???

Rational decision

A₂₅, A₉₀, A₁₈₀, A₁₄₄₀, ... what is “the right thing to do?”

Obviously dependent on relative importance of goals (being in time vs minimizing waiting time) AND on their respective likelihood of being achieved.

Uncertain reasoning: diagnosing a patient, i.e., find the CAUSE for the symptoms displayed.

“Diagnostic” rule: Toothache \Rightarrow Cavity ??? No!

Complex rule: Toothache \Rightarrow Cavity \vee GumProblem \vee Abscess \vee ... ??? Too much!

“Causal” rule: Cavity \Rightarrow Toothache ??? Well... not always

Using logic?

Fixing such “rules” would mean to make them logically exhaustive, but that is bound to fail due to:

Laziness (too much work to list all options)

Theoretical ignorance (there is simply no complete theory)

Practical ignorance (might be impossible to test exhaustively)

⇒ better use **probabilities** to represent certain **knowledge states**

⇒ Rational decisions (decision theory) combine probability and utility theory

Outline

- **Uncertainty (chapter 13)**
 - Uncertainty
 - **Probability**
 - Syntax and Semantics
 - Inference
 - Independence and Bayes' Rule
- Bayesian Networks (chapter 14.1-3)
 - Syntax
 - Semantics
 - Efficient representation of conditional distributions (parameterised distributions)

Probability

Probabilistic assertions summarise effects of

laziness: failure to enumerate exceptions, qualifications, etc.

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

Subjective or Bayesian probability:

Probabilities relate propositions to one's state of knowledge

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

Not claims of a “probabilistic tendency” in the current situation, but maybe learned from past experience of similar situations.

Probabilities of propositions change with new evidence:

e.g., $P(A_{25} \mid \text{no reported accidents, it's 5:00 in the morning}) = 0.15$

Making decisions under uncertainty

Suppose the following believes (from past experience):

$$P(A_{25} \text{ gets me there on time} \mid \dots) = 0.04$$

$$P(A_{90} \text{ gets me there on time} \mid \dots) = 0.70$$

$$P(A_{120} \text{ gets me there on time} \mid \dots) = 0.95$$

$$P(A_{1440} \text{ gets me there on time} \mid \dots) = 0.9999$$

Which action to choose?

Depends on my preferences for “missing flight” vs. “waiting (with airport cuisine)”, etc.

Utility theory is used to represent and infer preferences

Decision theory = utility theory + probability theory

Outline

- Uncertainty (chapter 13)
 - Uncertainty
 - Probability
- Syntax and Semantics
 - Inference
 - Independence and Bayes' Rule
- Bayesian Networks (chapter 14.1-3)
 - Syntax
 - Semantics
 - Efficient representation of conditional distributions (parameterised distributions)

Probability basics

A set Ω - the sample space, e.g., the 6 possible rolls of a die.

$\omega \in \Omega$ is a sample point / possible world / atomic event

A probability space of probability model is a sample space with an assignment $P(\omega)$ for every $\omega \in \Omega$ so that:

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

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

An event A is any subset of Ω

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

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, e.g., the reals or Booleans,

e.g., $\text{Odd}(1) = \text{true}$.

P induces a *probability distribution* for any random variable X

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

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

Propositions

A proposition describes the event (set of sample points) where it (the proposition) holds, i.e.,

Given Boolean random variables 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 $B(\omega) = \text{true}$

Often in AI applications, the sample points are defined by the values of a set of random variables, i.e., the sample space is the Cartesian product of the ranges of the variables.

Prior probability

Prior or unconditional probabilities of propositions

e.g., $P(\text{Cavity} = \text{true}) = 0.2$ and

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

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

Probability distribution gives values for all possible assignments (normalised):

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

Joint probability distribution for a set of (independent) random variables gives the probability of every atomic event on those random variables (i.e., every sample point):

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

Weather Cavity	sunny	rain	cloudy	snow
true	0.144	0.02	0.016	0.02
false	0.576	0.08	0.064	0.08

Posterior probability

Most often, there is *some* information, i.e., *evidence*, that one can base their belief on:

e.g., $P(\text{cavity}) = 0.2$ (prior, no evidence for anything), but

$$P(\text{cavity} \mid \text{toothache}) = 0.6$$

corresponds to belief *after the arrival of some evidence*
(also: *posterior* or *conditional probability*).

OBS: NOT “if *toothache*, then 60% chance of cavity”

THINK “given that *toothache* is all I know” instead!

Posterior probability

Most often, there is *some* information, i.e., *evidence*, that one can base their belief on:

e.g., $P(\text{cavity}) = 0.2$ (prior, no evidence for anything), but

$$P(\text{cavity} \mid \text{toothache}) = 0.6$$

corresponds to belief *after the arrival of some evidence*
(also: *posterior* or *conditional probability*).

OBS: NOT “if *toothache*, then 60% chance of cavity”

THINK “given that *toothache* is all I know” instead!

Evidence remains valid after more evidence arrives, but it might become less useful

Evidence may be completely useless, i.e., irrelevant.

$$P(\text{cavity} \mid \text{toothache, sunny}) = P(\text{cavity} \mid \text{toothache})$$

Domain knowledge lets us do this kind of inference.

Posterior probability (2)

Definition of conditional / posterior probability:

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

or as *Product rule* (for *a* and *b* being true, we need *b* true and then *a* true, given *b*):

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

and in general for whole distributions (e.g.):

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

(gives a 4x2 set of equations)

Chain rule (successive application of product rule):

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

Outline

- **Uncertainty (chapter 13)**
 - Uncertainty
 - Probability
 - Syntax and Semantics
- **Inference**
 - Independence and Bayes' Rule
- Bayesian Networks (chapter 14.1-3)
 - Syntax
 - Semantics
 - Efficient representation of conditional distributions (parameterised distributions)

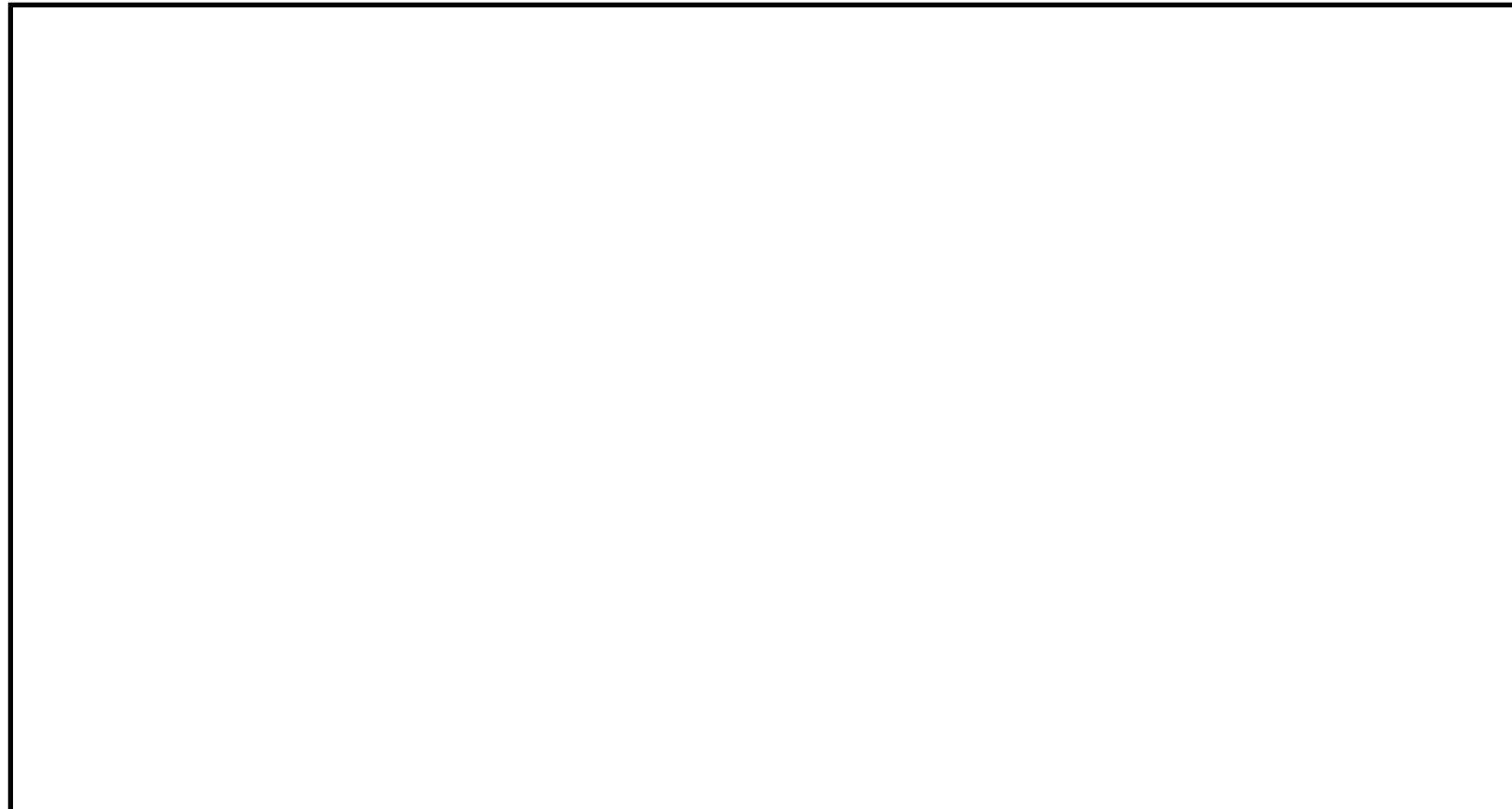
The suicidal student

A young student kills herself. Her diary is found. In the diary she speculates about her childhood and the possibility of her father abusing her during childhood. She had reported headaches to her friends and therapist, and started the diary due to the therapist's recommendation.

The father ends up in court, since

“headaches are caused by PTSD, and PTSD is caused by abuse”

What went wrong here?



The suicidal student

A young student kills herself. Her diary is found. In the diary she speculates about her childhood and the possibility of her father abusing her during childhood. She had reported headaches to her friends and therapist, and started the diary due to the therapist's recommendation.

The father ends up in court, since

“headaches are caused by PTSD, and PTSD is caused by abuse”

What went wrong here?

Psychologist knowing the math argues:

$P(\text{headache} \mid \text{PTSD}) = \text{high (statistics)}$

$P(\text{PTSD} \mid \text{abuse in childhood}) = \text{high (statistics)}$

but:

You do not know anything (in this case) of

$P(\text{PTSD} \mid \text{headache})$

$P(\text{abuse in childhood} \mid \text{headache})$

with only the evidence of headache and a speculation!

Inference

Probabilistic inference:

Computation of posterior probabilities given observed evidence

starting out with the full joint distribution as “knowledge base”:

Inference by enumeration

	toothache		¬ toothache	
	catch	¬ catch	catch	¬ catch
cavity	0.108	0.012	0.072	0.008
¬ cavity	0.016	0.064	0.144	0.576

For any proposition Φ , sum the atomic events where it is true:

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

Inference

Probabilistic inference:

Computation of posterior probabilities given observed evidence

starting out with the full joint distribution as “knowledge base”:

Inference by enumeration

	toothache		\neg toothache	
	catch	\neg catch	catch	\neg catch
cavity	0.108	0.012	0.072	0.008
\neg cavity	0.016	0.064	0.144	0.576

For any proposition Φ , sum the atomic events where it is true:

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

$$P(\text{toothache}) = 0.108 + 0.012 + 0.016 + 0.064 = 0.2$$

Inference

Probabilistic inference:

Computation of posterior probabilities given observed evidence

starting out with the full joint distribution as “knowledge base”:

Inference by enumeration

	toothache		\neg toothache	
	catch	\neg catch	catch	\neg catch
cavity	0.108	0.012	0.072	0.008
\neg cavity	0.016	0.064	0.144	0.576

For any proposition Φ , sum the atomic events where it is true:

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

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

Inference

Probabilistic inference:

Computation of posterior probabilities given observed evidence

starting out with the full joint distribution as “knowledge base”:

Inference by enumeration

	toothache		\neg toothache	
	catch	\neg catch	catch	\neg catch
cavity	0.108	0.012	0.072	0.008
\neg cavity	0.016	0.064	0.144	0.576

Can also compute posterior probabilities:

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

Normalisation

	toothache		\neg toothache	
	catch	\neg catch	catch	\neg catch
cavity	0.108	0.012	0.072	0.008
\neg cavity	0.016	0.064	0.144	0.576

Denominator can be viewed as a *normalisation constant*:

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

Normalisation

	toothache		\neg toothache	
	catch	\neg catch	catch	\neg catch
cavity	0.108	0.012	0.072	0.008
\neg cavity	0.016	0.064	0.144	0.576

Denominator can be viewed as a *normalisation constant*:

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

And the good news:

We can compute $P(\text{Cavity} \mid \text{toothache})$ without knowing the value of $P(\text{toothache})$!

... but

n Boolean variables give us an input table of size $O(2^n)$...

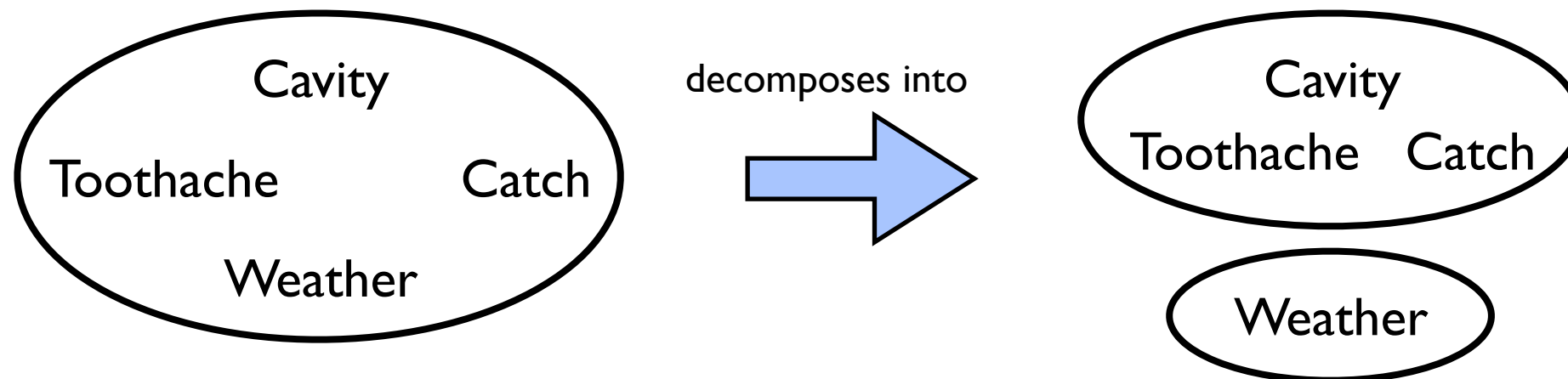
Outline

- **Uncertainty (chapter 13)**
 - Uncertainty
 - Probability
 - Syntax and Semantics
 - Inference
- **Independence and Bayes' Rule**
- **Bayesian Networks (chapter 14.1-3)**
 - Syntax
 - Semantics
 - Efficient representation of conditional distributions (parameterised distributions)

Independence

A and B are independent iff

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



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

32 entries reduced to 8 + 4. This absolute independence is powerful but rare!

Some fields (like dentistry) have still a lot, maybe hundreds, of variables, none of them being independent.

What can be done to overcome this mess...?

Conditional independence

$P(\text{Toothache}, \text{Cavity}, \text{Catch})$ has $2^3 - 1 = 7$ independent entries (must sum up to 1)

But: If there is a cavity, the probability for “catch” does not depend on whether there is a toothache:

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

The same holds when there is no 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*:

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

Writing out full joint distribution using chain rule:

$$\begin{aligned} &P(\text{Toothache}, \text{Catch}, \text{Cavity}) \\ &= P(\text{Toothache} \mid \text{Catch}, \text{Cavity}) P(\text{Catch}, \text{Cavity}) \\ &= P(\text{Toothache} \mid \text{Catch}, \text{Cavity}) P(\text{Catch} \mid \text{Cavity}) P(\text{Cavity}) \\ &= P(\text{Toothache} \mid \text{Cavity}) P(\text{Catch} \mid \text{Cavity}) P(\text{Cavity}) \end{aligned}$$

gives thus $2 + 2 + 1 = 5$ independent entries

Conditional independence (2)

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 .

Hence:

Conditional independence is our most basic and robust form of knowledge about uncertain environments

Bayes' Rule

Recap *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:

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

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

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

E.g., with M “meningitis”, S “stiff neck”:

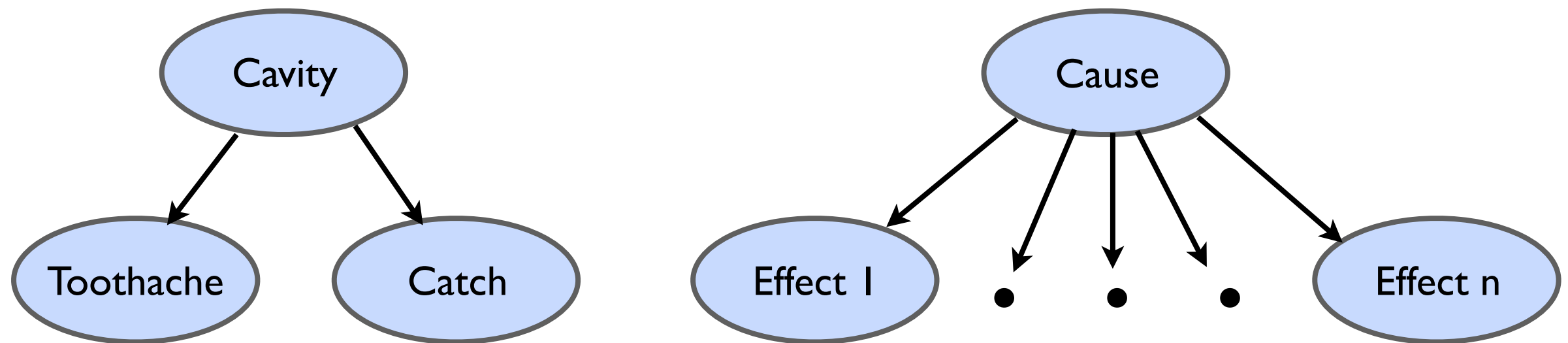
$$P(m | s) = \frac{P(s | m) P(m)}{P(s)} = \frac{0.8 * 0.0001}{0.1} = 0.0008 \quad (\text{not too bad, really!})$$

Bayes' Rule and conditional independence

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

An example of a *naive Bayes* model:

$$P(\text{Cause}, \text{Effect}_1, \dots, \text{Effect}_n) = P(\text{Cause}) \prod_i P(\text{Effect}_i \mid \text{Cause})$$



The total number of parameters is *linear* in n

Wumpus World

1,4	2,4	3,4	4,4
1,3	2,3	3,3	4,3
1,2 B ok	2,2	3,2	4,2
1,1 ok	2,1 B ok	3,1	4,1

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

$B_{ij} = \text{true}$ 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 $P(P_{1,1}, \dots, P_{4,4}, B_{1,1}, B_{1,2}, B_{2,1})$

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

(getting $P(\text{Effect} \mid \text{Cause})$.)

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

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

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

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 $P(P_{1,3} \mid known, b)$

Define: $Unknown = P_{i,j}$ s other than $P_{1,3}$ and $Known$

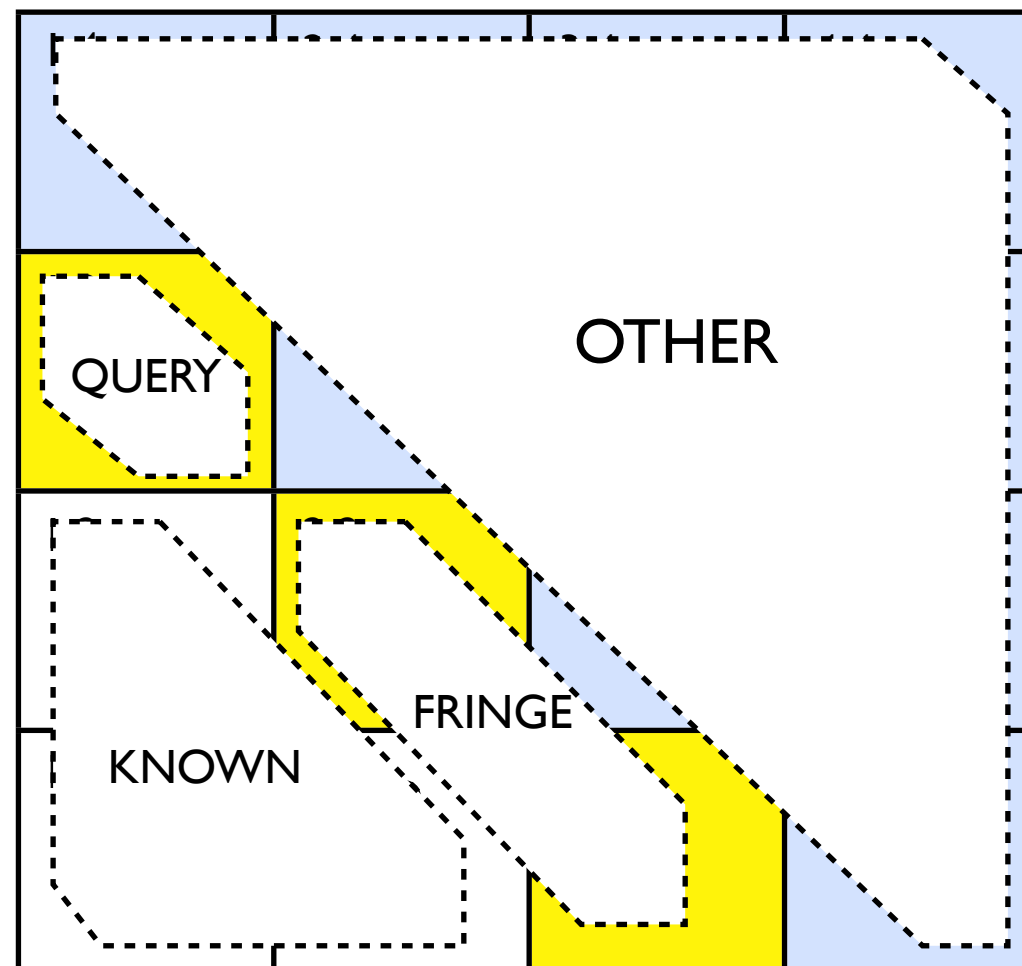
For inference by enumeration, we have

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

Grows exponentially with number of squares!

Using conditional independence

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



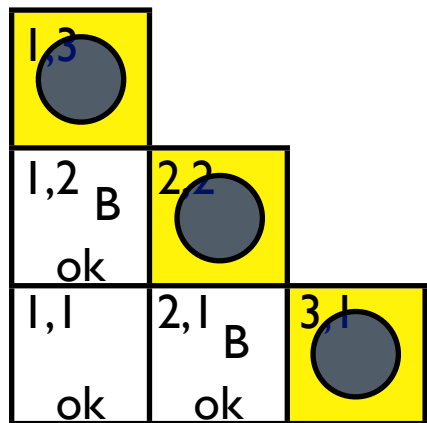
Define $Unknown = Fringe \cup Other$

$$\mathbf{P}(b \mid P_{I,3}, \text{Known}, \text{Unknown}) = \mathbf{P}(b \mid P_{I,3}, \text{Known}, \text{Fringe})$$

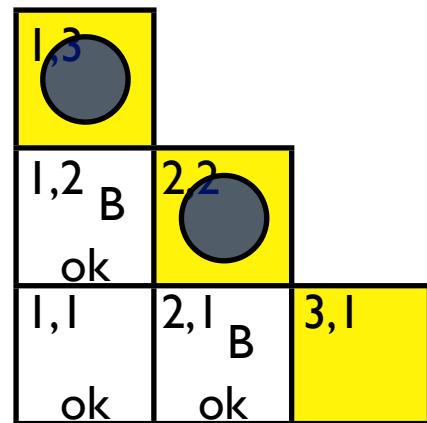
Using conditional independence (2)

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

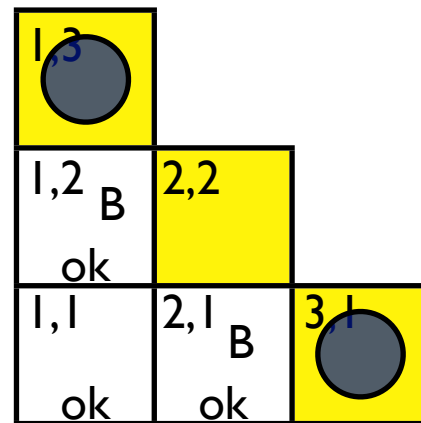
Wumpus World



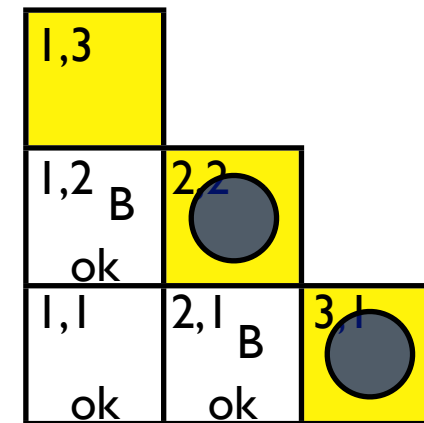
$$0.2 * 0.2 = 0.04$$



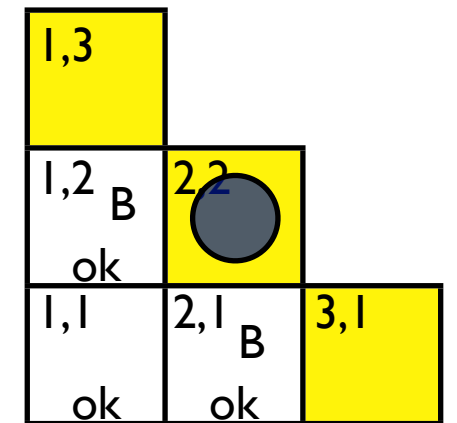
$$0.2 * 0.8 = 0.16$$



$$0.8 * 0.2 = 0.16$$



$$0.2 * 0.2 = 0.04$$



$$0.2 * 0.8 = 0.16$$

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

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

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

Summary

Probability is a way to formalise and represent uncertain knowledge

The *joint probability distribution* specifies probability over every *atomic event*

Queries can be answered by *summing* over atomic events

For *nontrivial* domains, we must find a way to *reduce* the joint size

Independence and *conditional independence* provide the tools

Bayes' rule can be applied to compute posterior probabilities so that *diagnostic* probabilities can be assessed from *causal* ones

Outline

- Uncertainty (chapter 13)
 - Uncertainty
 - Probability
 - Syntax and Semantics
 - Inference
 - Independence and Bayes' Rule
- Bayesian Networks (chapter 14.1-3)
 - Syntax
 - Semantics
 - Efficient representation of conditional distributions (parameterised distributions)

Outline

- Uncertainty (chapter 13)
 - Uncertainty
 - Probability
 - Syntax and Semantics
 - Inference
 - Independence and Bayes' Rule
- Bayesian Networks (chapter 14.1-3)
 - Syntax
 - Semantics
 - Efficient representation of conditional distributions (parameterised distributions)

Bayesian networks

A simple, graphical notation for conditional independence assertions and hence for compact specification of full joint distributions

Syntax:

- a set of nodes, one per random variable

- a directed, acyclic graph (link \approx “directly influences”)

- a conditional distribution for each node given its parents:

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

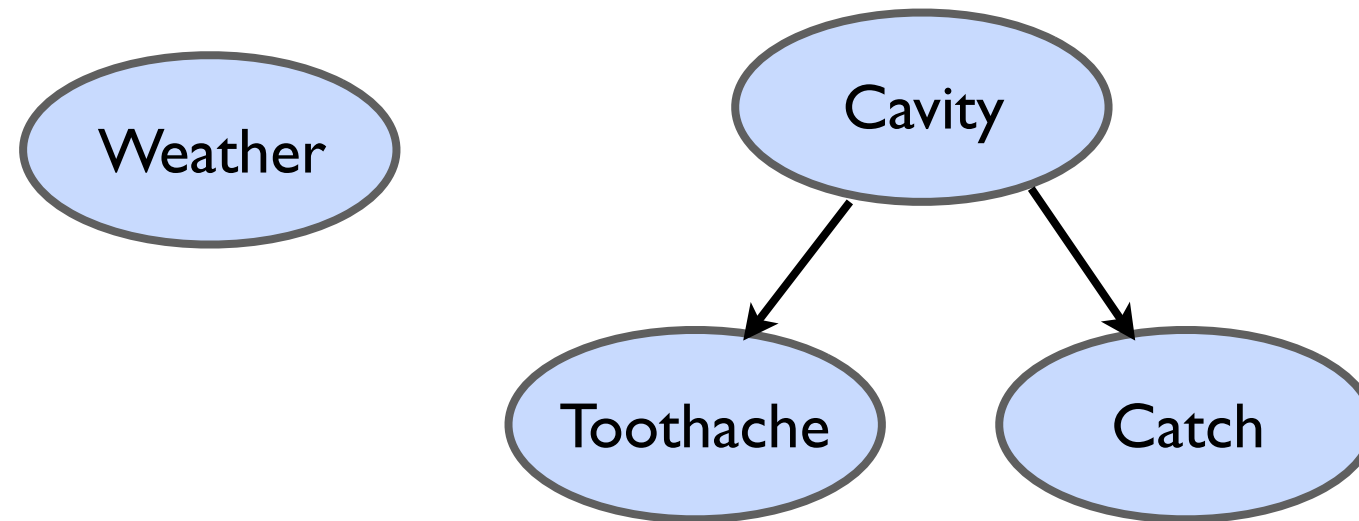
In the simplest case, conditional distribution represented as a

conditional probability table (CPT)

giving the distribution over X_i for each combination of parent values

Example

Topology of network encodes conditional independence assertions:



Weather is independent of the other variables

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

Example 2

I am at work, my neighbour John calls to say my alarm is ringing, but neighbour Mary does not call.

Sometimes the alarm is set off by minor earthquakes.

Is there a burglar?

Variables: *Burglar*, *Earthquake*, *Alarm*, *JohnCalls*, *MaryCalls*

Network topology reflects “causal” knowledge:

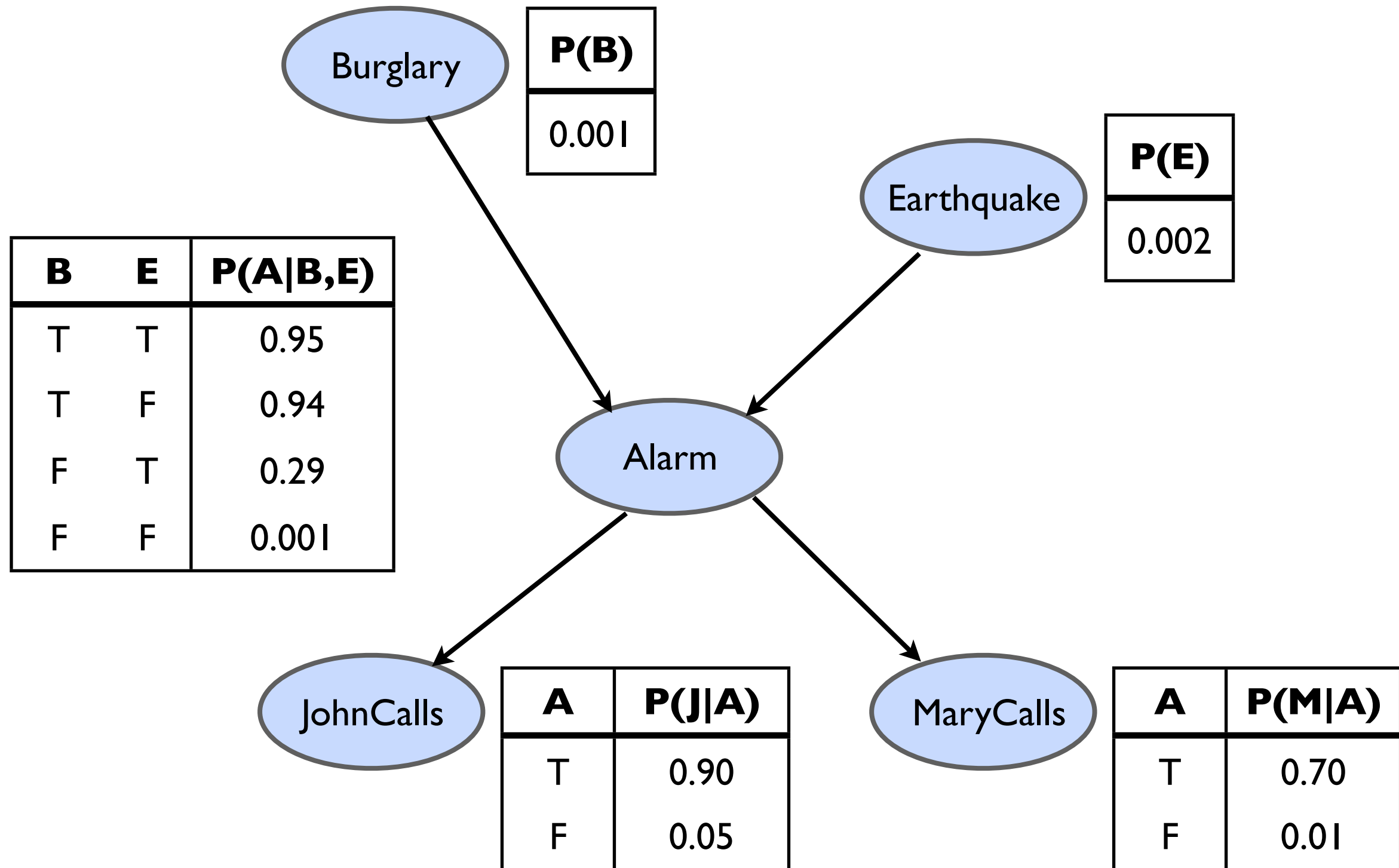
A burglar can set the alarm off

An earthquake can set the alarm off

The alarm can cause John to call

The alarm can cause Mary to call

Example 2 (2)



Example 2

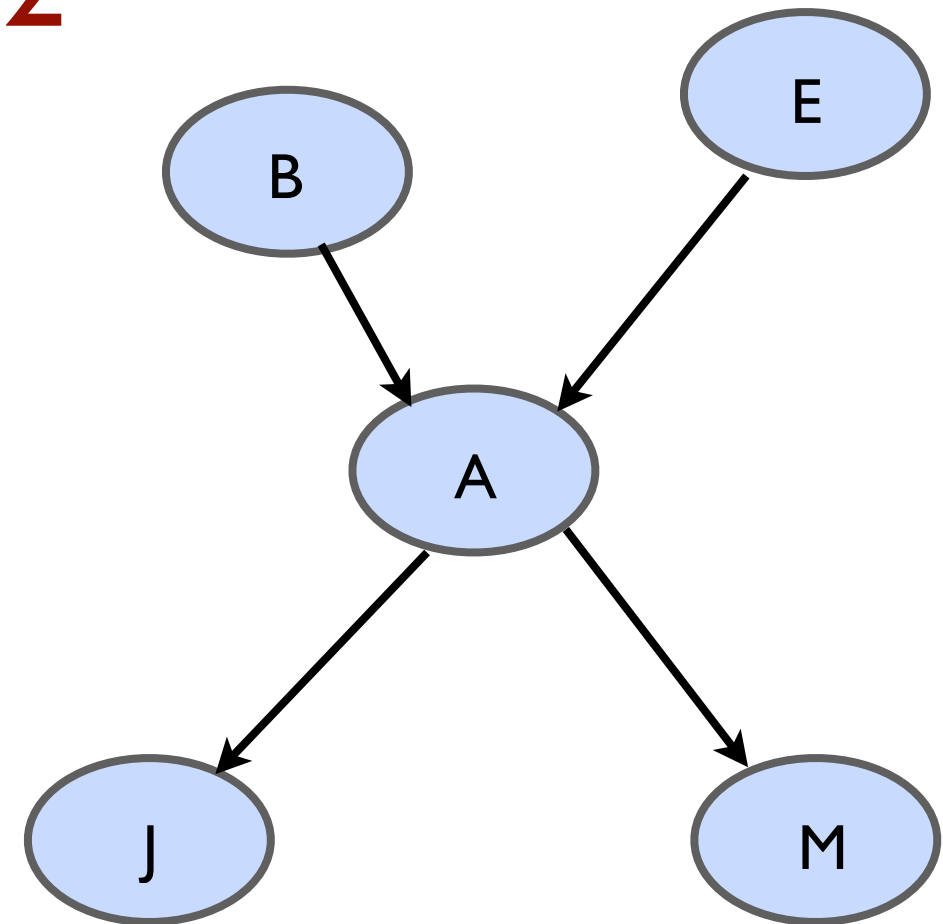
A CPT for Boolean X_i with k Boolean parents has 2^k rows for the combinations of parent values

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

If each variable has no more than k parents, the complete network requires $O(n 2^k)$ 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. $2^5 - 1 = 31$)



Outline

- Uncertainty (chapter 13)
 - Uncertainty
 - Probability
 - Syntax and Semantics
 - Inference
 - Independence and Bayes' Rule
- Bayesian Networks (chapter 14.1-3)
 - Syntax
 - Semantics
 - Efficient representation of conditional distributions (parameterised distributions)

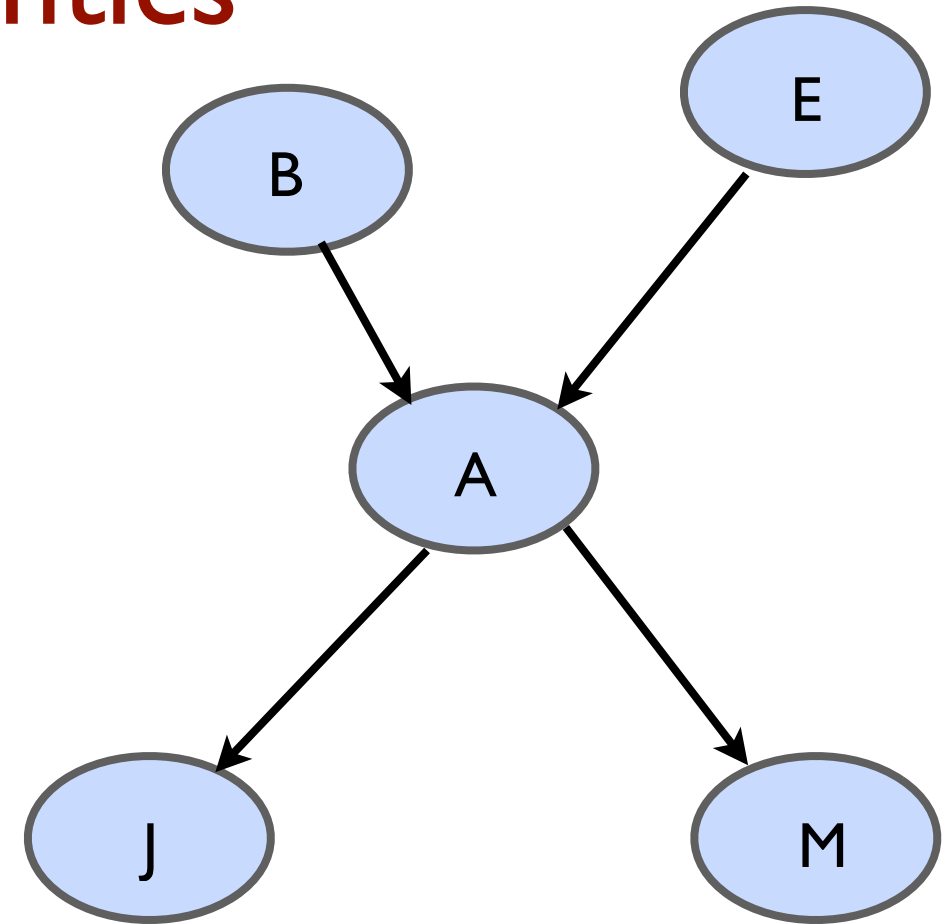
Global semantics

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

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

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

=



Global semantics

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

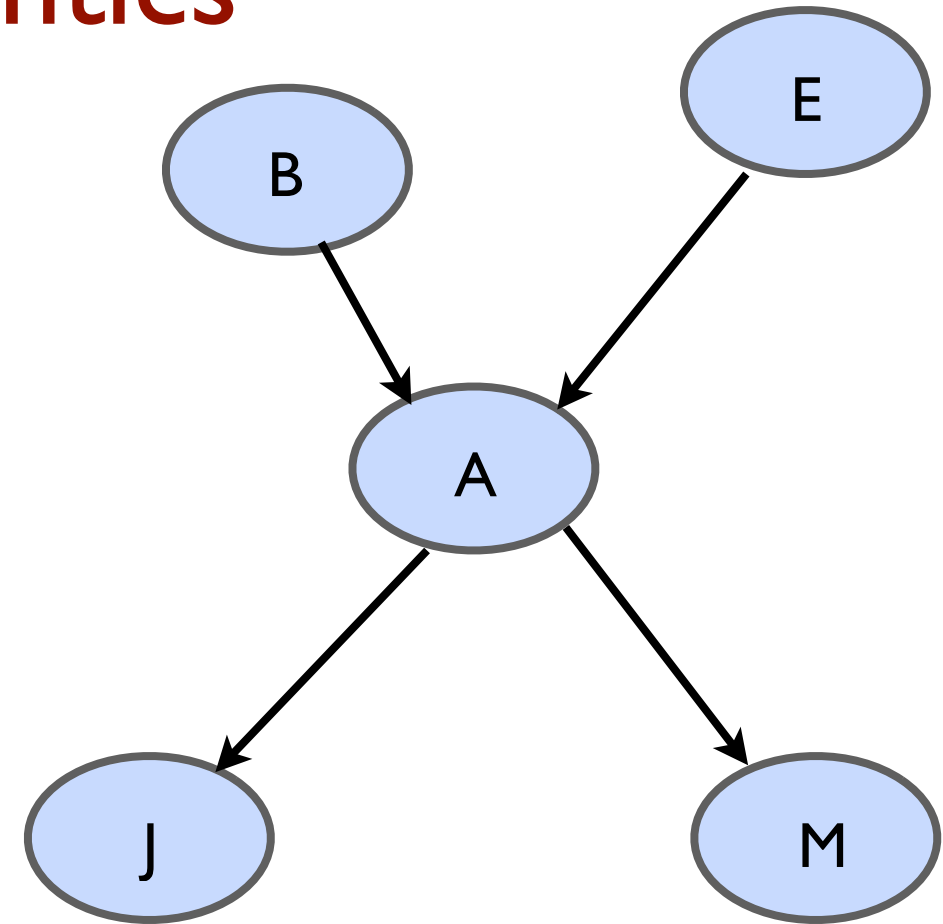
$$P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i \mid \text{parents}(X_i))$$

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

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

$$= 0.9 * 0.7 * 0.001 * 0.999 * 0.998$$

$$\approx 0.000628$$



Constructing Bayesian networks

We need a method such that 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

$$P(X_i \mid \text{Parents}(X_i)) = P(X_i \mid X_1, \dots, X_{i-1})$$

This choice of parents guarantees the global semantics:

$$P(X_1, \dots, X_n) = \prod_{i=1}^n P(X_i \mid X_1, \dots, X_{i-1}) \quad (\text{chain rule})$$

$$= \prod_{i=1}^n P(X_i \mid \text{Parents}(X_i)) \quad (\text{by construction})$$

Construction example

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

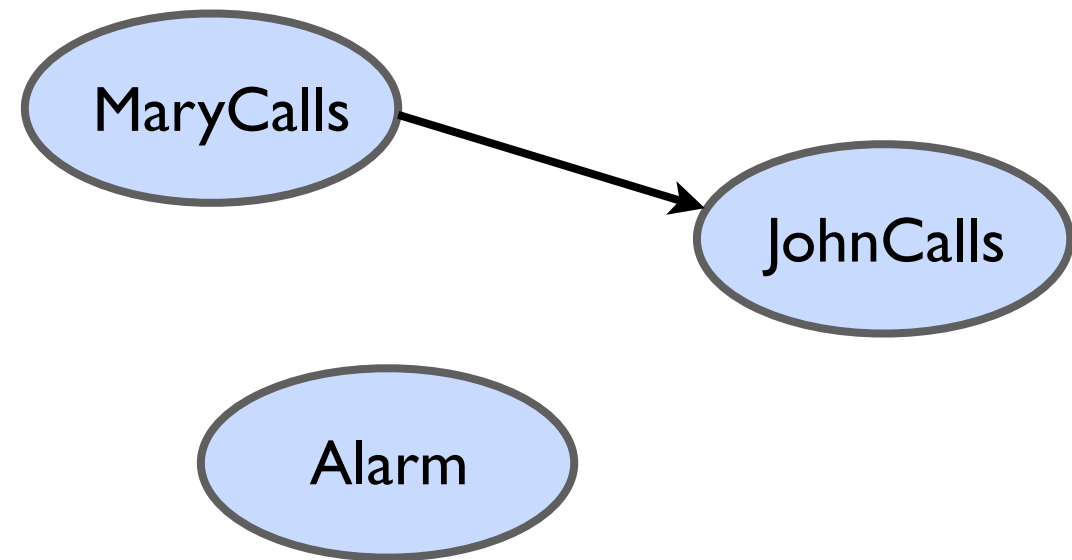
MaryCalls

JohnCalls

$$P(J \mid M) = P(J) ?$$

Construction example

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

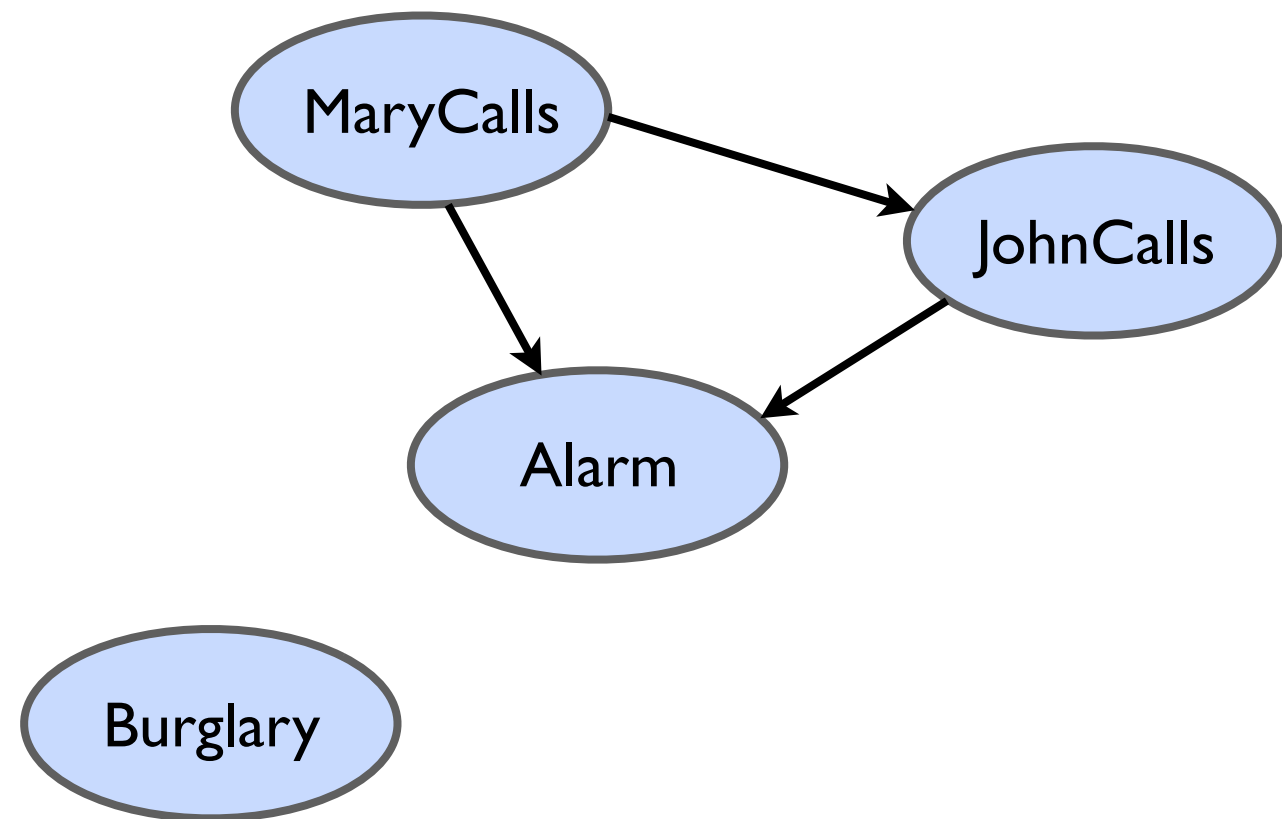


$P(J | M) = P(J)$? No

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

Construction example

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



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

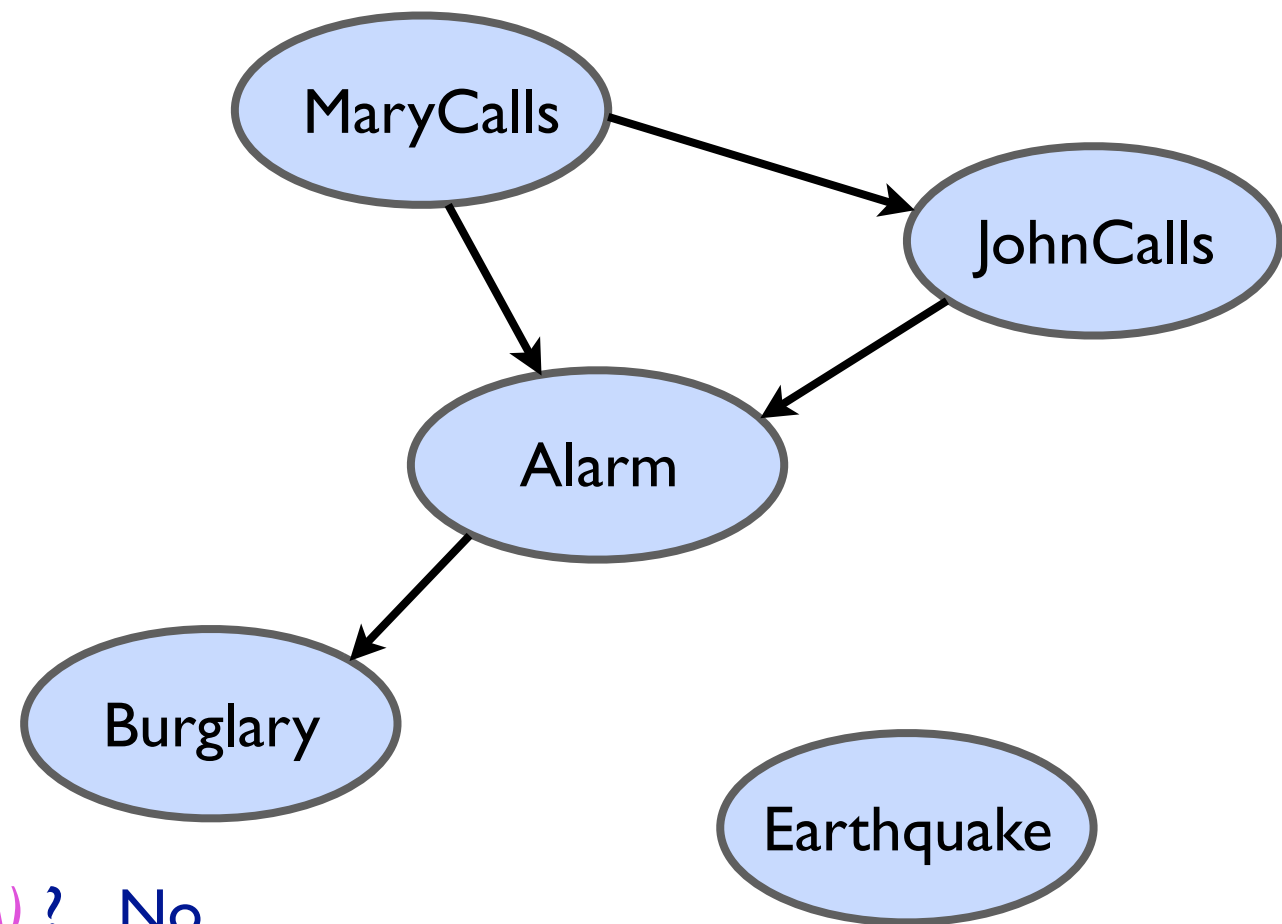
$P(A | J, M) = P(A | J)$? $P(A | J, M) = P(A)$? No

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

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

Construction example

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



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

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

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

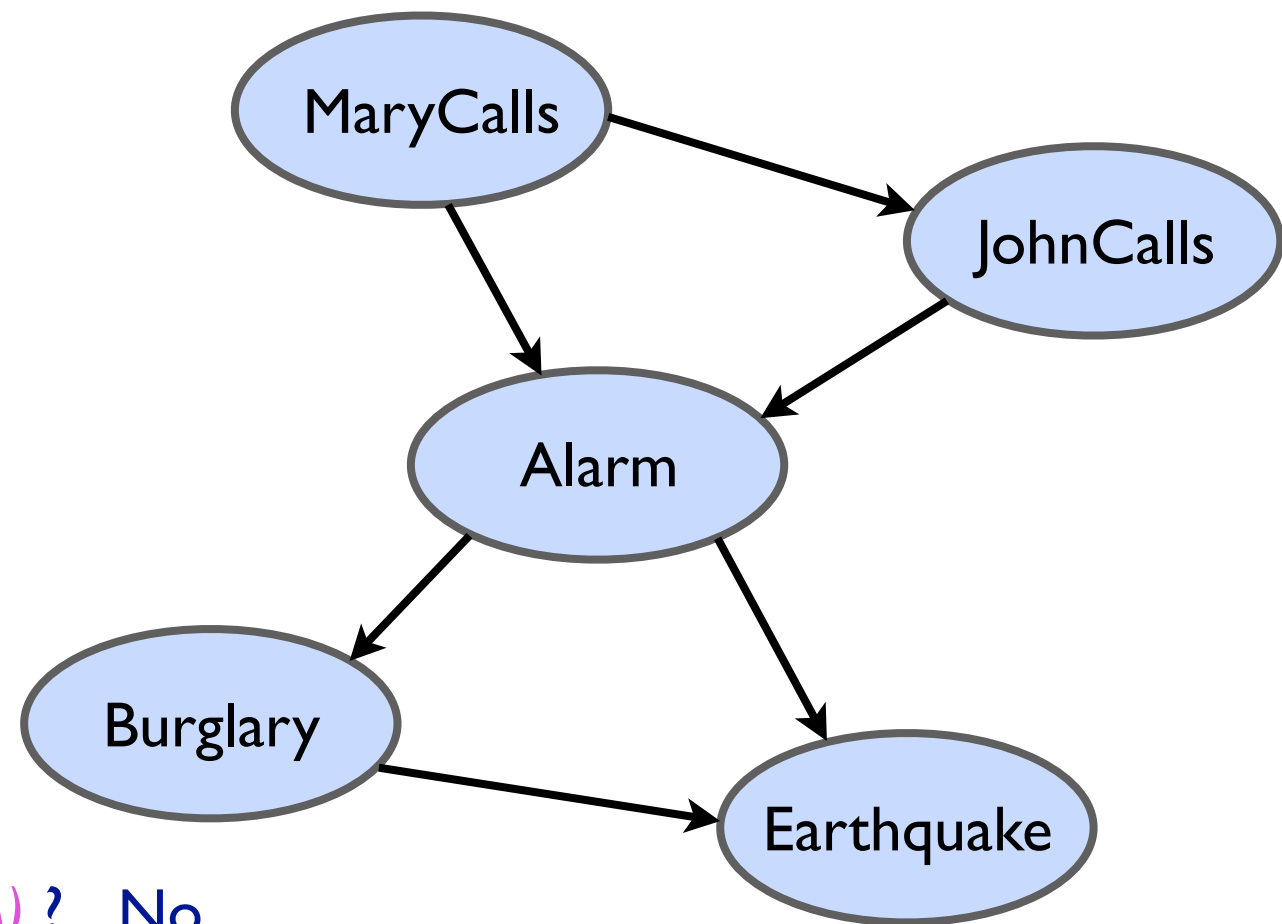
$P(B | A, J, M) = P(B)$? No

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

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

Construction example

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



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

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

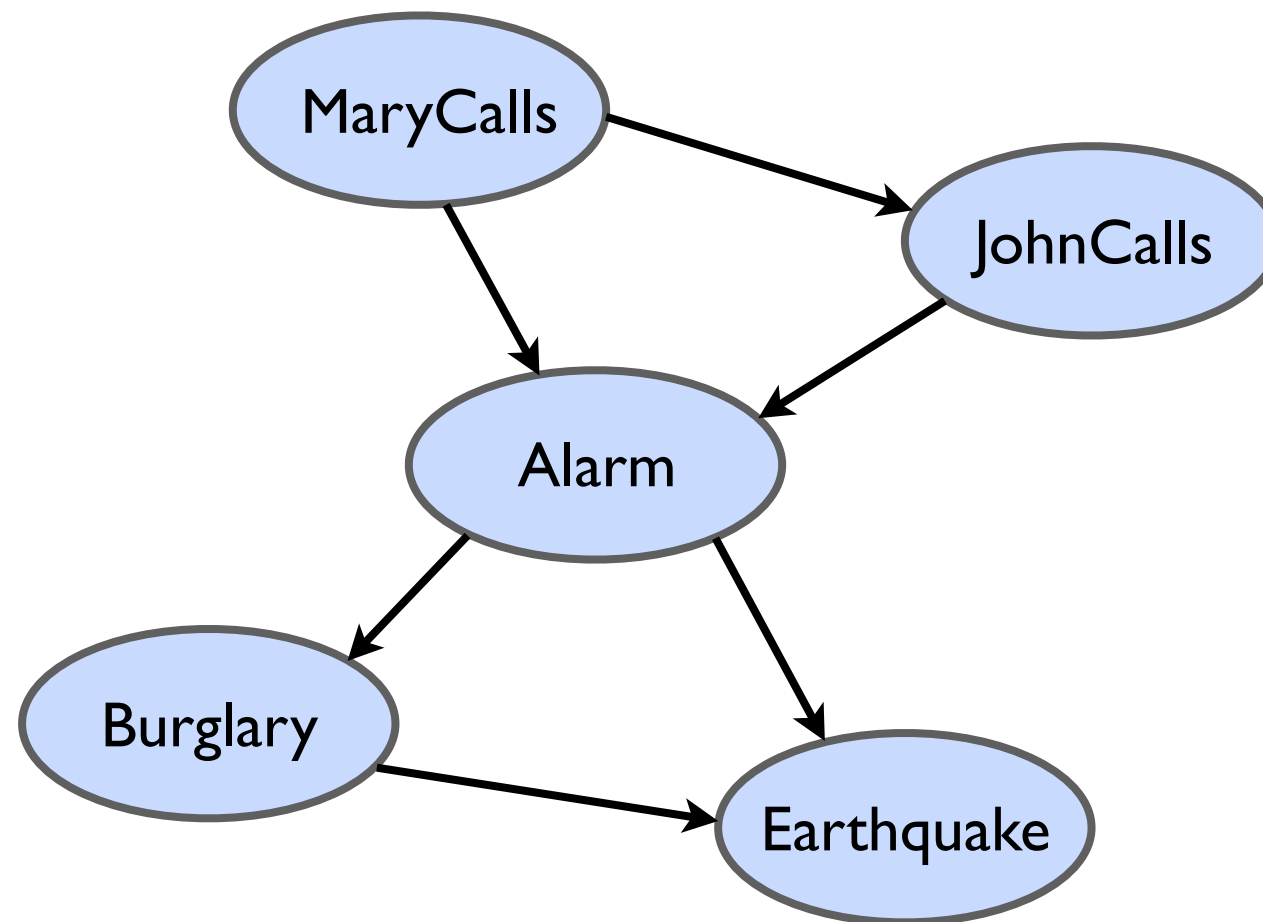
$P(B | A, J, M) = P(B | A)$? Yes

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

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

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

Construction example



Deciding conditional independence is hard in noncausal directions

(Causal models and conditional independence seem hardwired for humans!)

Assessing conditional probabilities is hard in noncausal directions

Network is less compact: $1 + 2 + 4 + 2 + 4 = 13$ numbers

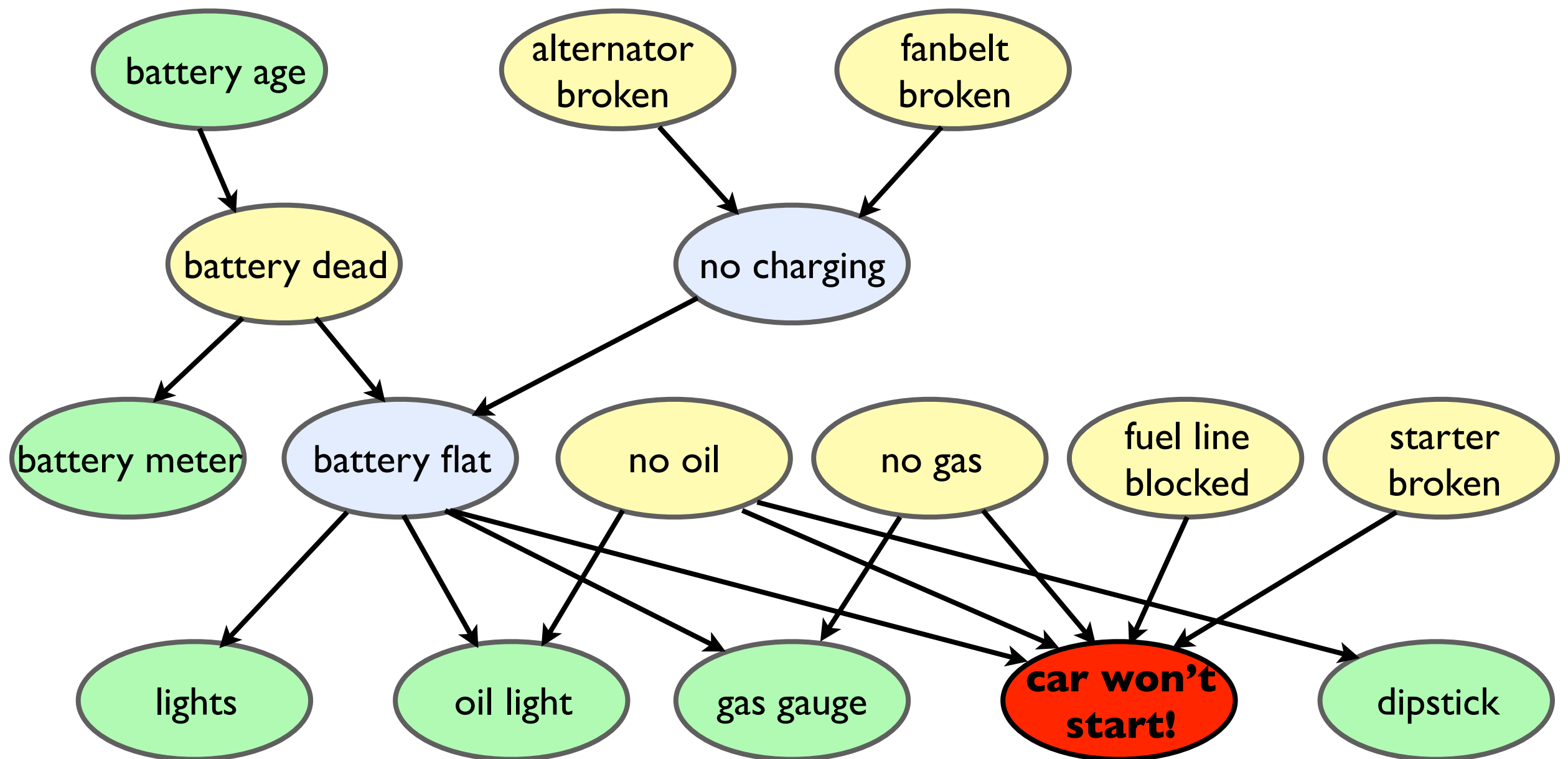
Hence: Choose preferably an order corresponding to the cause \rightarrow effect “chain”

Locally structured (sparse): Car diagnosis

Initial evidence: The *** car won't start!

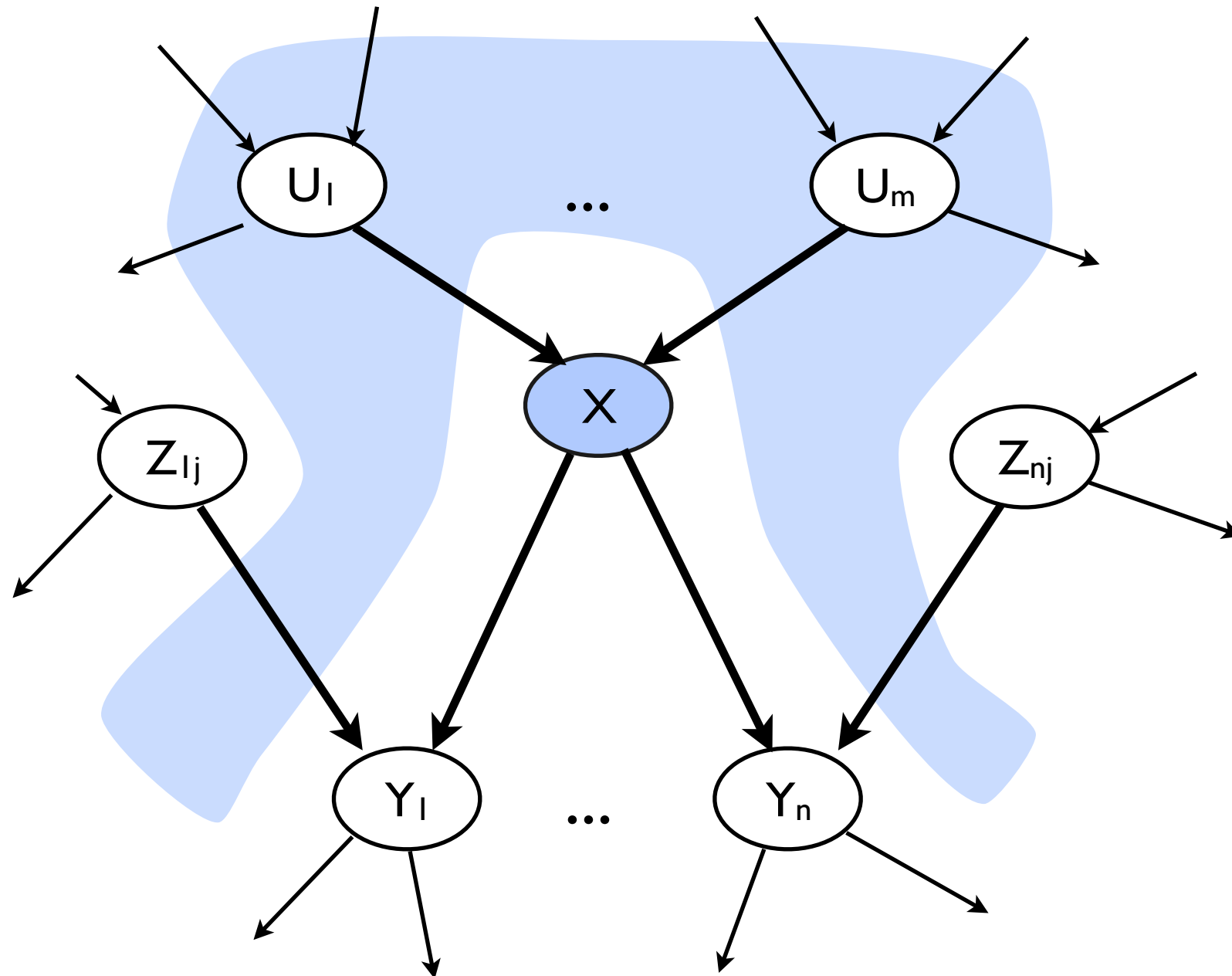
Testable variables (green), “broken, so fix it” variables (yellow)

Hidden variables (blue) ensure sparse structure / reduce parameters



Local semantics

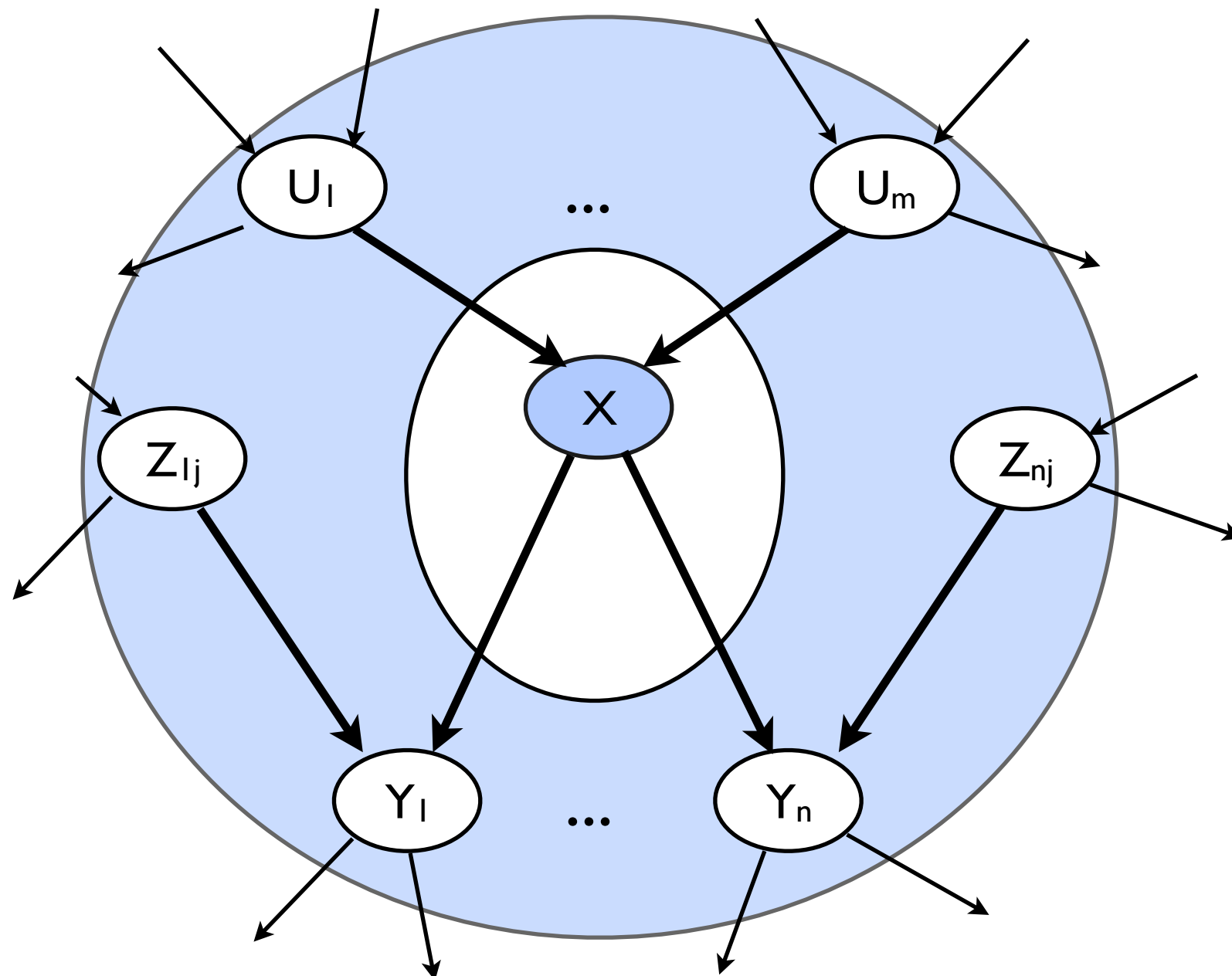
Local semantics: each node is conditionally independent of its non-descendants given its parents



Markov blanket

Each node is conditionally independent of all others given its

Markov blanket: parents + children + children's parents



Outline

- Uncertainty (chapter 13)
 - Uncertainty
 - Probability
 - Syntax and Semantics
 - Inference
 - Independence and Bayes' Rule
- Bayesian Networks (chapter 14.1-3)
 - Syntax
 - Semantics
 - Efficient representation of conditional distributions (parameterised distributions)

Compact conditional distributions

CPT grows exponentially with numbers of parents (i.e., causes to the effect)

CPT becomes infinite with continuous-valued parent or child

Solution: *canonical* distributions that are defined compactly

Deterministic nodes are the simplest case:

$$X = f(\text{Parents}(X)) \text{ for some function } f$$

E.g., Boolean functions

$$\text{NorthAmerican} \Leftrightarrow \text{Canadian} \vee \text{US} \vee \text{Mexican}$$

E.g., numerical relationships among continuous variables

$$\frac{\delta \text{Level}}{\delta t} = \text{inflow} + \text{precipitation} - \text{outflow} - \text{evaporation}$$

Compact conditional distributions (2)

Noisy-OR distributions model multiple noninteracting causes

1) Parents $U_1 \dots U_k$ include all causes (add *leak node* for “miscellaneous” ones)

2) Independent failure probability q_i for each cause alone

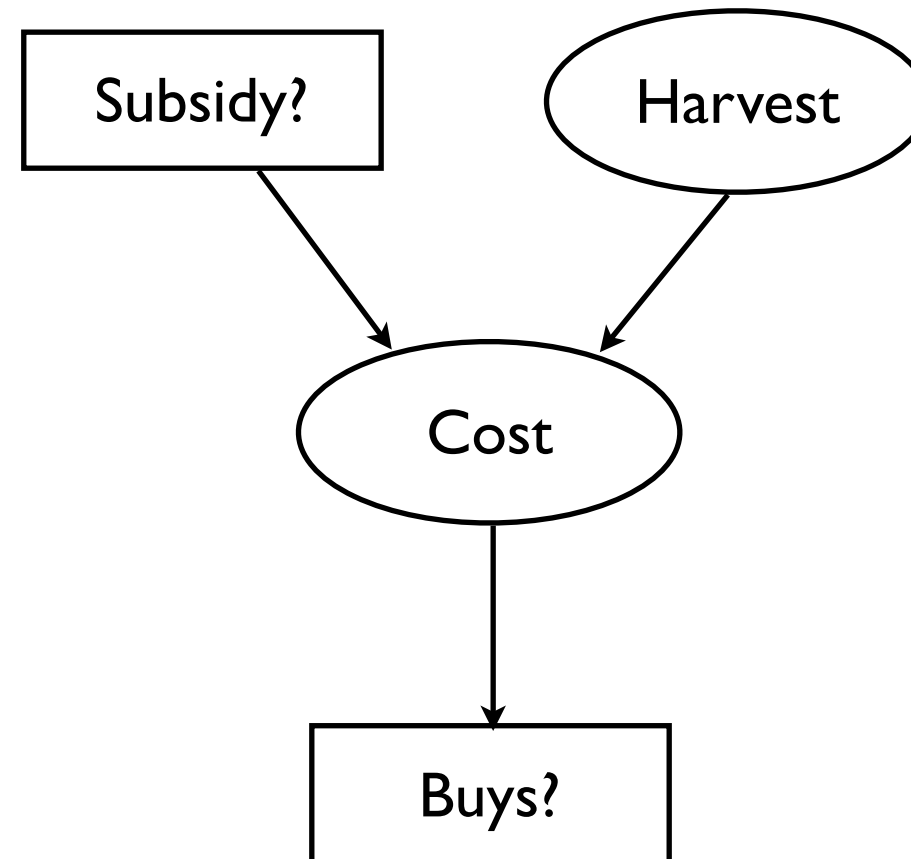
$$\Rightarrow P(X \mid U_1, \dots, U_j, \neg U_{j+1}, \dots, \neg U_k) = 1 - \prod_{i=1}^j q_i$$

<i>Cold</i>	<i>Flu</i>	<i>Malaria</i>	<i>P(Fever)</i>	<i>P(¬Fever)</i>
F	F	F	0.0	<i>1.0</i>
F	F	T	<i>0.9</i>	0.1
F	T	F	<i>0.8</i>	0.2
F	T	T	<i>0.98</i>	<i>0.02 = 0.2 * 0.1</i>
T	F	F	<i>0.4</i>	0.6
T	F	T	<i>0.94</i>	<i>0.06 = 0.6 * 0.1</i>
T	T	F	<i>0.88</i>	<i>0.12 = 0.6 * 0.2</i>
T	T	T	<i>0.988</i>	<i>0.012 = 0.6 * 0.2 * 0.1</i>

Number of parameters **linear** in number of parents

Hybrid (discrete + continuous) networks

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



Option 1: discretisation - possibly large errors, large CPTs

Option 2: finitely parameterised canonical families

1) Continuous variable, discrete + continuous parents (e.g., *Cost*)

2) Discrete variable, continuous parents (e.g., *Buys?*)

Continuous child variables

Need one *conditional density* function for child variable given continuous parents, for each possible assignment to discrete parents

Most common is the *linear Gaussian* model (following *G. normal distribution*), e.g.:

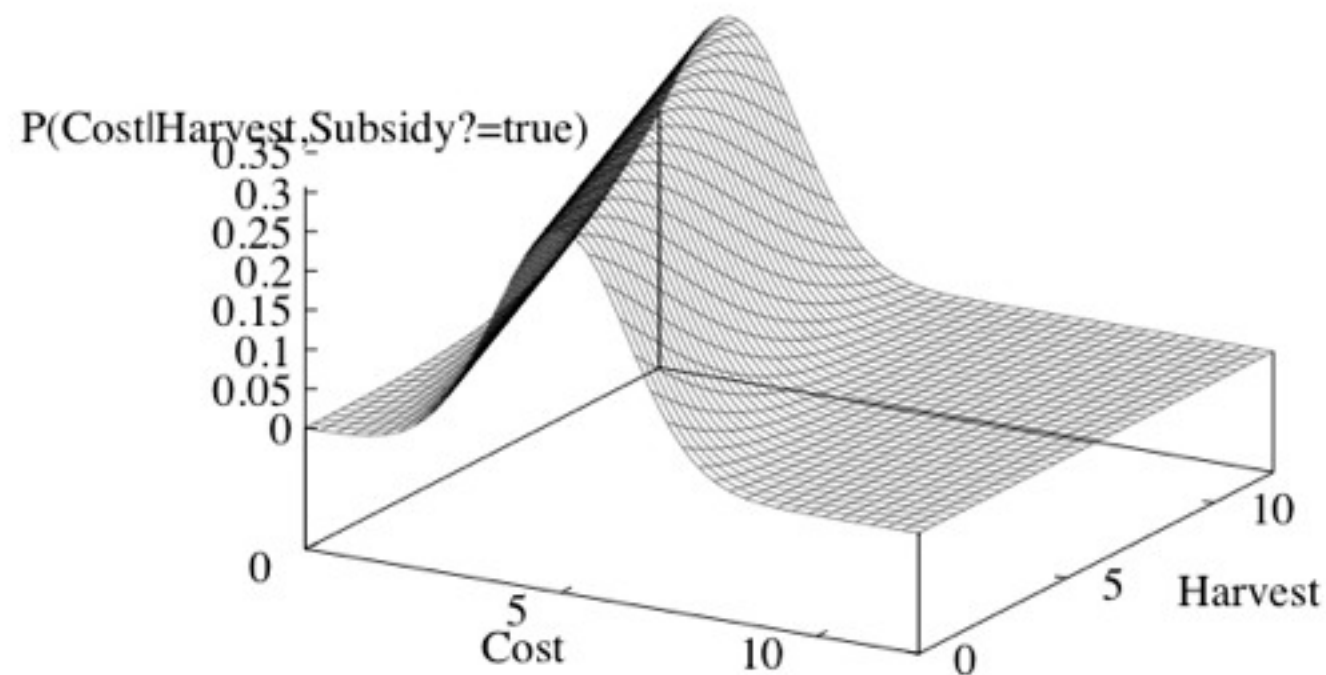
$$P(\text{Cost} = c \mid \text{Harvest} = h, \text{Subsidy?} = \text{true})$$

$$= N(a_th + b_t, \sigma_t)(c)$$

Mean *Cost* varies linearly with *Harvest*, variance is fixed

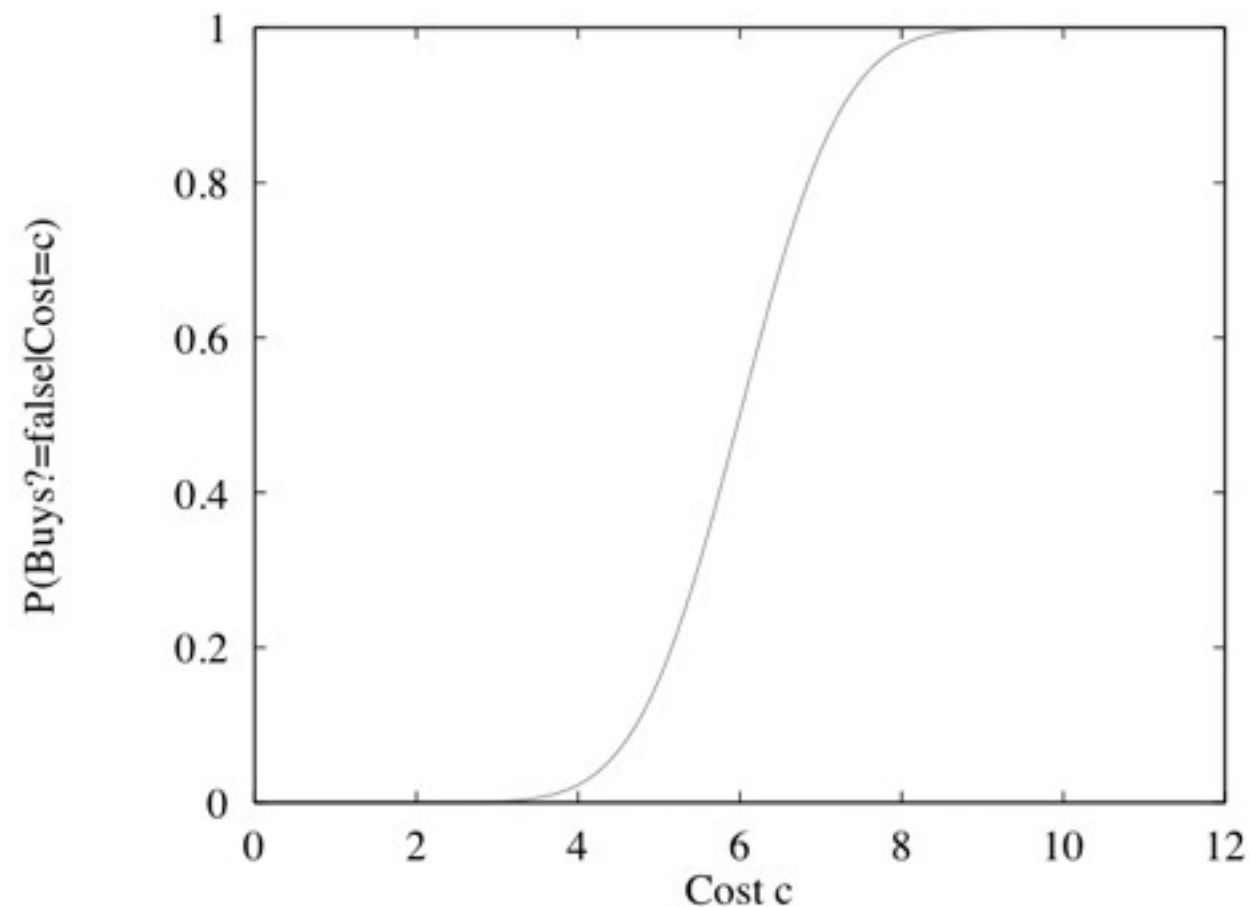
Linear variation is unreasonable over the full range

but works OK if the *likely* range of *Harvest* is narrow



Discrete variable with continuous parents

Probability of *Buys?* given *Cost* should be a soft threshold



E.g., Probit (probability unit), the integral of the standard normal distribution, or “logit”, the logistic function, also representing a type of sigmoid.

Summary

Bayesian networks provide a natural representation for (causally induced) conditional independence

Topology + CPTs = compact representation of joint distribution

Generally easy for (non)experts to construct

Canonical distributions (e.g., noisy-OR) = compact representation of CPTs

Continuous variables \Rightarrow parameterised distributions (e.g., linear Gaussians)