

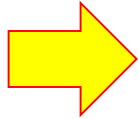
Lecture 9b: Bayesian Networks

CSCI 360

Introduction to Artificial Intelligence

USC

Here is where we are...



	3/1		Project 2 Out	
9	3/4	3/5	Quantifying Uncertainty	[Ch 13.1-13.6]
	3/6	3/7	Bayesian Networks	[Ch 14.1-14.2]
10	3/11	3/12	(spring break, no class)	
	3/13	3/14	(spring break, no class)	
11	3/18	3/19	Inference in Bayesian Networks	[Ch 14.3-14.4]
	3/20	3/21	Decision Theory	[Ch 16.1-16.3 and 16.5]
	3/23		Project 2 Due	
12	3/25	3/26	Advanced topics (Chao traveling to NSF)	
	3/27	3/28	Advanced topics (Chao traveling to NSF)	
	3/29		Homework 2 Out	
13	4/1	4/2	Markov Decision Processes	[Ch 17.1-17.2]
	4/3	4/4	Decision Tree Learning	[Ch 18.1-18.3]
	4/5		Homework 2 Due	
	4/5		Project 3 Out	
14	4/8	4/9	Perceptron Learning	[Ch 18.7.1-18.7.2]
	4/10	4/11	Neural Network Learning	[Ch 18.7.3-18.7.4]
15	4/15	4/16	Statistical Learning	[Ch 20.2.1-20.2.2]
	4/17	4/18	Reinforcement Learning	[Ch 21.1-21.2]
16	4/22	4/23	Artificial Intelligence Ethics	
	4/24	4/25	Wrap-Up and Final Review	
	4/26		Project 3 Due	
	5/3	5/2	Final Exam (2pm-4pm)	

Outline

- What is AI?
- Problem-solving agent (search)
- Knowledge-based agent (logical reasoning)
- **Probabilistic reasoning**
 - Quantifying Uncertainty
 - **Bayesian Networks**
 - Inference in Bayesian Networks
 - Decision Theory
 - Markov Decision Processes
- Machine learning

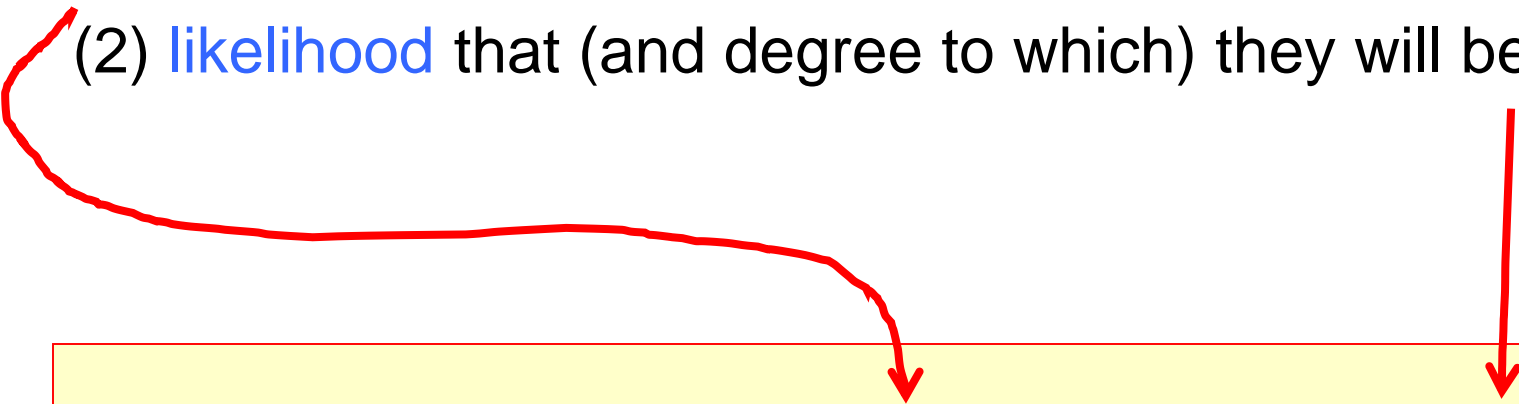
What we have learned so far...

- Early AI researchers largely rejected using probability in their systems
 - “People don’t think that way...”
- However, neither **problem-solving** nor **logical reasoning** agents tolerate approximation well...
 - Need probabilistic modeling/reasoning
 - Represent KB as relationships among random variables
 - KB can be learned from data
 - Given the KB, make inference about variables of interest
 - Example applications
 - *Medical diagnosis*: symptoms and diseases as random variables
 - *Business decision making*, e.g., predicting customer’s behavior
 - *Bio-informatics*, e.g., gene expression levels as random variables
 - *Computer science*: vision, speech recognition, spam filtering, ...

Recap: *Making decision*

Rational decision depends on

- (1) The **relative importance** of various goals and
- (2) **likelihood** that (and degree to which) they will be reached



A diagram illustrating the components of decision theory. Two red arrows point from the list above to a yellow box. The first arrow starts at '(1) The relative importance...' and points to 'Utility theory'. The second arrow starts at '(2) likelihood...' and points to 'Probability theory'. The yellow box contains the equation: **Decision theory** = Utility theory + Probability theory.

Decision theory = Utility theory + Probability theory

Choose the action that yields the **highest expected utility**, averaged over all the possible outcomes of the action

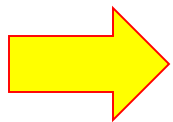
Recap: *Probability axioms*

- A numerical probability $P(\omega)$ for each possible world

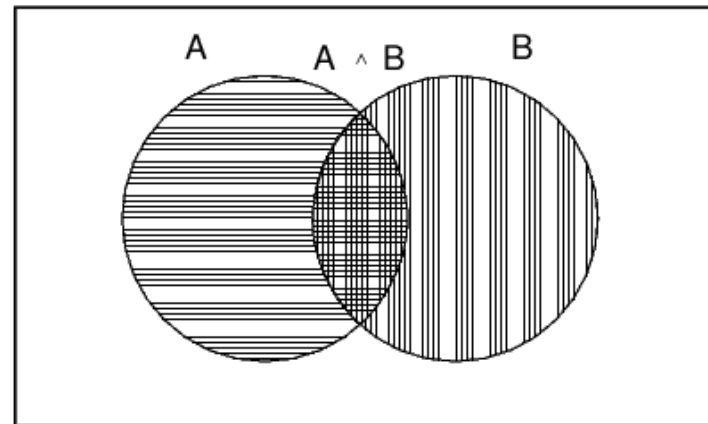
$$0 \leq P(\omega) \leq 1 \text{ for every } \omega$$

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

$$P(\neg a) = 1 - P(a)$$



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



Recap: Conditional (or posterior) probability

- For any propositions a and b , we have

$$P(a | b) = \frac{P(a \wedge b)}{P(b)} \quad \text{whenever } P(b) > 0.$$



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

joint probability

conditional probability

prior probability

Recap: *Probability distribution*

- Probabilities of all possible values of a random variable

$$P(\text{Weather} = \text{sunny}) = 0.6$$

$$P(\text{Weather} = \text{rain}) = 0.1$$

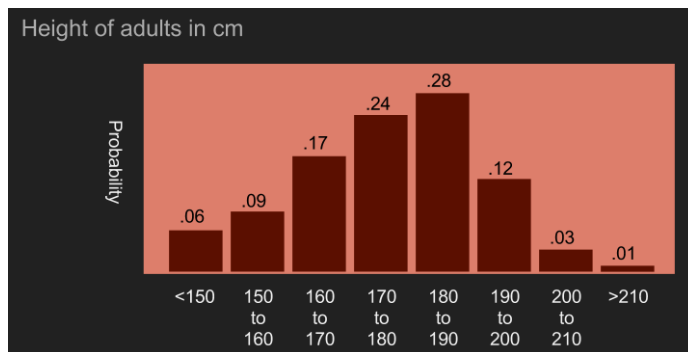
$$P(\text{Weather} = \text{cloudy}) = 0.29$$

$$P(\text{Weather} = \text{snow}) = 0.01 ,$$

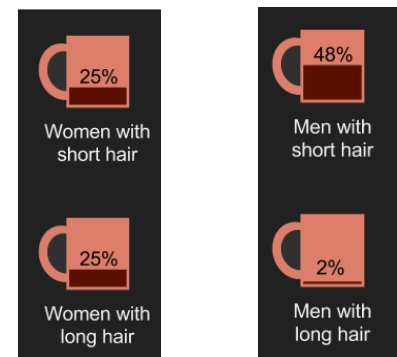
- In a vector format

$$\mathbf{P}(\text{Weather}) = \langle 0.6, 0.1, 0.29, 0.01 \rangle$$

- Other examples

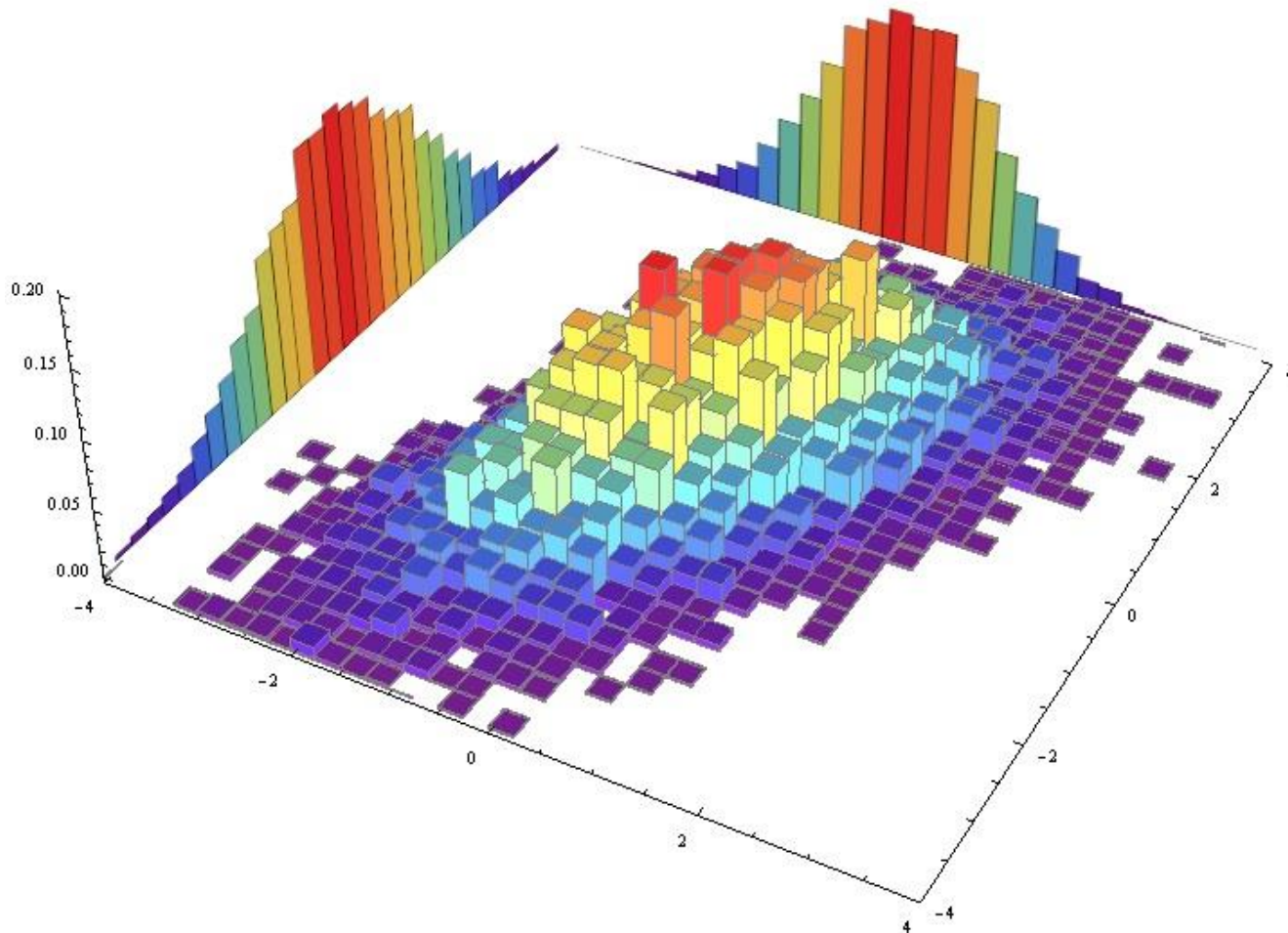


<https://brohrer.github.io>



Recap: *Joint probability distribution*

- Probabilities of all possible values of **multiple** variables



Recap: *Marginal probability*

- Extracting the distribution **over a subset** of variables from the full 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>	0.108	0.012	0.072	0.008
\neg <i>cavity</i>	0.016	0.064	0.144	0.576

- Example

$$P(cavity) =$$

Recap: *Marginal probability*

- Extracting the distribution **over a subset** of variables from the full 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>	0.108	0.012	0.072	0.008
\neg <i>cavity</i>	0.016	0.064	0.144	0.576

- Example

$$P(\text{cavity}) = 0.108 + 0.012 + 0.072 + 0.008 = 0.2$$

Recap: Normalization

- The probability of **cavity**, or **no cavity**, given **toothache**

	<i>toothache</i>		\neg <i>toothache</i>	
	<i>catch</i>	\neg <i>catch</i>	<i>catch</i>	\neg <i>catch</i>
<i>cavity</i>	0.108	0.012	0.072	0.008
\neg <i>cavity</i>	0.016	0.064	0.144	0.576

- Example

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

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

Sum of the two
is always 1.0

Recap: Normalization

- The probability of **cavity**, or **no cavity**, given **toothache**

	<i>toothache</i>		\neg <i>toothache</i>	
	<i>catch</i>	\neg <i>catch</i>	<i>catch</i>	\neg <i>catch</i>
<i>cavity</i>	0.108	0.012	0.072	0.008
\neg <i>cavity</i>	0.016	0.064	0.144	0.576

- Example

$$P(\text{cavity} \mid \text{toothache}) = \frac{P(\text{cavity} \wedge \text{toothache})}{\cancel{P(\text{toothache})}}$$

$= \frac{0.108 + 0.012}{\cancel{0.108 + 0.012 + 0.016 + 0.064}} = 0.6$

No need to compute $P(\text{toothache})$ any more

$$P(\neg \text{cavity} \mid \text{toothache}) = \frac{P(\neg \text{cavity} \wedge \text{toothache})}{\cancel{P(\text{toothache})}}$$

$= \frac{0.016 + 0.064}{\cancel{0.108 + 0.012 + 0.016 + 0.064}} = 0.4$

Recap: Normalization

- The probability of **cavity**, or **no cavity**, given **toothache**

	<i>toothache</i>		\neg <i>toothache</i>	
	<i>catch</i>	\neg <i>catch</i>	<i>catch</i>	\neg <i>catch</i>
<i>cavity</i>	0.108	0.012	0.072	0.008
\neg <i>cavity</i>	0.016	0.064	0.144	0.576

- Example

$$\mathbf{P}(Cavity \mid toothache) = \alpha \mathbf{P}(Cavity, toothache)$$

$$= \alpha \langle 0.12, 0.08 \rangle = \langle 0.6, 0.4 \rangle .$$

Assume that

$$\begin{aligned}\alpha &= 1 / (0.12 + 0.08) \\ &= 1 / 0.2 \\ &= 5\end{aligned}$$

Recap: *Independence to reduce table size*

- Consider $P(\textit{Toothache}, \textit{Catch}, \textit{Cavity}, \textit{Weather})$, which has 32 entries in the **full joint distribution** table

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

- Applying the product rule

$$\begin{aligned} P(\textit{toothache}, \textit{catch}, \textit{cavity}, \textit{cloudy}) \\ = P(\textit{cloudy} \mid \textit{toothache}, \textit{catch}, \textit{cavity}) P(\textit{toothache}, \textit{catch}, \textit{cavity}) \end{aligned}$$

- But weather is not influenced by dentistry!

$$P(\textit{cloudy} \mid \textit{toothache}, \textit{catch}, \textit{cavity}) = P(\textit{cloudy})$$

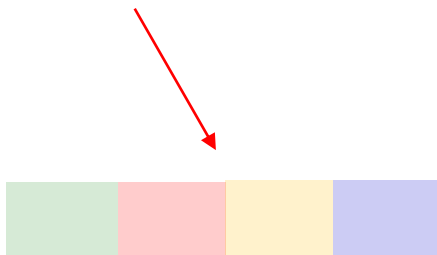
$$P(\textit{toothache}, \textit{catch}, \textit{cavity}, \textit{cloudy}) = P(\textit{cloudy}) P(\textit{toothache}, \textit{catch}, \textit{cavity})$$

Recap: *Independence to reduce table size*

- Consider $P(\textit{Toothache}, \textit{Catch}, \textit{Cavity}, \textit{Weather})$, which has 32 entries in the **full joint distribution** table

	<i>toothache</i>		<i>¬toothache</i>		<i>toothache</i>		<i>¬toothache</i>		<i>toothache</i>		<i>¬toothache</i>		<i>toothache</i>		<i>¬toothache</i>	
	<i>catch</i>	<i>¬catch</i>	<i>catch</i>	<i>¬catch</i>	<i>catch</i>	<i>¬catch</i>	<i>catch</i>	<i>¬catch</i>	<i>catch</i>	<i>¬catch</i>	<i>catch</i>	<i>¬catch</i>	<i>catch</i>	<i>¬catch</i>	<i>catch</i>	<i>¬catch</i>
<i>cavity</i>	0.108	0.012	0.072	0.008	0.108	0.012	0.072	0.008	0.108	0.012	0.072	0.008	0.108	0.012	0.072	0.008
<i>¬cavity</i>	0.016	0.064	0.144	0.576	0.016	0.064	0.144	0.576	0.016	0.064	0.144	0.576	0.016	0.064	0.144	0.576

- The 32-element table can be reduced to a 8-element table and a 4-element table

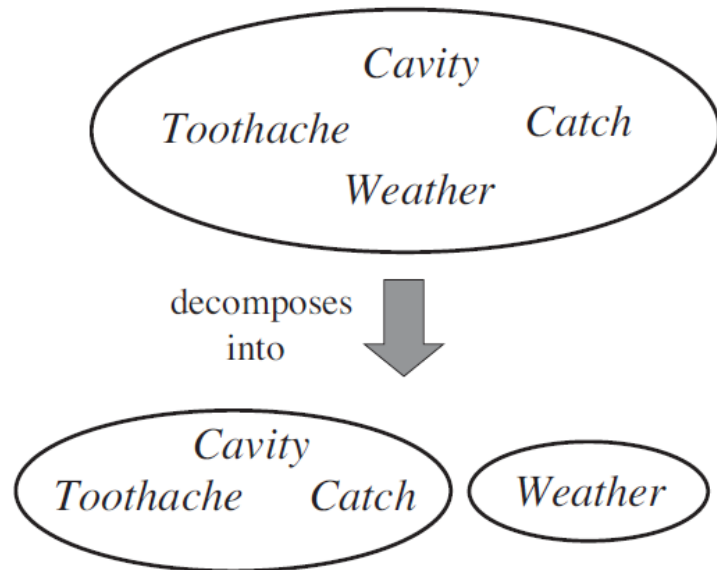


A diagram showing the 8-element table. A red arrow points from the second block (red) to the second block of the 8-element table below.

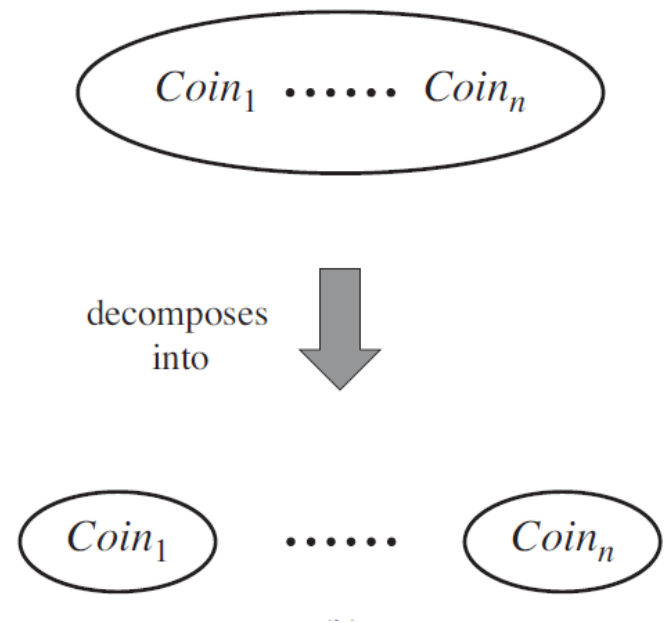
	<i>toothache</i>		<i>¬toothache</i>	
	<i>catch</i>	<i>¬catch</i>	<i>catch</i>	<i>¬catch</i>
<i>cavity</i>	0.108	0.012	0.072	0.008
<i>¬cavity</i>	0.016	0.064	0.144	0.576

Recap: *Independence to reduce table size*

Leveraging the (absolute) independence



Weather and dentistry are independent



Coin flips are independent

Recap: *Conditional independence*

- Variables X and Y are **conditional independent**, given a third variable Z

$$\mathbf{P}(X, Y \mid Z) = \mathbf{P}(X \mid Z)\mathbf{P}(Y \mid Z)$$

- Alternatively, we have

$$\mathbf{P}(X \mid Y, Z) = \mathbf{P}(X \mid Z)$$

For X , variable Y doesn't provide any additional information, given Z

$$\mathbf{P}(Y \mid X, Z) = \mathbf{P}(Y \mid Z)$$

For Y , variable X doesn't provide any additional information, given Z

Recap: *Conditional independence to reduce table size*

- For n effects that are conditionally independent given the cause, the **full joint distribution** table size grows as $O(n)$ instead of $O(2^n)$

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

Recap: *Bayes' rule*

- Derive **Bayes' rule** from the **product rule** of conditional probability

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

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

- Equating the right-hand sides and dividing by $P(a)$

$$P(b | a) = \frac{P(a | b)P(b)}{P(a)}$$

Recap: *Bayes' rule*

- Derive **Bayes' rule** from the **product rule** of conditional probability

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

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

- Equating the right-hand sides and dividing by $P(a)$

$$P(b | a) = \frac{P(a | b)P(b)}{P(a)}$$

This equation underlies most **modern AI systems** for **probabilistic inference**...



Recap: Quiz 4 solution

- Assume that the doctor knows some unconditional facts:
 - Prior probability that a patient has a disease $P(d) = 0.01$
 - Probability of test positive given no disease $P(tp | \neg d) = 0.096$
 - Probability of test positive given disease $P(tp | d) = 0.8$
- Now, a patient has test positive; what is the probability that this particular patient has the disease?

$$\begin{aligned} P(d|tp) &= P(tp/d) P(d) / P(tp) \\ &= P(tp/d) P(d) / (P(tp/d) P(d) + P(tp/\neg d) P(\neg d)) \\ &= 0.8*0.01 / (0.8*0.01 + 0.096*0.99) \\ &= 0.008 / (0.008+0.09504) \\ &= 0.0776 \end{aligned}$$

Outline

- Representing knowledge in Bayesian networks
- Semantics of Bayesian networks
- Efficient representation of conditional distributions
- Exact inference in Bayesian networks

Why Bayesian networks?

- **Full joint distribution** can be used to answer any query about the world
 - But table size is exponential in the number of variables
- **Independence** and **conditional independence** relations are important in simplifying the table
 - But they are unnatural and tedious to specify
- **Bayesian networks** is a data structure to represent both joint distributions and the dependencies among variables

Example

- Medical diagnosis
 - $S1, S2, \dots$: symptoms (e.g. high temperature) or causes of diseases (e.g. age)
 - $D1, D2, \dots$: diseases (e.g. flu, kidney stone, ...)

S1	S2	S3	...	D1	D2	D3	...	$P(S1, S2, S3, \dots, D1, D2, D3, \dots)$
true	true	true	...	true	true	true	...	0.0000001
...	
false	false	false	...	false	false	false	...	0.0000002

Example (cont'd)

- Medical diagnosis
 - $S1, S2, \dots$: symptoms (e.g. high temperature) or causes of diseases (e.g. age)
 - $D1, D2, \dots$: diseases (e.g. flu, kidney stone, ...)

S1	S2	S3	...	D1	D2	D3	...	$P(S1, S2, S3, \dots, D1, D2, D3, \dots)$
true	true	true	...	true	true	true	...	0.0000001
...	
false	false	false	...	false	false	false	...	0.0000002

- When the doctor observes **presence of $S1$** and **absence of $S3$** , calculate
 - $P(D1 \mid S1, \neg S3) = P(D1, S1, \neg S3) / P(S1, \neg S3) = \dots$
 - $P(D2 \mid S1, \neg S3) = \dots$
 - $P(D3 \mid S1, \neg S3) = \dots$
 - ...

Example (cont'd)

- Medical diagnosis
 - $S1, S2, \dots$: symptoms (e.g. high temperature) or causes of diseases (e.g. age)
 - $D1, D2, \dots$: diseases (e.g. flu, kidney stone, ...)

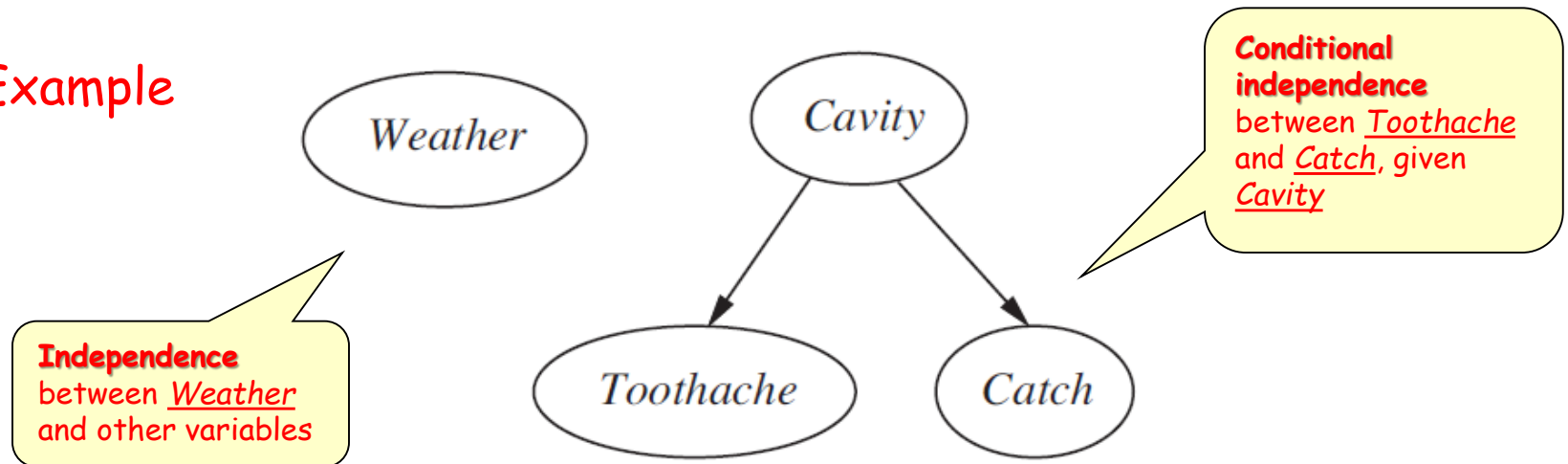
S1	S2	S3	...	D1	D2	D3	...	$P(S1, S2, S3, \dots, D1, D2, D3, \dots)$
true	true	true	...	true	true	true	...	0.0000001
...	
false	false	false	...	false	false	false	...	0.0000002

- We need to acquire **too many** probabilities from the expert.
- Many of the probabilities are **very close to zero** and thus **hard to specify** by experts.

Bayesian network

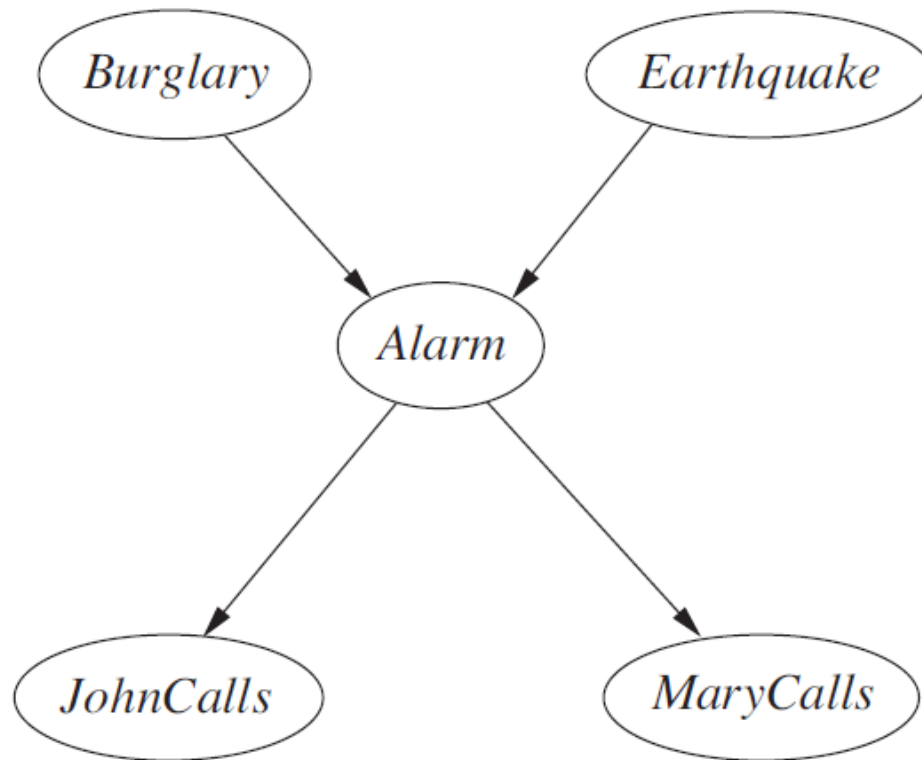
- A directed acyclic graph (DAG) where
 - Each node corresponds to a random variable,
 - Each edge from node X to node Y represents a direct influence of X on Y ,
 - Each node X_i has a conditional probability distribution $P(X_i | \text{Parents}(X_i))$ that quantifies the effect of the parents on the node.

- Example



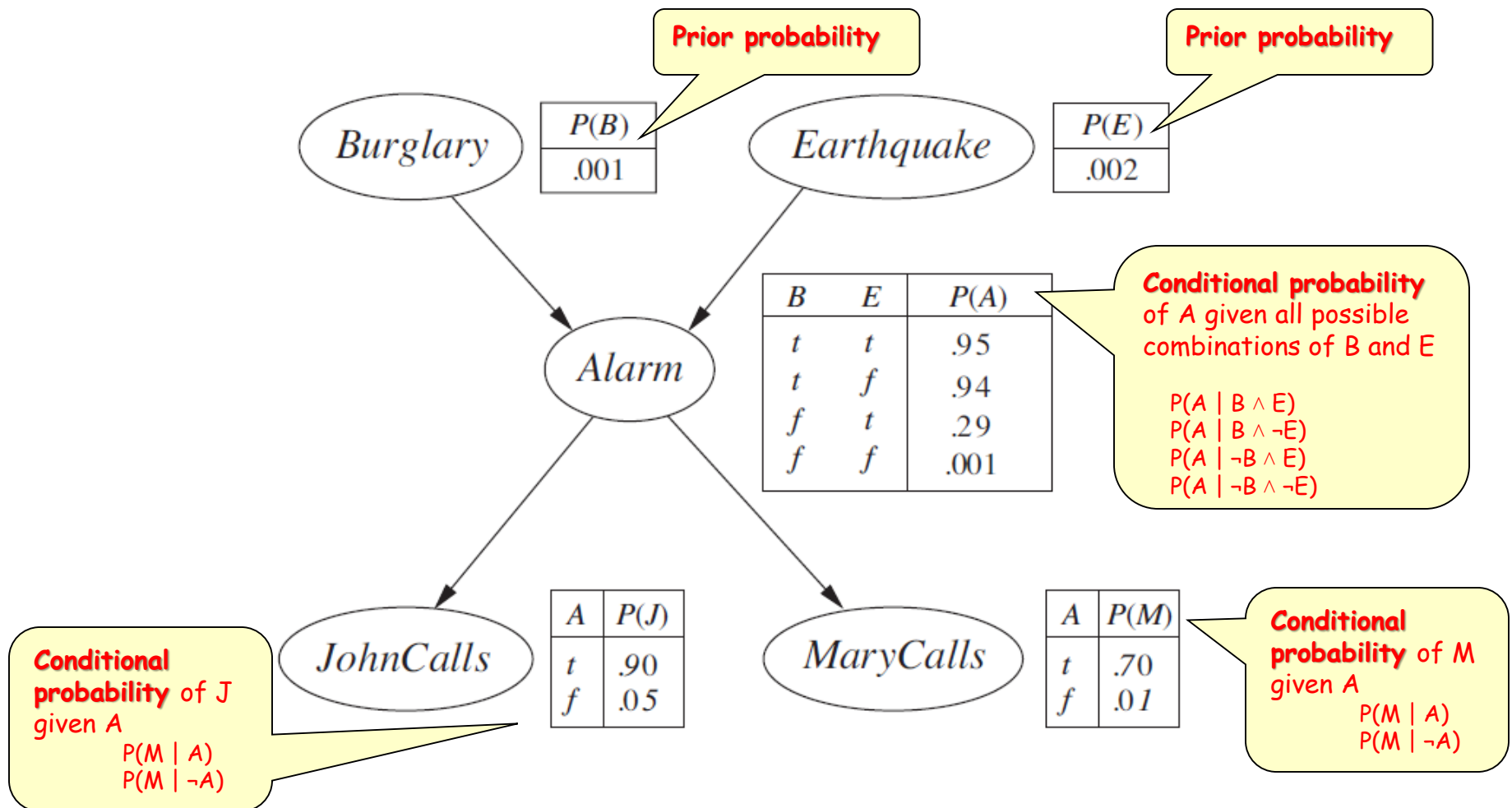
Bayesian network *(another example)*

- Both the **topology** and the **conditional probability tables** (CPTs)



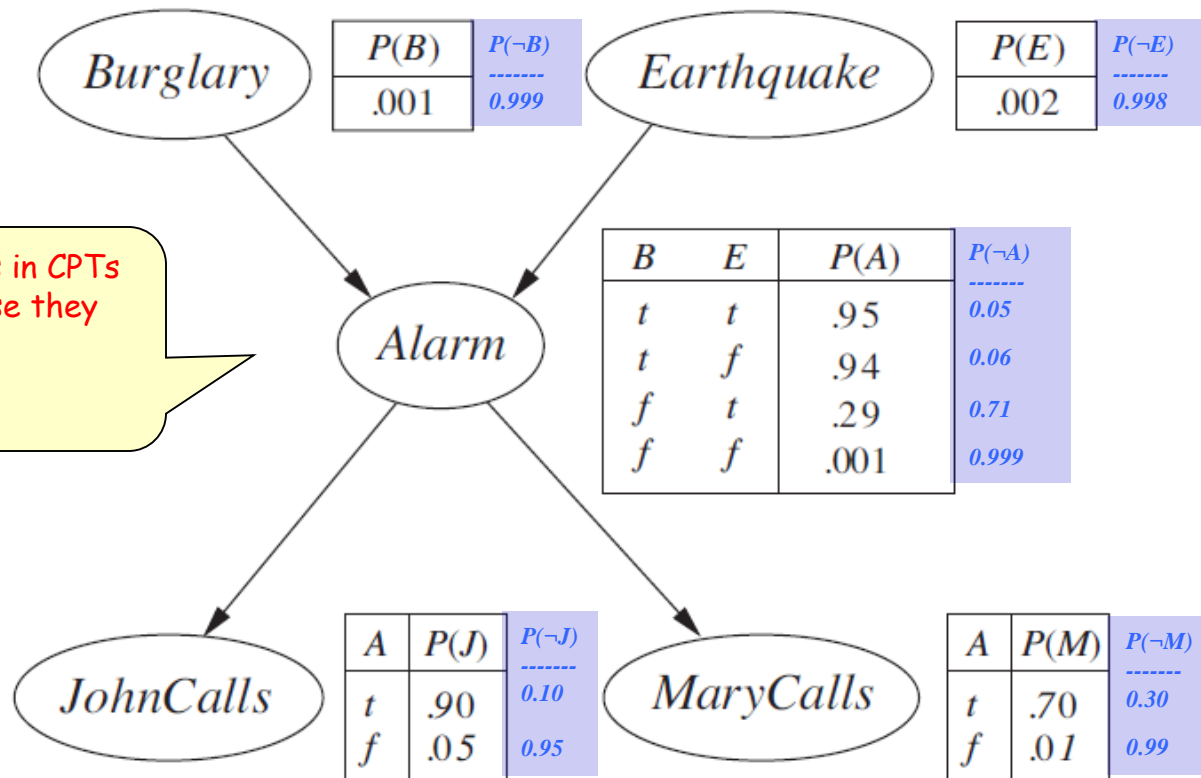
Bayesian network (*another example*)

- Both the **topology** and the **conditional probability tables (CPTs)**



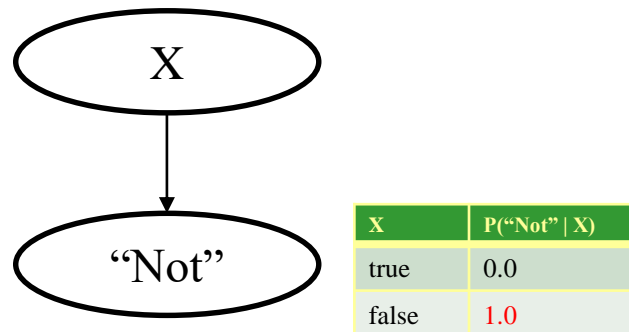
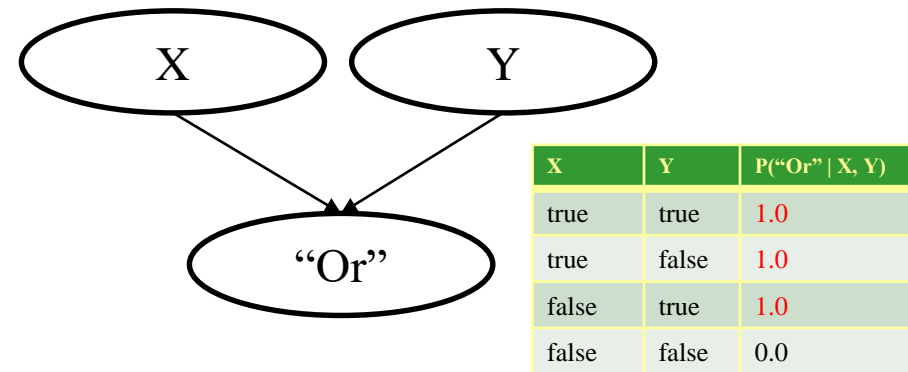
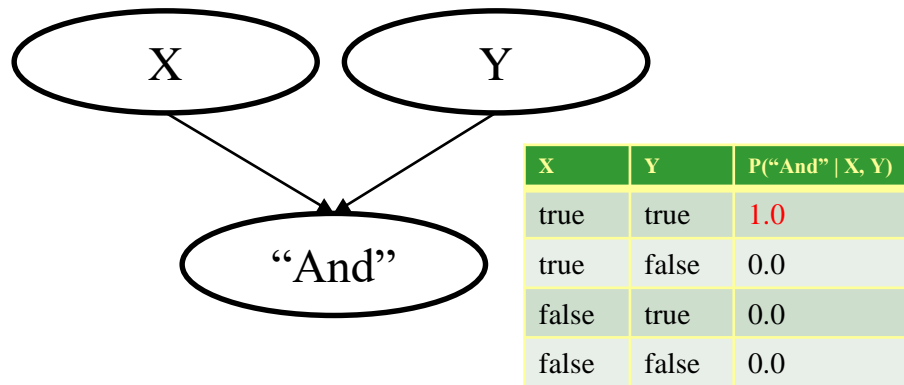
Bayesian network (*another example*)

- Both the **topology** and the **conditional probability tables** (CPTs)



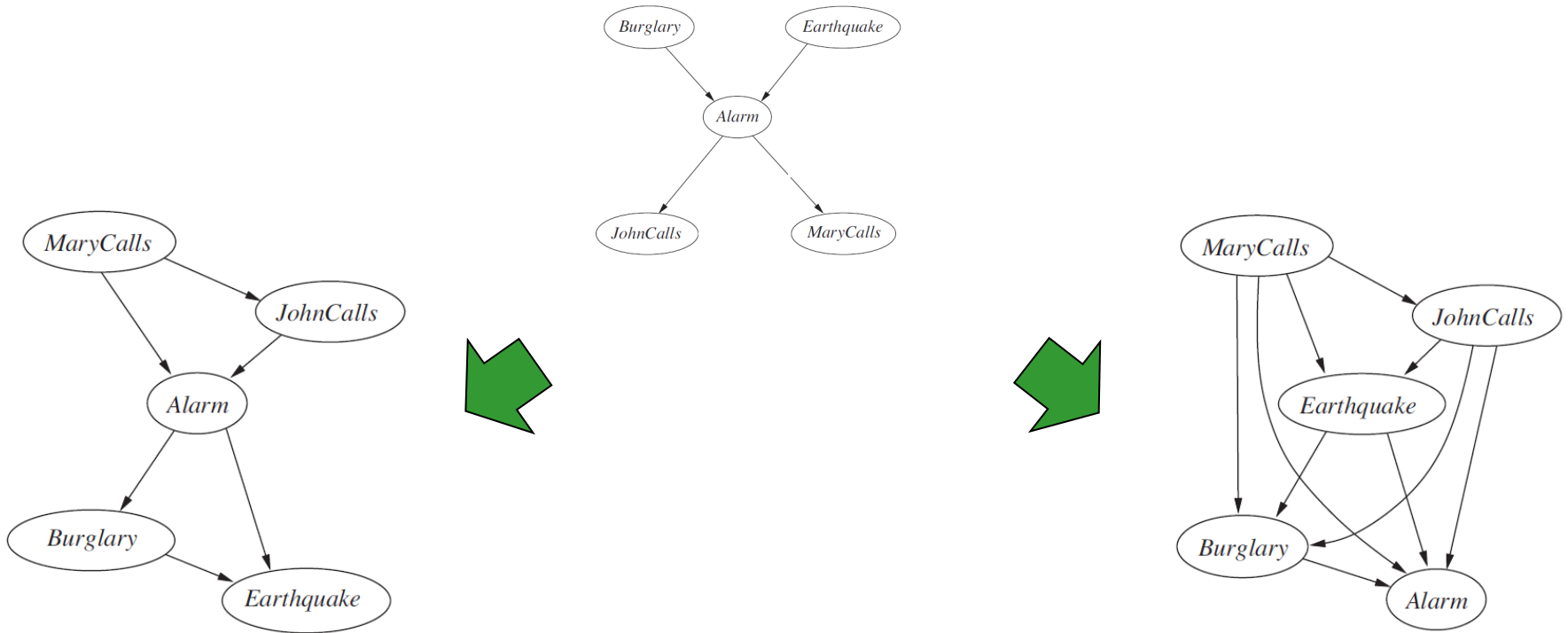
Bayesian networks for Boolean functions

- Can Bayesian networks represent all **Boolean functions**? – Yes.
 $f(\text{Feature_1}, \dots, \text{Feature_n}) \equiv \text{some propositional sentence}$



Compactness and node ordering

- There are multiple, **equivalent**, Bayesian networks, some of which are more **compact** than the others

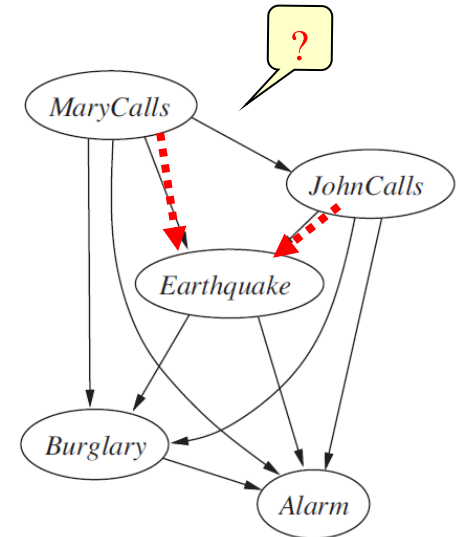
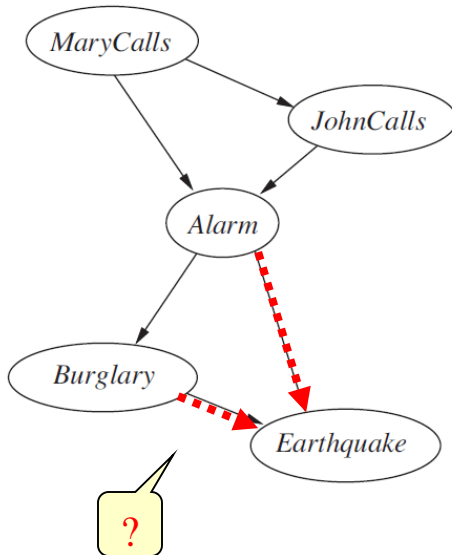
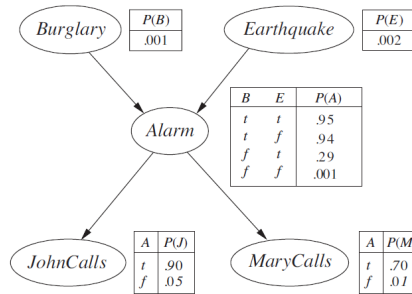


Compactness depend on the node ordering

Compactness and node ordering

- There are multiple, **equivalent**, Bayesian networks, some of which are more **compact** than the others

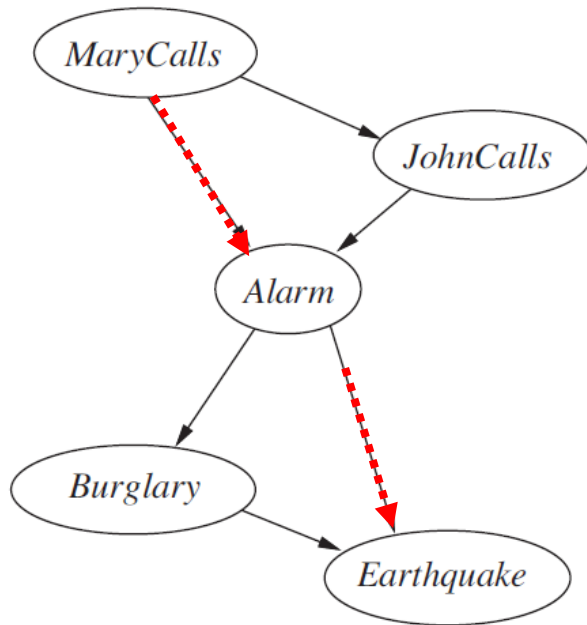
Not only more edges, but also tenuous relationships that require unnatural probability judgments



Compactness depend on the node ordering

Compactness and node ordering

- Distinction between **causal** model and **diagnostic** model
 - If we try to build a **diagnostic model**, with links **from symptoms to causes**, we have to specify additional dependencies between otherwise independent causes
 - Example: from *MaryCalls* to *Alarm*, or from *Alarm* to *Burglary*



Solution: stick to a causal model

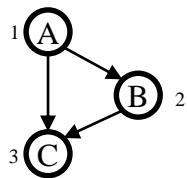
(1) Fewer dependences

(2) Easier to come up with probability

Let's go through an example

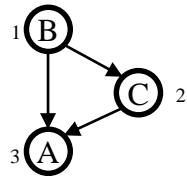
- Note: each way of **factoring joint distribution** corresponds to a **different Bayesian network**
- **Example:** 6 ways of factoring $P(A, B, C)$, including

- $P(A, B, C) = P(C \mid B, A) P(B, A) = P(C \mid B, A) P(B \mid A) P(A)$



(First picking A, then picking B, and finally picking C, each time conditioning the picked random variable on all random variables picked earlier)

- $P(A, B, C) = P(A \mid B, C) P(B, C) = P(A \mid B, C) P(C \mid B) P(B)$

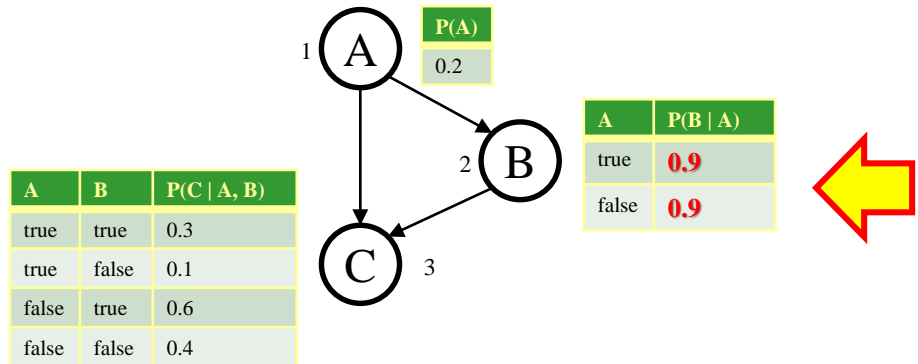


(First picking B, then picking C and finally picking A, each time conditioning the picked random variable on all random variables picked earlier)

Bayesian network is not unique

- The network topology determines how many probabilities need to be specified for the conditional probability tables.
 - Let's choose $P(A, B, C) = P(C \mid B, A) P(B \mid A) P(A)$.

A	B	C	P(A, B, C)
true	true	true	0.054
true	true	false	0.126
true	false	true	0.002
true	false	false	0.018
false	true	true	0.432
false	true	false	0.288
false	false	true	0.032
false	false	false	0.048

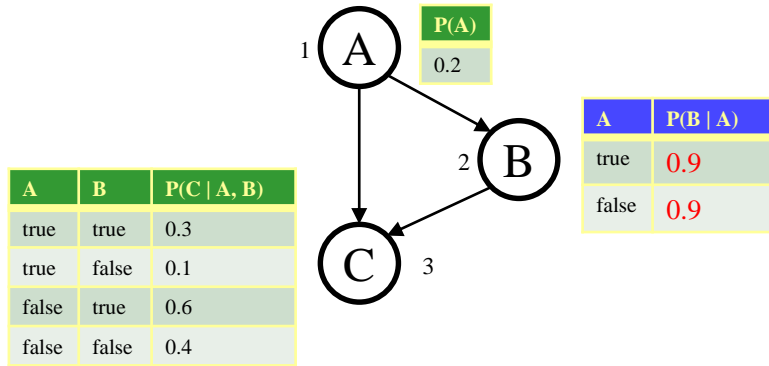


Independence detected

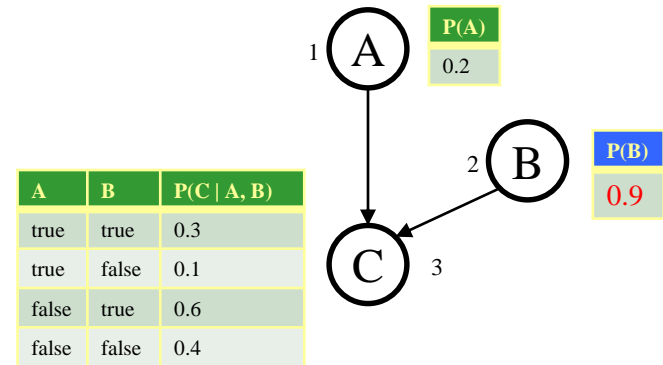
- Here: $P(B | A) = P(B | \neg A)$.
- Thus, A and B are **independent**
- Detailed explanation
 - $P(B) = P(B \wedge A) + P(B \wedge \neg A)$
 $= P(B | A) P(A) + P(B | \neg A) P(\neg A)$
 $= P(B | A) P(A) + P(B | A) P(\neg A)$
 $= P(B | A) (P(A) + P(\neg A))$
 $= P(B | A)$

Simplifying Bayesian network

- Independence allows us to **simplify** Bayesian network

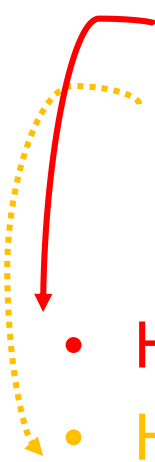


Need to **specify 7 probabilities** for all conditional probability tables



Need to specify **only 6 probabilities** for all conditional probability tables

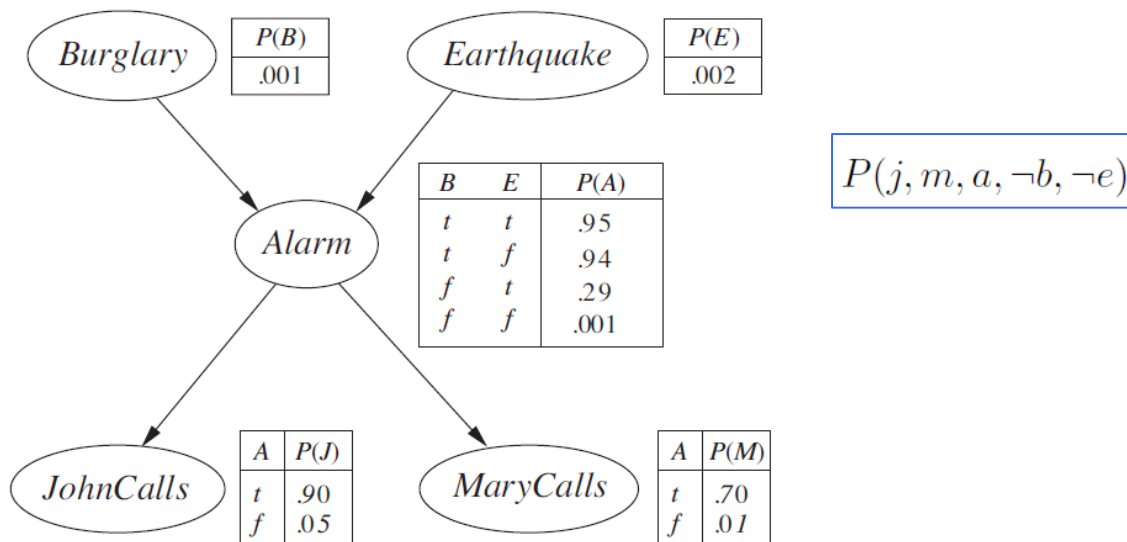
Semantics of Bayesian networks

- Two different, but equivalent **views**
 - Representation of the **joint probability distribution**
 - Encoding of a collection of **conditional independence** statements
 - How to construct networks
 - How to design inference procedures
- 
- A diagram consisting of two curved arrows. A solid red arrow starts from the first bullet point (Two different, but equivalent views) and points down to the third bullet point (How to construct networks). A dashed yellow arrow starts from the second bullet point (Encoding of a collection of conditional independence statements) and points down to the fourth bullet point (How to design inference procedures).

Bayesian network: *Representing the full joint distribution*

- Each entry $P(x_1, \dots, x_n)$ in the full joint distribution, which is the abbreviation of $P(X_1=x_1 \wedge \dots \wedge X_n=x_n)$, is the product of the elements of the CPTs defined as follows:

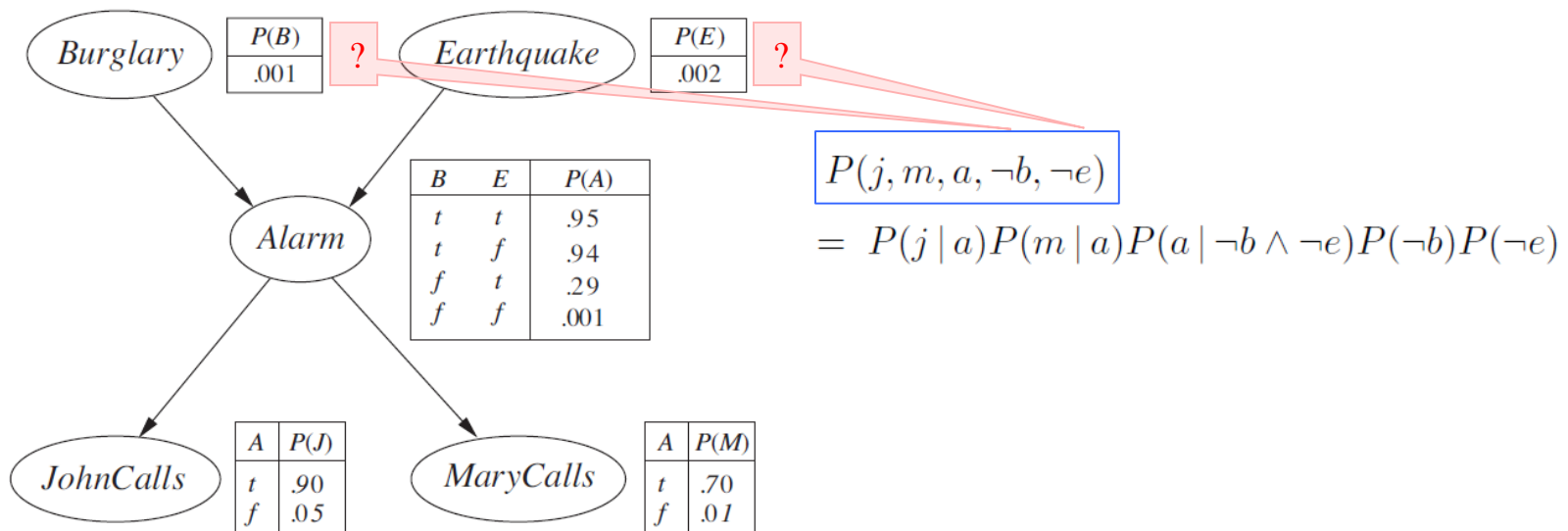
$$P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i \mid \text{parents}(X_i))$$



Bayesian network: *Representing the full joint distribution*

- Each entry $P(x_1, \dots, x_n)$ in the full joint distribution, which is the abbreviation of $P(X_1=x_1 \wedge \dots \wedge X_n=x_n)$, is the product of the elements of the CPTs defined as follows:

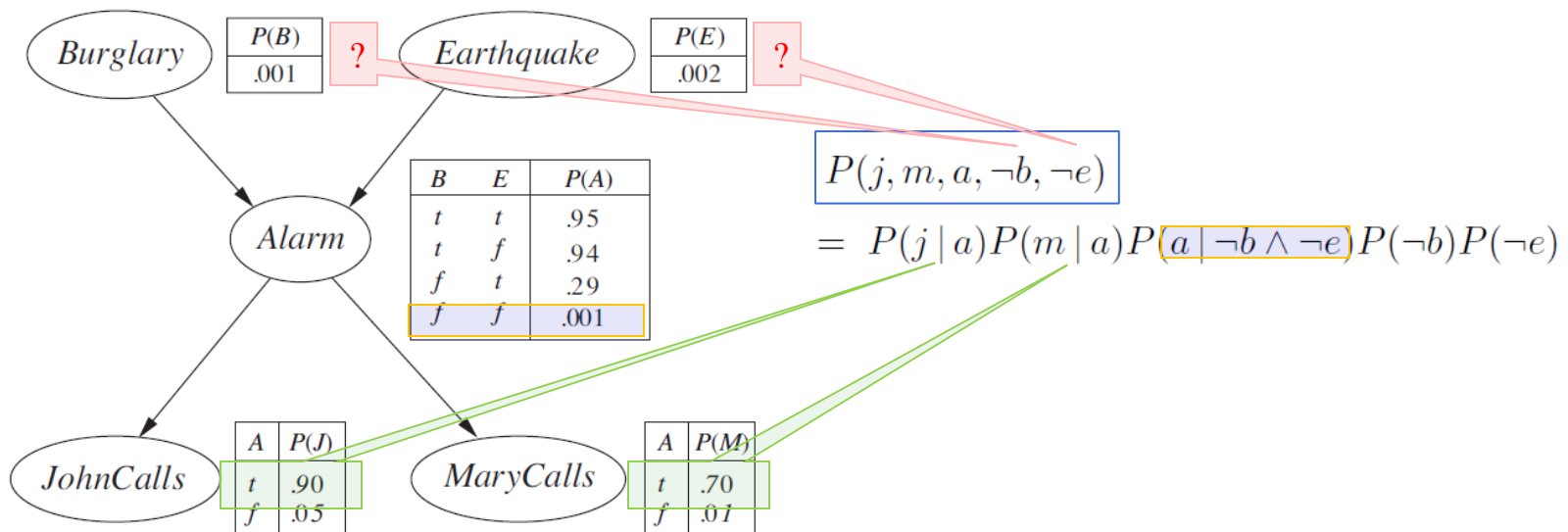
$$P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i \mid \text{parents}(X_i))$$



Bayesian network: *Representing the full joint distribution*

- Each entry $P(x_1, \dots, x_n)$ in the full joint distribution, which is the abbreviation of $P(X_1=x_1 \wedge \dots \wedge X_n=x_n)$, is the product of the elements of the CPTs defined as follows:

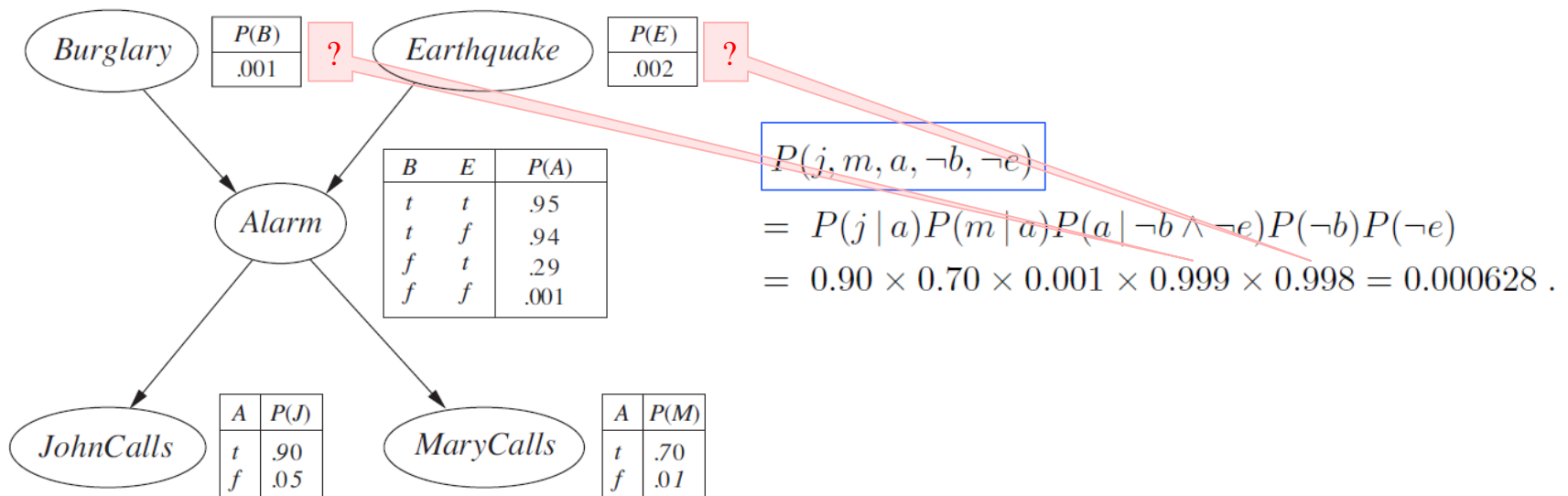
$$P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i \mid \text{parents}(X_i))$$



Bayesian network: *Representing the full joint distribution*

- Each entry $P(x_1, \dots, x_n)$ in the full joint distribution, which is the abbreviation of $P(X_1=x_1 \wedge \dots \wedge X_n=x_n)$ is the product of the elements of the CPTs defined as follows:

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



How to construct Bayesian networks

- Starting from the full joint distribution, first, we rewrite the entries in terms of conditional probability

$$P(x_1, \dots, x_n) = P(x_n | x_{n-1}, \dots, x_1) P(x_{n-1}, \dots, x_1)$$

- Then, we repeat the process

$$\begin{aligned} P(x_1, \dots, x_n) &= P(x_n | x_{n-1}, \dots, x_1) P(x_{n-1} | x_{n-2}, \dots, x_1) \cdots P(x_2 | x_1) P(x_1) \\ &= \prod_{i=1}^n P(x_i | x_{i-1}, \dots, x_1) . \end{aligned}$$

- Now, compare to Bayesian network

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

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

Each node must be **conditionally independent** of its other predecessors in the node ordering, given its parents

How to construct Bayesian networks *(cont'd)*

- Starting from the full joint distribution, first, we rewrite the entries in terms of conditional probability

$$\begin{aligned} P(x_1, \dots, x_n) &= P(x_n | x_{n-1}, \dots, x_1) P(x_{n-1} | x_{n-2}, \dots, x_1) \cdots P(x_2 | x_1) P(x_1) \\ &= \prod_{i=1}^n P(x_i | x_{i-1}, \dots, x_1) . \end{aligned}$$

- Now, compare to Bayesian network

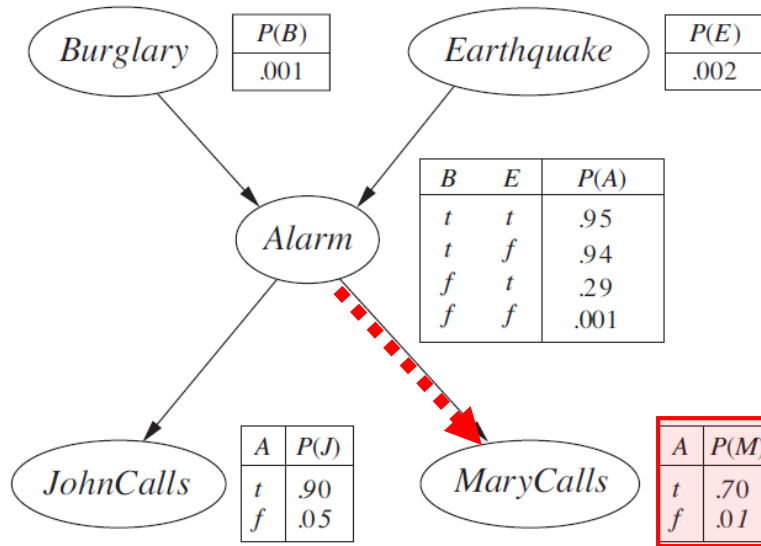
$$P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i | \text{parents}(X_i))$$

Each node must be **conditionally independent** of its other predecessors in the node ordering, given its parents

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

- For $i = 1$ to n do:
 - Choose, from X_1, \dots, X_{i-1} , a minimum set of parents for X_i , to satisfy the equation
 - Insert edges from the parents to X_i
 - CTPs: write down the conditional probability table, $P(X_i | \text{Parents}(X_i))$

How to construct Bayesian networks (cont'd)



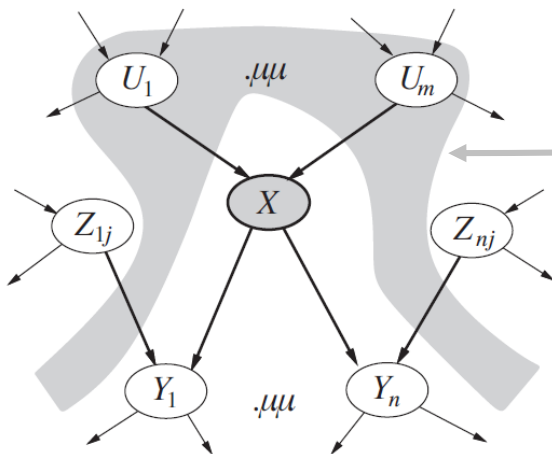
$$P(\text{MaryCalls} \mid \text{JohnCalls}, \text{Alarm}, \text{Earthquake}, \text{Burglary}) \\ = P(\text{MaryCalls} \mid \text{Alarm})$$

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

- For $i = 1$ to n do:
 - Choose, from X_1, \dots, X_{i-1} , a minimum set of parents for X_i , to satisfy the equation
 - Insert edges from the parents to X_i
 - CTPs: write down the conditional probability table, $P(X_i \mid \text{Parents}(X_i))$

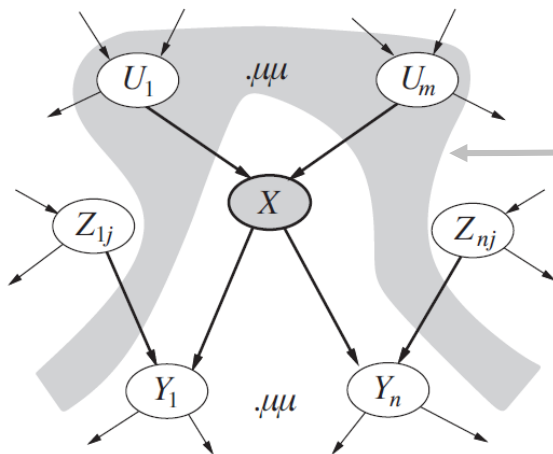
Semantics of Bayesian networks

- Two different, but equivalent views
 - Representation of the joint probability distribution
 - Encoding of a collection of **conditional independence** statements
- How to construct Bayesian networks
 - a node is conditionally independent of its other predecessors, given its parents
- **X is conditionally independent of its non-descendants, given its parents**



Semantics of Bayesian networks

- Two different, but equivalent views
 - Representation of the joint probability distribution
 - Encoding of a collection of **conditional independence** statements
- How to construct Bayesian networks
 - a node is conditionally independent of its other predecessors, given its parents
- **X is conditionally independent of its non-descendants, given its parents**



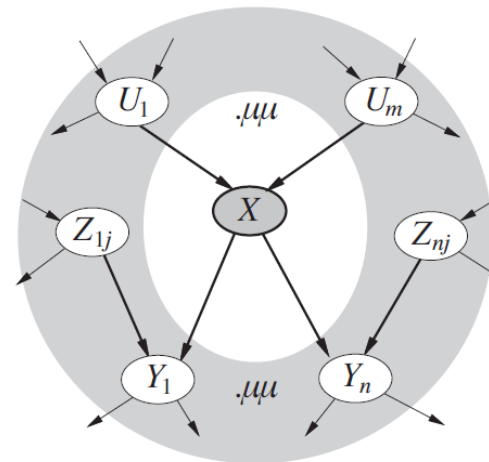
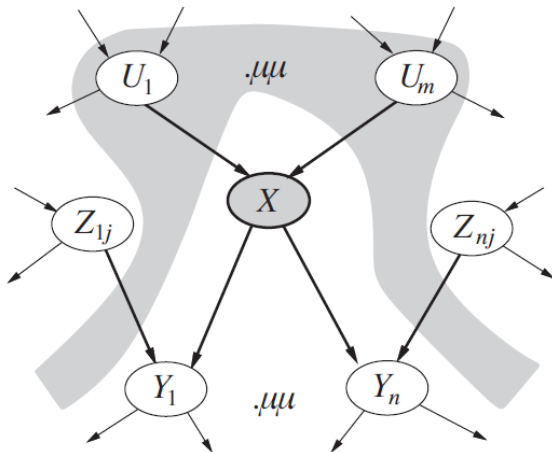
X may even be conditionally independent of some of these descendants

Semantics of Bayesian networks

- Two different, but equivalent views
 - Representation of the joint probability distribution
 - Encoding of a collection of conditional independence statements
- How to construct Bayesian networks
 - a node is conditionally independent of its other predecessors
- **X is conditionally independent of all other nodes, given Markov blanket**

Markov blanket of X
includes

- (1) its parents
- (2) its children,
- (3) its children's parents



Summary

- Bayesian network is a well-developed representation for uncertain knowledge
 - It is often **exponentially smaller** than full joint distribution
 - It plays similar role as propositional logic for definite knowledge
- Bayesian network is a DAG where
 - a node denotes a random variable, together with local conditional distribution of that variable, given its parents
 - It's a **concise representation** of conditional independence
 - It represents full joint distribution, as product of corresponding entries in the local conditional distributions