# Probabilistic representation

Applied artificial intelligence (EDA132)
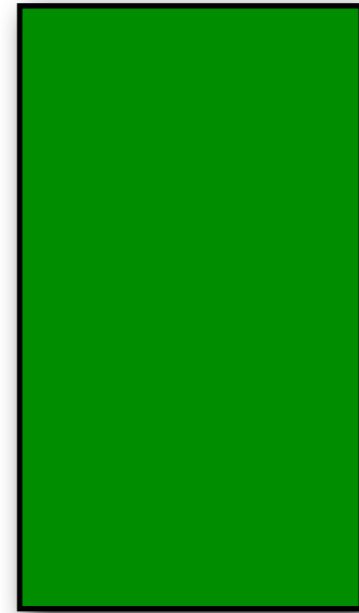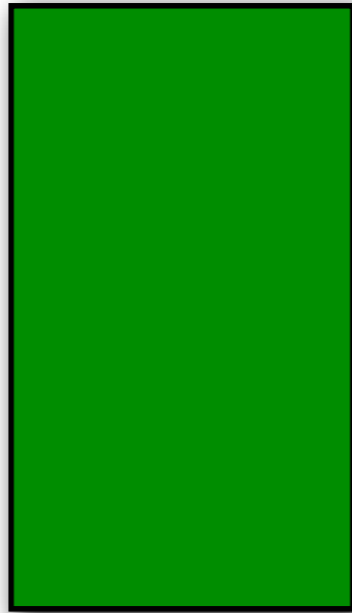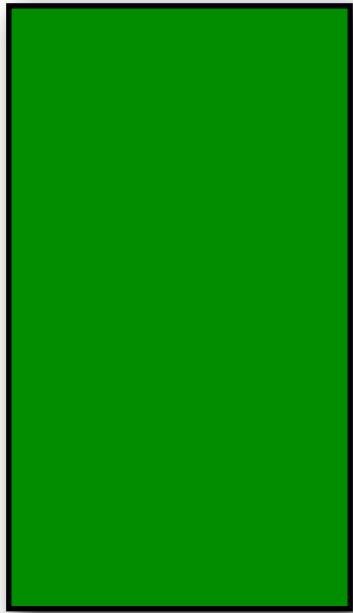Lecture 05
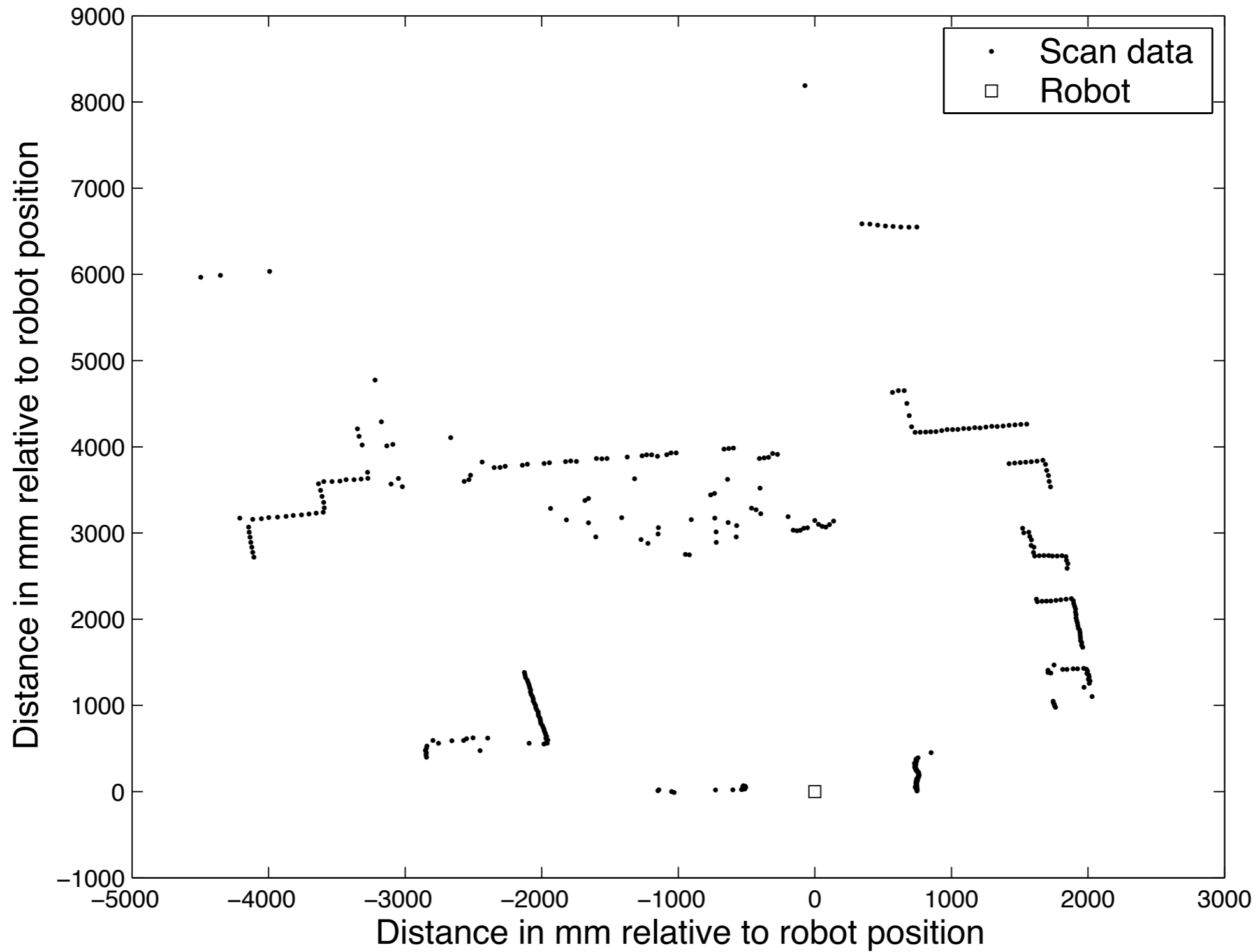2014-02-04
Elin A. Topp

1

# Show time!

Two boxes of chocolates, one luxury car.
Where is the car?



Philosopher: It does not matter whether I change my choice, I will either get chocolates or a car.

Mathematician: It is more likely to get the car when I change my choice - even though it is not certain!

# A robot's view of the world...

# Outline

- Uncertainty & probability (chapter 13)

  - Uncertainty

  - Probability

  - Syntax and Semantics

  - Inference

  - Independence and Bayes' Rule

- Bayesian Networks (chapter 14.1-3)

  - Syntax

  - Semantics

  - Efficient representation

# Outline

- Uncertainty & probability (chapter 13)

  - Uncertainty

  - Probability

  - Syntax and Semantics

  - Inference

  - Independence and Bayes' Rule

- Bayesian Networks (chapter 14.1-3)

  - Syntax

  - Semantics

  - Efficient representation

# 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                 ???          No!

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

"Causal" rule:  Cavity $\Rightarrow$ Toothache                 ???             Well... not always

x

# 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)

$\Rightarrow$ better use **probabilities** to represent certain **knowledge states**

$\Rightarrow$ Rational decisions (decision theory) combine probability and utility theory

X

# 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., $P(A_{25} \mid \textit{no reported accidents}) = 0.06$

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

# Making decisions under uncertainty

Suppose the following believes (from past experience):

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

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

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

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

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

X

# Probability basics

A set $\Omega$ - 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:

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

$$\Sigma_\omega P(\omega) = 1$$

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

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

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

x

# Random variables

A random variable is a function from sample points to some range, e.g., the reals or Booleans,

e.g., *Odd( 1 ) = 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(Odd = true) = P(1) + P(3) + P(5) = 1/6 + 1/6 + 1/6 = 1/2*

x

# 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(ω) = true*

*event ¬a = set of sample points where A(ω) = false*

*event a∧b = points where A(ω) = true and B(ω) = 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( Cavity = true) = 0.2* and

*P( Weather = sunny) = 0.72*

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

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

**P***(Weather) = ⟨0.72, 0.1, 0.08, 0.1⟩*

*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***(Weather, Cavity) =* a *4 x 2* matrix of values:

| Weather<br>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( cavity) = 0.2* (prior, no evidence for anything), but

*P( cavity | 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( cavity | toothache, sunny) = P( cavity | toothache)*

*Domain knowledge* lets us do this kind of inference.

# Posterior probability (2)

Definition of conditional / posterior probability:

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

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

$$P(\,a \wedge b) \quad = \quad P(\,a\mid b)\,P(\,b) \quad = \quad P(\,b\mid a)\,P(\,a)$$

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

$$\mathbf{P}(\,Weather, Cavity) \quad = \quad \mathbf{P}(\,Weather\mid Cavity)\,\mathbf{P}(\,Cavity)$$

(gives a *4x2* set of equations)

*Chain rule* (successive application of product rule):

$$\mathbf{P}(\,X_1, ..., X_n) \;= \mathbf{P}(\,X_1, ..., X_{n-1})\,\mathbf{P}(\,X_n\mid X_1, ..., X_{n-1})$$

$$= \mathbf{P}(\,X_1, ..., X_{n-2})\,\mathbf{P}(\,X_{n-1}\mid X_1, ..., X_{n-1})\,\mathbf{P}(\,X_n\mid X_1, ..., X_{n-1})$$

$$= ... = \prod_{i=1}^{n} \mathbf{P}(\,X_i\mid X_1, ..., X_{i-1})$$

# 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 $\Phi$, sum the atomic events where it is true:

Can also compute posterior probabilities

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

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

$$P(toothache) = 0.108 + 0.012 + 0.072 + 0.008 + 0.016 + 0.064 = 0.2$$

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

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

# Normalisation

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

Denominator can be viewed as a *normalisation constant*:

$P$*( Cavity | toothache) =* α $P$*( Cavity, toothache)*

$\quad$ = α [$P$*( Cavity, toothache, catch) +* $P$*( Cavity, toothache, ¬catch)]*

$\quad$ = α [⟨0.108, 0.016⟩ + ⟨0.012, 0.064⟩]

$\quad$ = α ⟨0.12, 0.08⟩ = ⟨0.6, 0.4⟩

And the good news:

We can compute $P$*( Cavity | toothache)* without knowing the value of *P( toothache)*!

# 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"

Would you agree?

> Psychologist knowing "the math" argues:
>
> $P(\text{ headache } | \text{ PTSD}) = high$ (statistics)
>
> $P(\text{ PTSD } | \text{ abuse in childhood}) = high$ (statistics)
>
> ok, yes, sure, but:
>
> You did not consider the relevant relations of
>
> $P(\text{ PTSD } | \text{ headache})$ or
>
> $P(\text{ abuse in childhood } | \text{ PTSD}),$
>
> i.e., you mixed up cause and effect in your argumentation!

# Bayes' Rule

Recap *product rule:* $P( a \land b) = P( a \mid b) P( b) = P( b \mid a) P(a)$

$\Rightarrow$ Bayes' Rule $P( a \mid b) = \dfrac{P( b \mid a) P( a)}{P( b)}$

or in distribution form:

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

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

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

E.g., with $M$ "meningitis", $S$ "stiff neck":

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

15

# All is well that ends well ...

We can model cause-effect relationships,

we can base our judgement on mathematically sound inference,

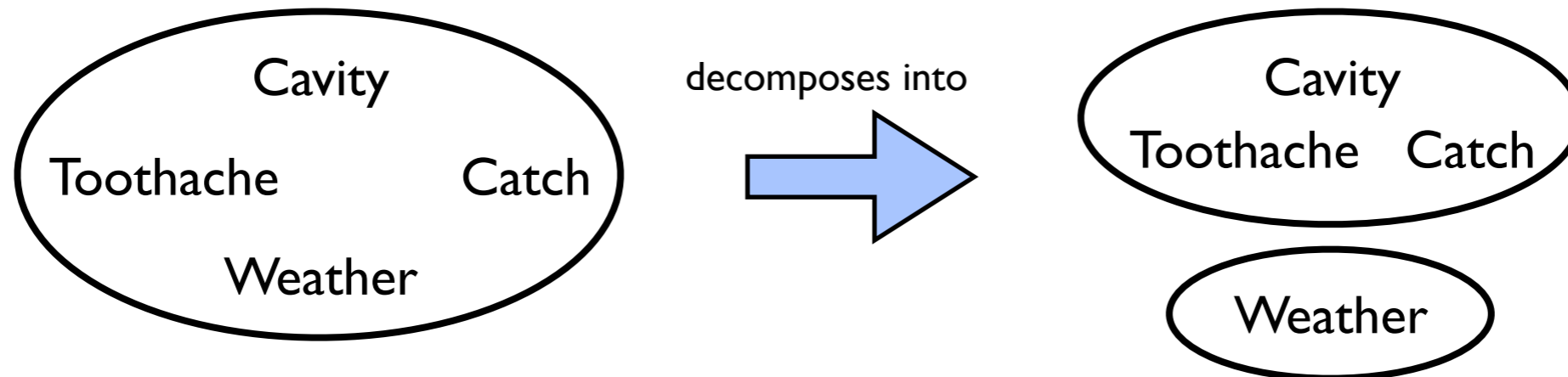we can even do this inference with only partial knowledge on the priors, ...

# ... but

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

(and for non-Booleans it gets even more nasty...)

# Independence

*A* and *B* are *independent* iff

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



decomposes into

*P( Toothache, Catch, Cavity, Weather)* = *P( Toothache, Catch, Cavity) P( 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( Toothache, Cavity, 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( catch | toothache, cavity) = P( catch | cavity)*

The same holds when there is no cavity:

(2) *P( catch | toothache, ¬cavity) = P( catch | ¬cavity)*

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

*P( Catch |Toothache, Cavity) = P( Catch | Cavity)*

Writing out full joint distribution using chain rule:

*P( Toothache, Catch, Cavity)*
*= P( Toothache | Catch, Cavity) P( Catch, Cavity)*
*= P( Toothache | Catch, Cavity) P( Catch | Cavity) P( Cavity)*
*= P( Toothache | Cavity) P( Catch | Cavity) P( Cavity)*

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 $n$.
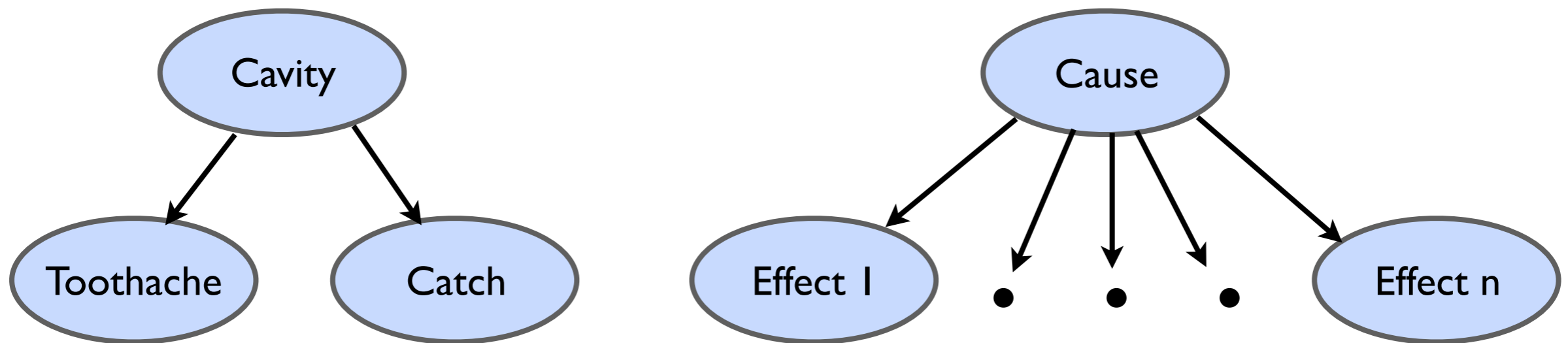
Hence:

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

# Bayes' Rule and conditional independence

*P( Cavity | toothache ∧ catch)*

= α  *P( toothache ∧ catch | Cavity) P( Cavity)*

= α  *P( toothache | Cavity) P( catch | Cavity) P( Cavity)*

An example of a *naive Bayes* model:

$$P( Cause, Effect_1, ...., Effect_n) =  P( Cause) \prod_i P( Effect_i | Cause)$$



The total number of parameters is *linear* in *n*

# 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

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

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

*Independence* and *conditional independence* provide the tools

# Outline

- Uncertainty & probability (chapter 13)

  - Uncertainty

  - Probability

  - Syntax and Semantics

  - Inference

  - Independence and Bayes' Rule

- **Bayesian Networks (chapter 14.1-3)**

  - Syntax

  - Semantics

  - Efficient representation

# Outline

- Uncertainty (chapter 13)

  - Uncertainty

  - Probability

  - Syntax and Semantics

  - Inference

  - Independence and Bayes' Rule

- **Bayesian Networks (chapter 14.1-3)**

  - **Syntax**

  - Semantics

  - Efficient representation

# 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:
          *$P( X_i | Parents( X_i))$*

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
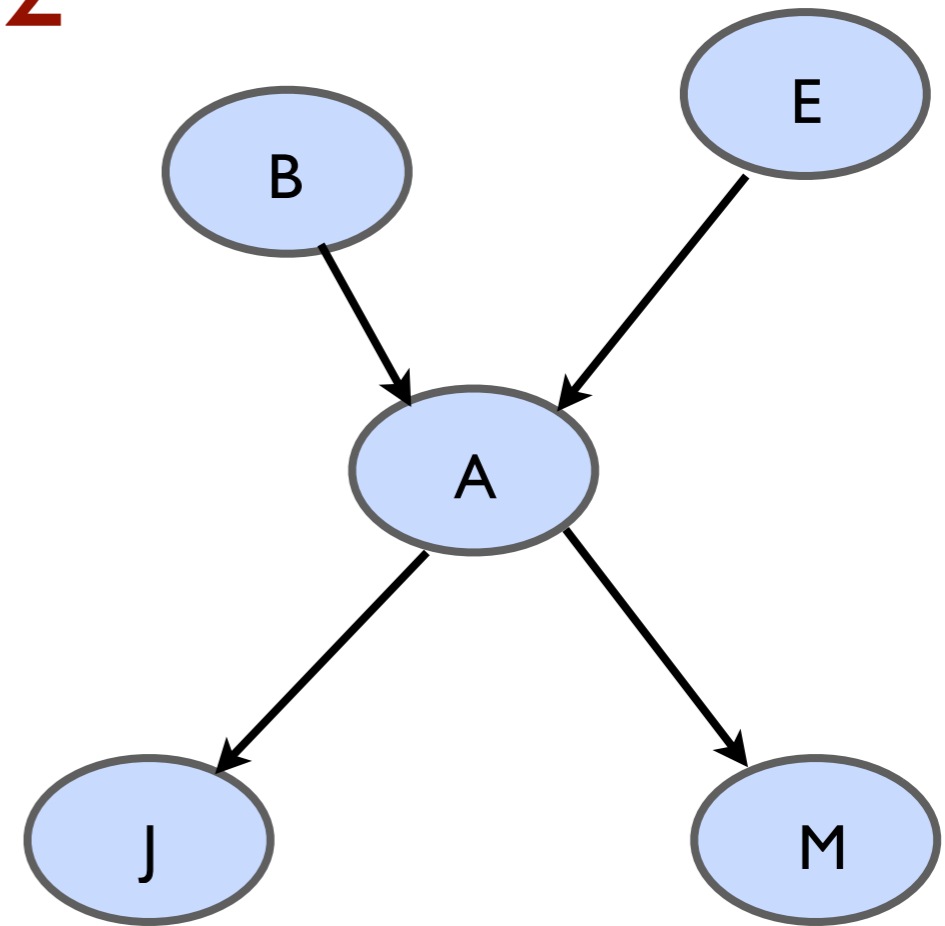
The alarm can cause Mary to call

# Example 2 (2)



| B | E | P(A\|B,E) |
|---|---|---------|
| T | T | 0.95 |
| T | F | 0.94 |
| F | T | 0.29 |
| F | F | 0.001 |

| P(B) |
|------|
| 0.001 |

| P(E) |
|------|
| 0.002 |

| A | P(J\|A) |
|---|--------|
| T | 0.90 |
| F | 0.05 |

| A | P(M\|A) |
|---|--------|
| T | 0.70 |
| F | 0.01 |

# 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 $X_i$ = *true*
(the number for $X_i$ = *false* is just $1-p$)

If each variable has no more than $k$ parents,
the complete network requires $O(n\ 2^k)$ 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 & probability (chapter 13)

  - Uncertainty

  - Probability

  - Syntax and Semantics

  - Inference

  - Independence and Bayes' Rule

- Bayesian Networks (chapter 14.1-3)

  - Syntax

  - Semantics

  - Efficient representation

# Global semantics

*Global* semantics defines the full joint distribution as the product of the local conditional distributions:
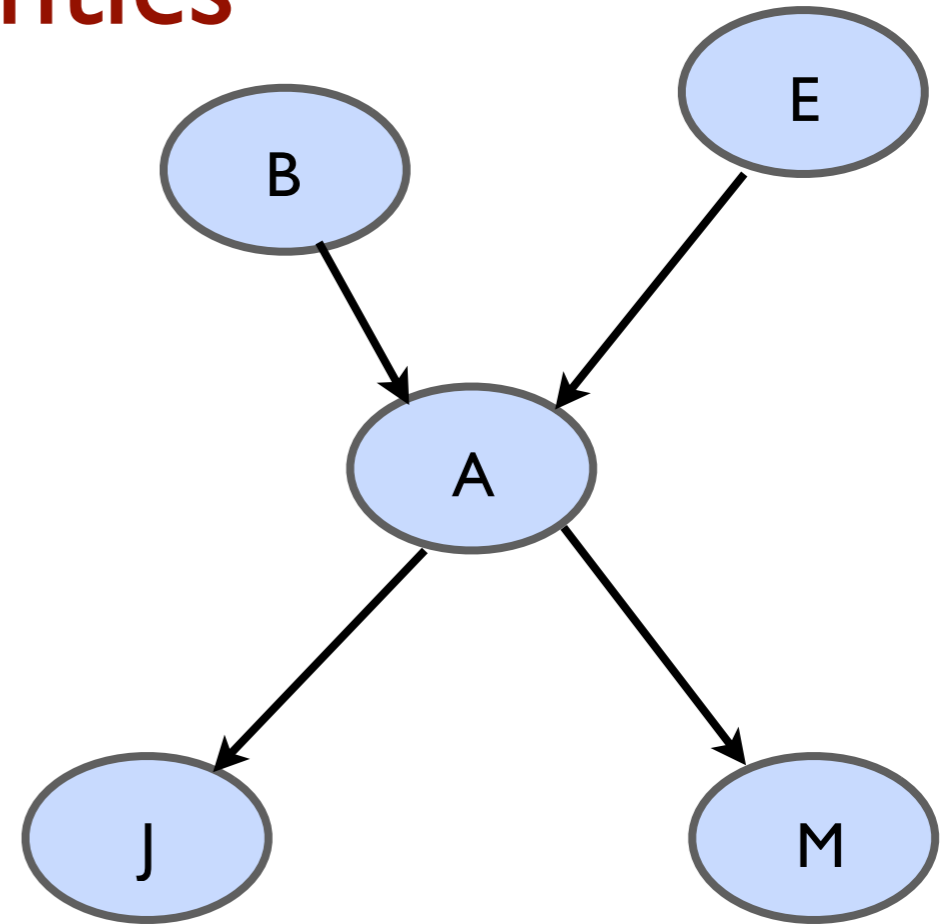
$$P(x_1, ..., x_n) = \prod_{i=1}^{n} P(x_i \mid parents(X_i))$$

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 $X_1, ..., X_n$

2. For $i = 1$ to $n$

      add $X_i$ to the network

      select parents from $X_1, ..., X_{i-1}$ such that

$$P( X_i | Parents( X_i)) = P( X_i | X_1, ..., X_{i-1} )$$
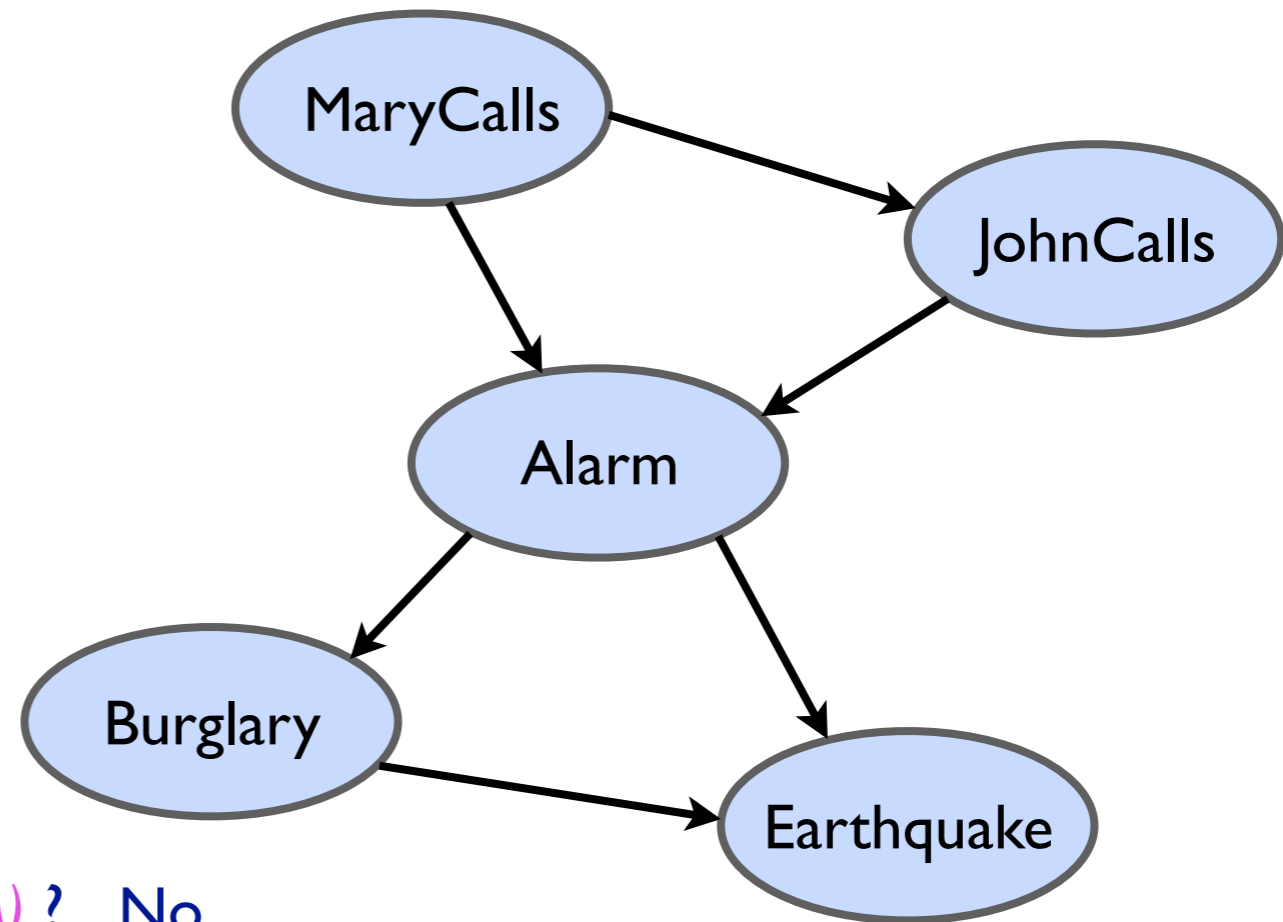
This choice of parents guarantees the global semantics:

$$P( X_1, ..., X_n )  =  \prod_{i=1}^{n} P( X_i | X_1, ..., X_{i-1} ) \qquad \text{(chain rule)}$$

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

# Construction example

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



$P( J \mid M) = P( J)$ ?  No

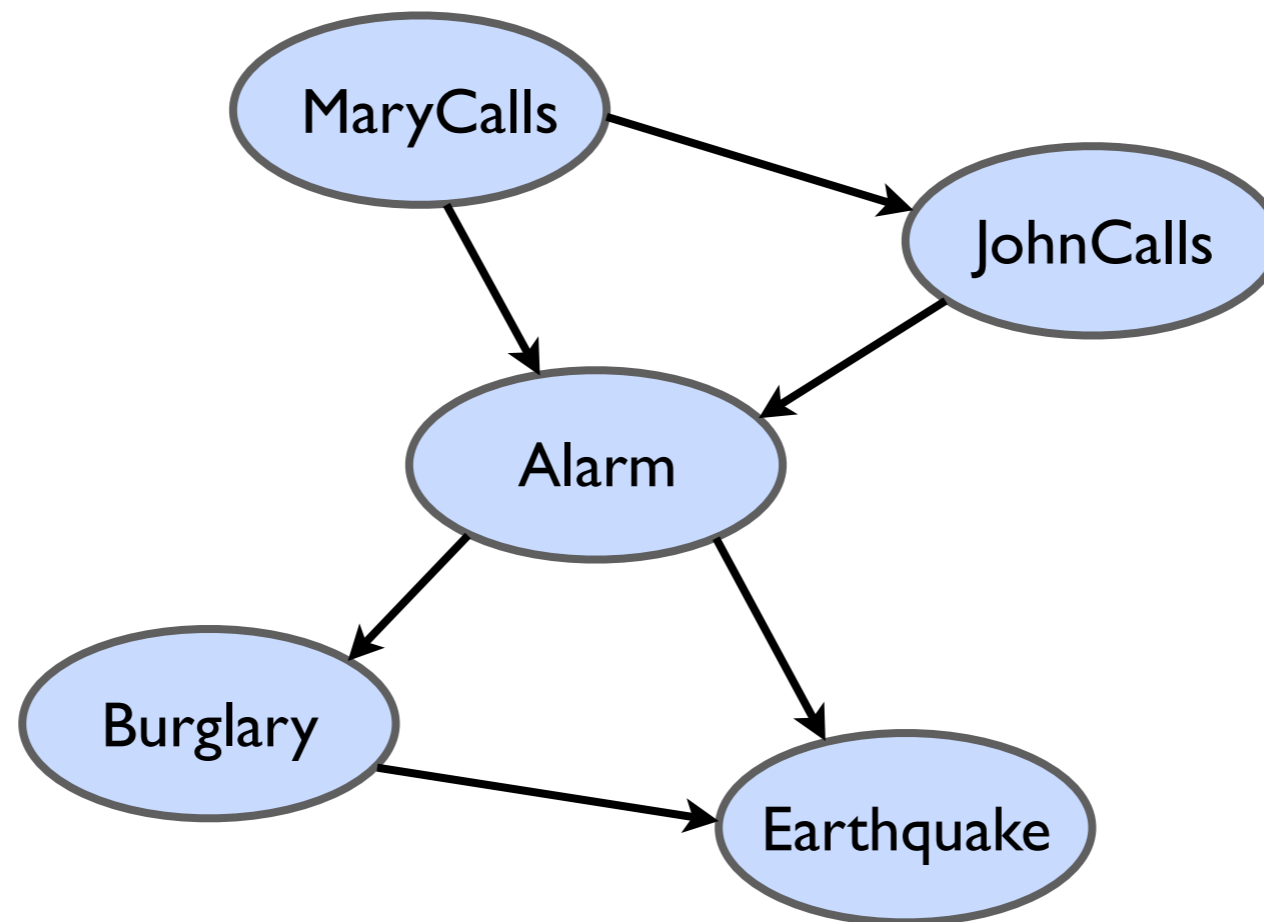$P( A \mid J, M) = P( A \mid J)$ ? $P( A \mid J, M) = P( A)$ ?   No

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

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

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

$P( E \mid B, A, J, M) = P( E \mid A, B)$ ?  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
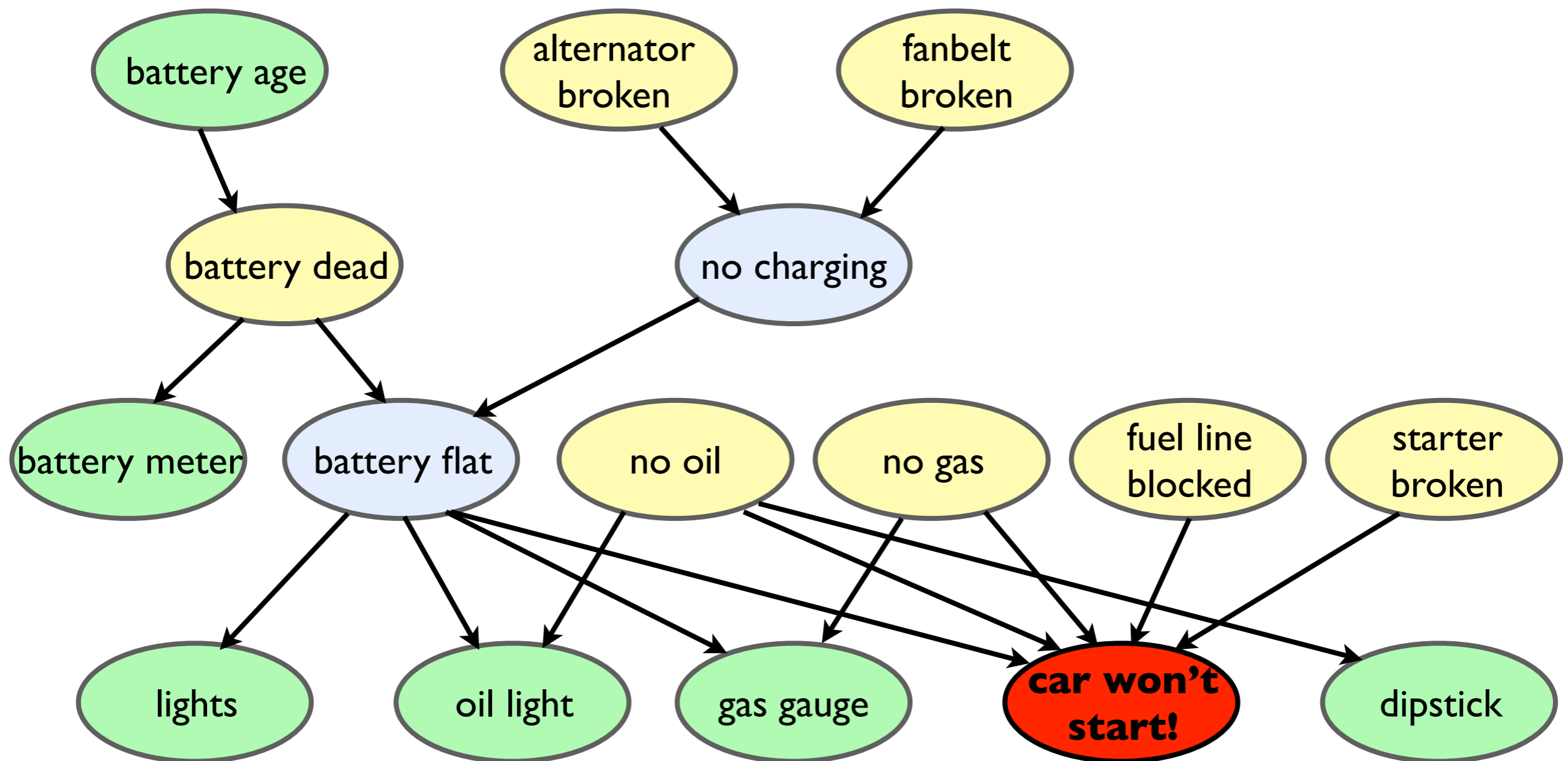
Network is less compact: *1 + 2 + 4 +2 +4 = 13* numbers

Hence: Choose preferably an order corresponding to the cause → effect "chain"

# Locally structured (sparse): Car diagnosis

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

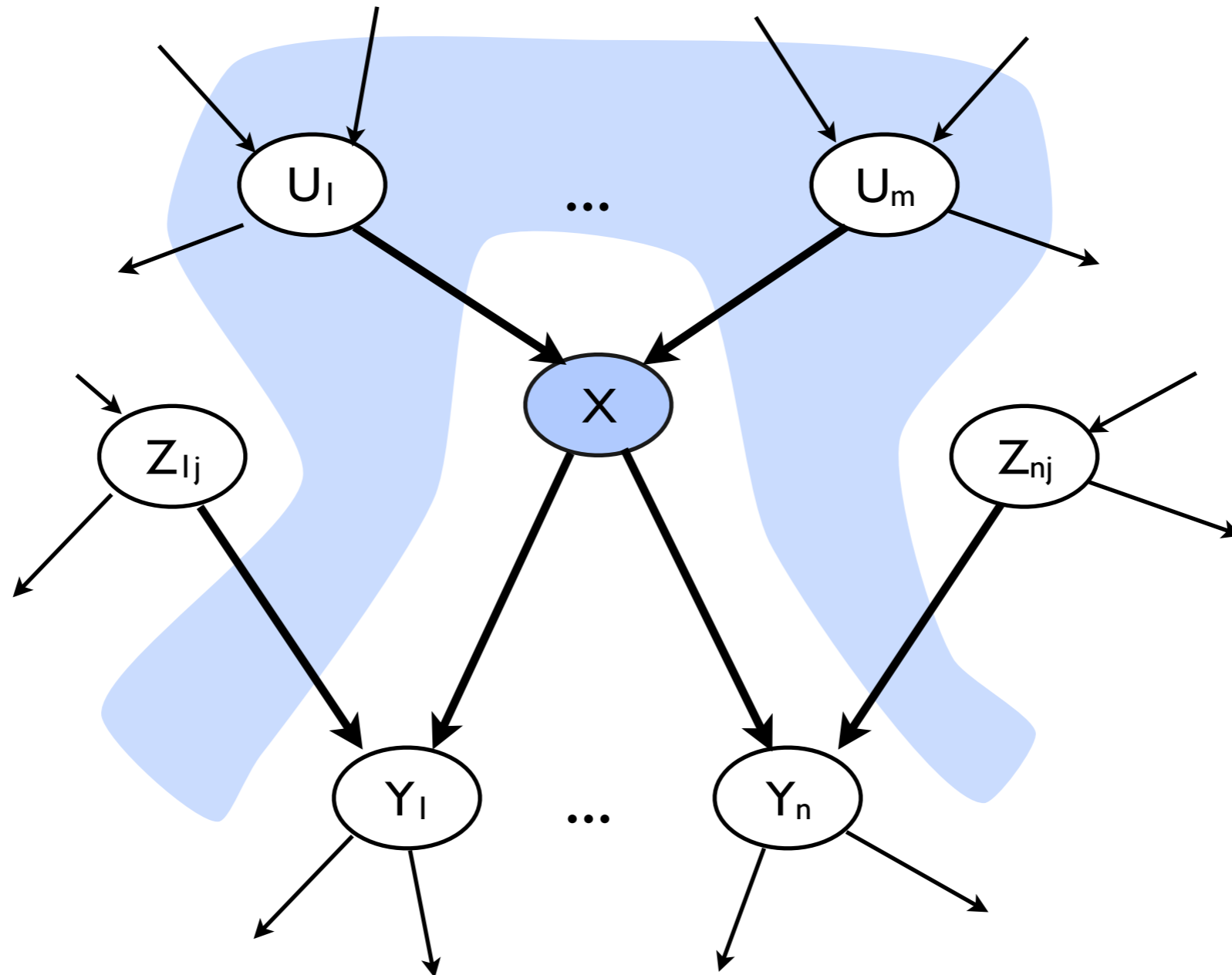Testable variables (green), "broken, so fix it" variables (yellow)

Hidden variables (blue) ensure sparse structure / reduce parameters
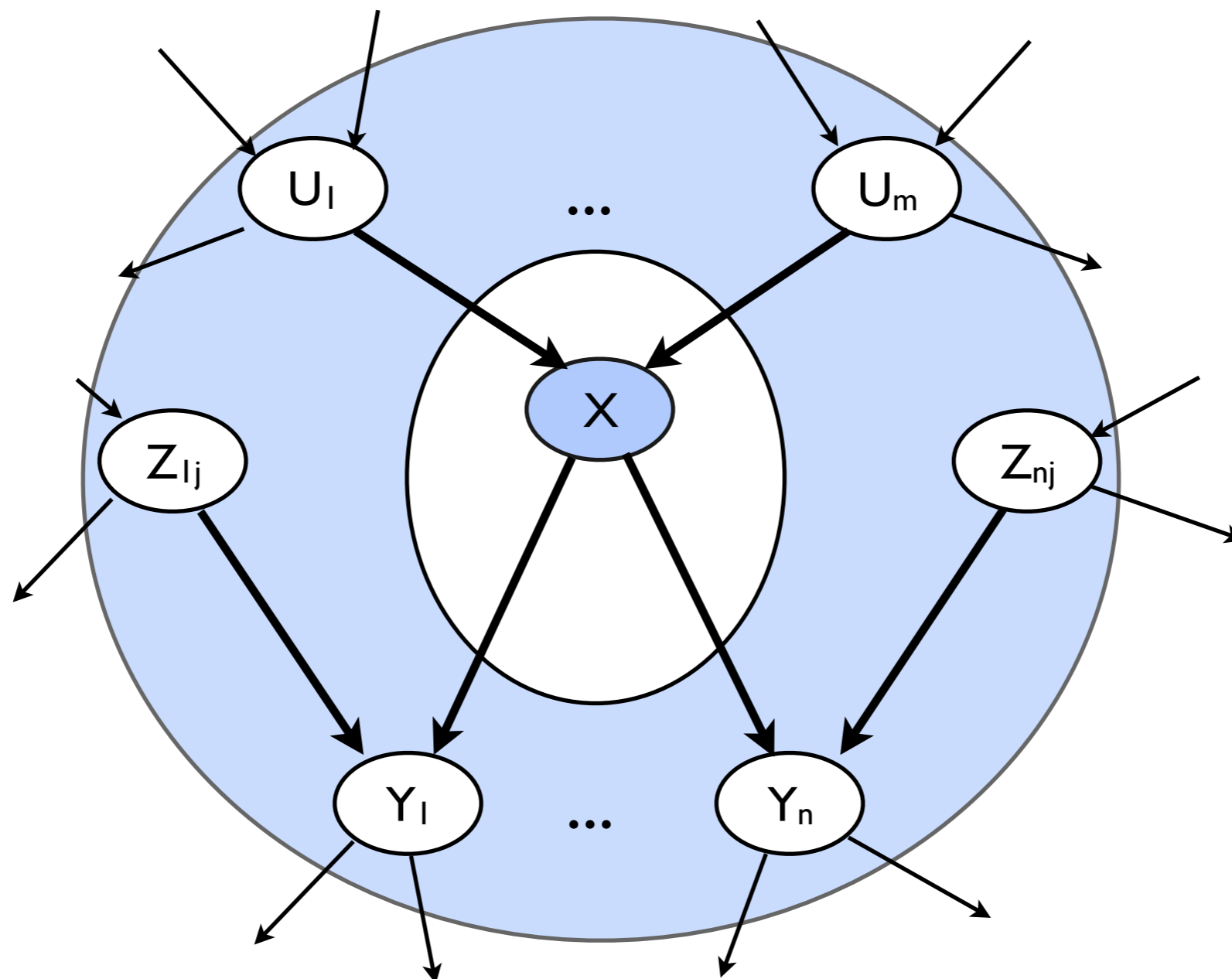
X

# 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 & probability (chapter 13)

  - Uncertainty

  - Probability

  - Syntax and Semantics

  - Inference

  - Independence and Bayes' Rule

- Bayesian Networks (chapter 14.1-3)

  - Syntax

  - Semantics

  - Efficient representation

X

# 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( Parents( X))* for some function *f*

E.g., Boolean functions

*NorthAmerican ⇔ Canadian ∨ US ∨ Mexican*

E.g., numerical relationships among continuous variables

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

X

# Compact conditional distributions (2)

*Noisy-OR* distributions model multiple noninteracting causes

1) Parents $U_1 \ldots U_k$ include all causes ( add *leak node* for "miscellaneous" ones)

2) Independent failure probability $q_i$ for each cause alone

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

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

Number of parameters linear in number of parents

X

# 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)