

Probabilistic Inference and Belief Networks

Uncertain Knowledge and Reasoning

[related reading: chapter 8 of the AI textbook]

Probabilities everywhere

- Not just for games of chance!
 - I'm snuffling: am I sick?
 - Email contains "FREE!": is it spam?
 - Tooth hurts: have cavity?
 - Safe to cross street?
 - 60 min enough to get to the airport?
 - Robot rotated wheel three times, how far did it advance?
- Reasons for uncertainty and randomness:
 - Theoretical and modeling limitations
 - Coin toss example: we may not have a complete model for physics of the environment (e.g. molecular structure of the coin, micro air movements, ...)
 - Sensory and measurement limitations
 - Coin toss example: we have a physical model that requires as input very precise measurements that are not available (e.g. how much energy exactly the coin received, the current state of the environment).
 - Computational limitations
 - Coin toss example: assuming that the model and input data are available, calculations might be too time consuming.

Random Variables

- A random variable is some aspect of the world about which we have uncertainty
 - R = Is it raining?
 - D = How long will it take to drive to work?
 - L = Where am I?
- We denote random variables with capital letters. For their values we use lower case.
- Each random variable has a domain
 - R in $\{\text{true}, \text{false}\}$
 - D in $[0, \infty]$
 - L in possible locations

Probability distributions

- Unobserved random variables have distributions

$P(T)$

T	P
warm	0.5
cold	0.5

$P(W)$

W	P
sun	0.6
rain	0.1
fog	0.3

- A distribution is a TABLE of probabilities of values
- A probability (lower case value) is a single number

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

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

- Must have: $\forall x P(x) \geq 0$

$$\sum_x P(x) = 1$$

Joint distributions

- A *joint distribution* over a set of random variables: X_1, X_2, \dots, X_n is a map from assignments (or *outcomes*, or *atomic events*) to reals:

$$P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n)$$

$$P(x_1, x_2, \dots, x_n)$$

- Size of distribution if n variables with domain sizes d ?

- Must obey: $0 \leq P(x_1, x_2, \dots, x_n) \leq 1$

$$\sum_{(x_1, x_2, \dots, x_n)} P(x_1, x_2, \dots, x_n) = 1$$

T	S	P
warm	sun	0.4
warm	rain	0.1
cold	sun	0.2
cold	rain	0.3

- For all but the smallest distributions, impractical to write out

Events

- An *event* is a set E of assignments (or outcomes)

$$P(E) = \sum_{(x_1 \dots x_n) \in E} P(x_1 \dots x_n)$$

- From a joint distribution, we can calculate the probability of any event
- Probability that it's warm AND sunny?
- Probability that it's warm?
- Probability that it's warm OR sunny?

T	S	P
warm	sun	0.4
warm	rain	0.1
cold	sun	0.2
cold	rain	0.3

Marginalization

- Marginalization (or summing out) is *projecting* a joint distribution to a sub-distribution over subset of variables

$$P(X_1 = x_1) = \sum_{x_2} P(X_1 = x_1, X_2 = x_2)$$

$P(T, S)$				$P(T)$	
T	S	P		T	P
warm	sun	0.4	$\xrightarrow{P(t) = \sum_s P(t, s)}$	warm	0.5
warm	rain	0.1		cold	0.5
cold	sun	0.2	$\xrightarrow{P(s) = \sum_t P(t, s)}$	$P(S)$	
cold	rain	0.3		S	P
				sun	0.6
				rain	0.4

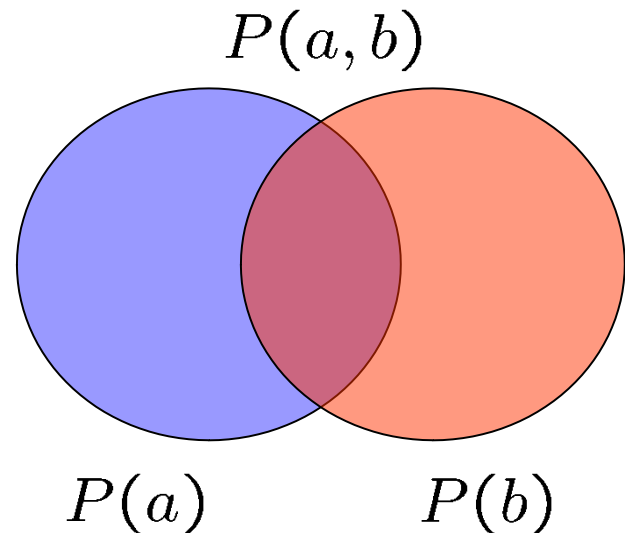
Conditional Probabilities

- A conditional probability is the probability of an event given another event (usually called evidence)

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

$P(T, S)$

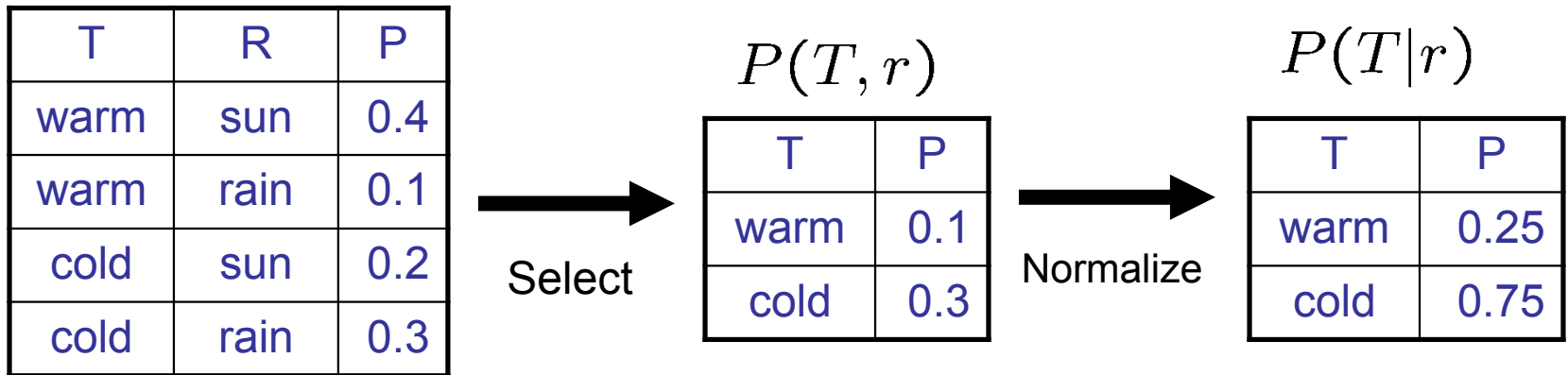
T	S	P
warm	sun	0.4
warm	rain	0.1
cold	sun	0.2
cold	rain	0.3



$P(\text{rain} \mid \text{cold}) = ???$

Normalization Trick

- A trick to get the whole conditional distribution at once:
 - Select the joint probabilities matching the evidence
 - Normalize the selection (make it sum to one)
- Example: find $P(T|\text{rain})$



$$P(x_1|x_2) = \frac{P(x_1, x_2)}{P(x_2)} = \frac{P(x_1, x_2)}{\sum_{x_1} P(x_1, x_2)}$$

The Product Rule

- Sometimes joint $P(X,Y)$ is easy to get
- Sometimes easier to get conditional $P(X|Y)$

$$P(x|y) = \frac{P(x,y)}{P(y)} \quad \longleftrightarrow \quad P(x,y) = P(x|y)P(y)$$

The Chain Rule

- More generally, can always write any joint distribution as an incremental product of conditional distributions

$$P(x_1, x_2, x_3) = P(x_1)P(x_2|x_1)P(x_3|x_1, x_2)$$

$$P(x_1, x_2, \dots x_n) = \prod_i P(x_i|x_1 \dots x_{i-1})$$

Probabilistic Inference

- Probabilistic inference: compute a desired probability from other known probabilities (e.g. conditional from joint)
- We generally compute conditional probabilities
 - $P(\text{on time} \mid \text{no reported accidents}) = 0.90$
 - These represent the agent's *beliefs* given the evidence
- Probabilities change with new evidence:
 - $P(\text{on time} \mid \text{no accidents, 5 a.m.}) = 0.95$
 - $P(\text{on time} \mid \text{no accidents, 5 a.m., raining}) = 0.80$
 - Observing new evidence causes *beliefs to be updated*

Simple inference with Bayes' Rule

- A simple case of inference with two random variables where the joint distribution is written as a product:

$$P(x, y) = P(x|y)P(y) = P(y|x)P(x)$$

- Dividing by the marginal, we get Bayes' rule:

$$P(x|y) = \frac{P(y|x)P(x)}{P(y)}$$



- Why is this at all helpful?
 - Lets us invert a conditional distribution
 - Often one conditional is tricky but the other simple
 - Foundation of many systems

[More about these points when we see Bayes nets and Machine Learning]

- One of the most commonly-used equations in AI

Example: Inference with Bayes' rule

- Example: Diagnostic probability from causal probability:

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

- Suppose

- m means meningitis
 - s means stiff neck

$$\left. \begin{array}{l} P(s|m) = 0.8 \\ P(m) = 0.0001 \\ P(s) = 0.1 \end{array} \right\} \text{Example gives}$$

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

- Note: The posterior probability of meningitis is very small. However, you should still get stiff necks checked out. Why?
 - A low probability multiplied by a large cost can still create a significant risk!
- Note: The given probabilities implicitly define a joint distribution:
 - $p(s,m) = p(s|m)p(m)$
 - $p(\neg s,m) = p(\neg s|m)p(m) = (1 - p(s|m)) p(m)$
 - $p(s,\neg m) = p(s) - p(s,m)$
 - $p(\neg s,\neg m) = 1 - \text{sum of above entries}$

Inference by Enumeration

- A more general procedure. We want:

$$P(Y_1 \dots Y_m | e_1 \dots e_k)$$

- Evidence variables: $(E_1 \dots E_k) = (e_1 \dots e_k)$
 - Query variables: $Y_1 \dots Y_m$
 - Hidden variables: $H_1 \dots H_r$
- $\left. \begin{array}{l} \\ \\ \end{array} \right\} \begin{array}{l} X_1, X_2, \dots, X_n \\ \text{All variables} \end{array}$

- First, select the entries consistent with the evidence
- Second, sum out H:

$$P(Y_1 \dots Y_m, e_1 \dots e_k) = \sum_{h_1 \dots h_r} P(\underbrace{Y_1 \dots Y_m, h_1 \dots h_r, e_1 \dots e_k}_{X_1, X_2, \dots, X_n})$$

- Finally, normalize the remaining entries.
- Obvious problems:
 - Worst-case time complexity $O(d^n)$
 - Space complexity $O(d^n)$ to store the joint distribution

Inference by Enumeration

- $P(R)$?

- $P(\text{sun} \mid \text{winter})$?

or equivalently $P(R=\text{sun} \mid \text{winter})$?

S	T	R	P
summer	warm	sun	0.30
summer	warm	rain	0.05
summer	cold	sun	0.10
summer	cold	rain	0.05
winter	warm	sun	0.10
winter	warm	rain	0.05
winter	cold	sun	0.15
winter	cold	rain	0.20

- $P(R \mid \text{winter}, \text{warm})$?

Probabilistic Models

- Models describe how (a portion of) the world works
- **Models are always simplifications**
 - May not account for every variable
 - May not account for all interactions between variables
 - “All models are wrong; but some are useful.”
 - George E. P. Box
- A joint distribution is a probabilistic model.
- What do we do with probabilistic models?
 - We (or our agents) need to reason about unknown variables, given evidence
 - Example: explanation (diagnostic reasoning)
 - Example: prediction (causal reasoning)
 - Example: value of information

Complexity of Models

- Engineers and designers are interested in *simple* and *compact* models (as long as the model is sufficiently good)
 - Simple models are *easier to build*
 - Simple models are *easier to explain* (e.g. why/how they work)
 - Compact models take *less space*
 - Simplicity/Compactness usually implies more efficient computation (*lower time complexity*)
- One way of measuring the complexity of a probabilistic model is to count the number of (free) parameters (values) that must be specified.

A joint distribution over n variables whose domain sizes are d (each can take d distinct values) requires d^n entries in the table.

The number of (*free*) *parameters* is $d^n - 1$. [Because once you specify $d^n - 1$ of the entries, the last one is $(1 - \text{sum of those specified})$ because the entries should add up to 1.]

Complexity of Models (cont'd)

If a probabilistic model has multiple tables (distributions), the number of its free parameters is the sum of the number of free parameters of the tables.

Question: How many free parameters would we need if instead of a full joint distribution, we used multiple smaller distributions using the chain rule?

$$P(x_1, x_2, \dots, x_n) = \prod_i P(x_i | x_1 \dots x_{i-1})$$

Concretely, counting the number of free parameters accounting for that we know probabilities sum to one:

$$\begin{aligned} & (d-1) + d(d-1) + d^2(d-1) + \dots + d^{n-1}(d-1) \\ &= (d^n - 1)/(d - 1) (d - 1) \\ &= d^n - 1 \end{aligned}$$

It doesn't make a difference. (i.e. using the chain rule alone doesn't reduce complexity.)

Independence

- Two variables are *independent* if:

$$P(X, Y) = P(X)P(Y)$$

$$X \perp\!\!\!\perp Y$$

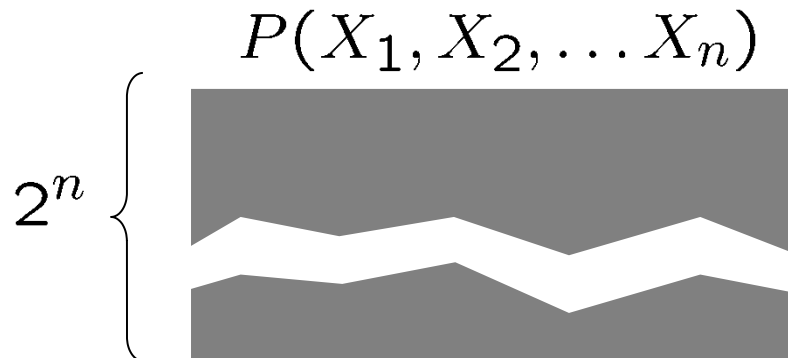
$$\forall x, y \ P(x, y) = P(x)P(y)$$

- This says that their joint distribution *factors* into a product of two simpler distributions
- We can use independence as a *modeling assumption*
 - Independence can be a simplifying assumption
 - What could we assume for {Weather, Traffic, Meningitis}?
- How many parameters in the joint model using one table?
- How many parameters when assuming independence?

Example: Independence

- N fair, independent coin flips:

$P(X_1)$		$P(X_2)$		\dots		$P(X_n)$	
H	0.5	H	0.5			H	0.5
T	0.5	T	0.5			T	0.5



Example: Independence?

- Arbitrary joint distributions can be poorly modeled by independent factors

$P(T)$

T	P
warm	0.5
cold	0.5

$P(S)$

S	P
sun	0.6
rain	0.4

$P(T, S)$

T	S	P
warm	sun	0.4
warm	rain	0.1
cold	sun	0.2
cold	rain	0.3

$P(T)P(S)$

T	S	P
warm	sun	0.3
warm	rain	0.2
cold	sun	0.3
cold	rain	0.2

Conditional Independence

- Unconditional (absolute) independence is very rare.
- Conditional independence is our most basic and robust form of knowledge about uncertain environments:

$$\forall x, y, z : P(x, y|z) = P(x|z)P(y|z)$$

$$\forall x, y, z : P(x|z, y) = P(x|z)$$

$$X \perp\!\!\!\perp Y | Z$$

Conditional Independence: Example

- If I have a cavity, the probability that the probe catches in it doesn't depend on whether I have a toothache:
 - $P(\text{catch} \mid \text{toothache}, \text{cavity}) = P(\text{catch} \mid \text{cavity})$
- The same independence holds if I don't have a cavity:
 - $P(\text{catch} \mid \text{toothache}, \neg \text{cavity}) = P(\text{catch} \mid \neg \text{cavity})$
- Catch is *conditionally independent* of Toothache given Cavity:
 - $P(\text{Catch} \mid \text{Toothache}, \text{Cavity}) = P(\text{Catch} \mid \text{Cavity})$
- Equivalent statements:
 - $P(\text{Toothache} \mid \text{Catch}, \text{Cavity}) = P(\text{Toothache} \mid \text{Cavity})$
 - $P(\text{Toothache}, \text{Catch} \mid \text{Cavity}) = P(\text{Toothache} \mid \text{Cavity}) P(\text{Catch} \mid \text{Cavity})$
 - One can be derived from the other easily

Chain Rule and Conditional Independence

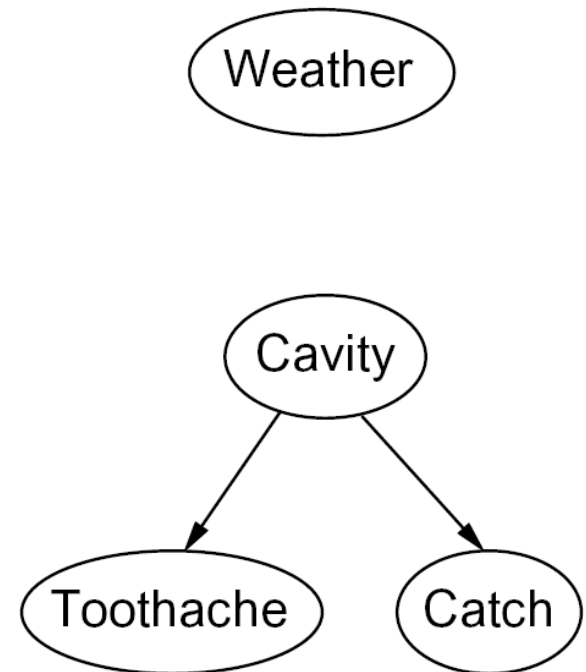
- How many entries/parameters do we need for the joint distribution $P(\text{Toothache}, \text{Cavity}, \text{Catch})$?
7 independent entries ($2^3 - 1$)
- Write out full joint distribution using chain rule:
 - $P(\text{Toothache}, \text{Catch}, \text{Cavity})$
= $P(\text{Toothache} \mid \text{Catch}, \text{Cavity}) P(\text{Catch}, \text{Cavity})$
= $P(\text{Toothache} \mid \text{Catch}, \text{Cavity}) P(\text{Catch} \mid \text{Cavity}) P(\text{Cavity})$
= $P(\text{Toothache} \mid \text{Cavity}) P(\text{Catch} \mid \text{Cavity}) P(\text{Cavity})$
- How many (free) parameters does it need now?
 $2 + 2 + 1 = 5$
(i.e. by assuming conditional independence, the complexity is reduced.)

Belief Networks: Big Picture

- Two problems with using full joint distribution tables for probabilistic models:
 - Unless there are only a few variables, the joint is WAY too big to represent explicitly
 - Hard to learn (estimate) anything empirically about more than a few variables at a time
- **Belief nets:** a technique for describing complex joint distributions (models) using simple, local distributions (conditional probabilities)
 - More properly called **graphical models**
 - We describe how variables locally interact
 - Local interactions chain together to give global, indirect interactions
 - Also known as **Bayes' nets** or **Bayesian networks**

Graphical Model Notation

- **Nodes: variables (with domains)**
 - Can be assigned (observed) or unassigned (unobserved)
- **Arcs: interactions**
 - Indicate “direct influence” between variables
- For now: imagine that arrows mean causation (in general, they don't have to)



Example: Coin Flips

- N independent coin flips

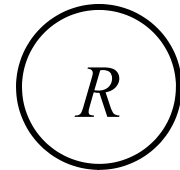


- No interactions between variables: **absolute independence**

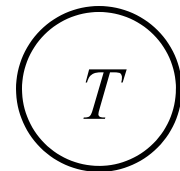
Example: Traffic

- Variables:

- R : It rains
- T : There is traffic

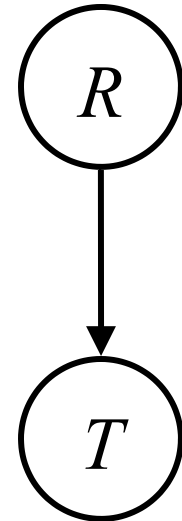


- Model 1: independence



Example: Traffic

- Variables:
 - R: It rains
 - T: There is traffic
- Model 2: rain causes traffic

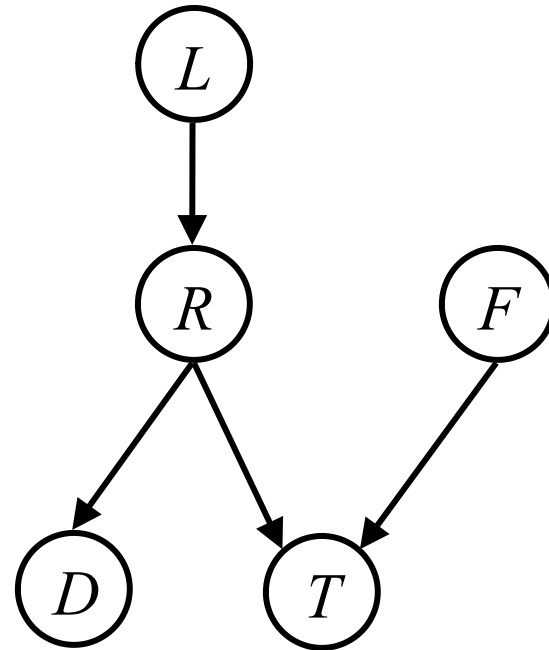


- An agent using model 2 may perform better.

Example: Traffic II

■ Variables

- T: Traffic
- R: It rains
- L: Low pressure
- D: Roof drips
- F: Festival



Example: Alarm network

You have a new burglar alarm installed at home. It is fairly reliable at detecting a burglary, but also responds on occasion to minor earthquakes.

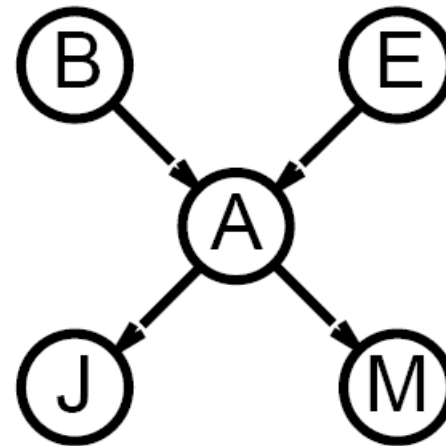
You also have two neighbors, John and Mary, who have promised to call you at work when they hear the alarm. John nearly always calls when he hears the alarm, but sometimes confuses the telephone ringing with the alarm and calls then, too. Mary, on the other hand, likes rather loud music and often misses the alarm altogether.

- Example inference task in this domain:
 - Given the evidence of who has or has not called, we would like to estimate the probability of a burglary.
 - We would like to know the probability of the alarm going off.

Example: Alarm network

■ Variables

- B: Burglary
- A: Alarm goes off
- M: Mary calls
- J: John calls
- E: Earthquake

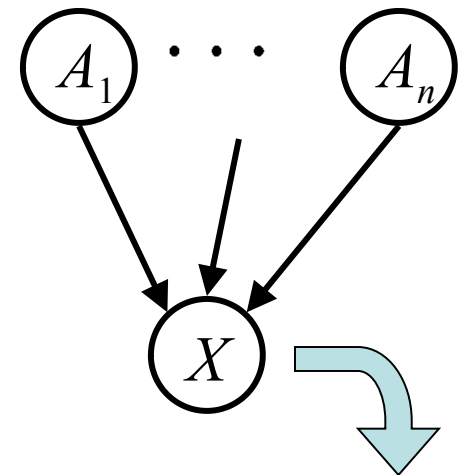


Belief Net Semantics

A belief network is:

- A set of nodes, one per random variable
- A directed, acyclic graph
- A collection of distributions (CPTs) over each node, one for each combination of parents' values

$$P(X|a_1 \dots a_n)$$



$$P(X|A_1 \dots A_n)$$

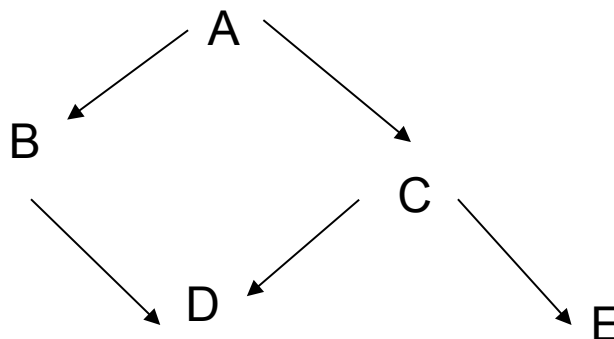
- CPT: conditional probability table
- Description of a noisy “causal” process

A belief net = Topology (graph) + Local Conditional Probabilities

Topological semantics

- A node is **conditionally independent** of its **non-descendants** given its **parents**
- A node is **conditionally independent** of all other nodes in the network given its parents, children, and children's parents (also known as its **Markov blanket**)
- The method called **d-separation** can be applied to decide whether a set of nodes X is independent of another set Y , given a third set Z

Joint Probabilities in BNs: Starting with an example



Computing the joint probability for all variables is easy:

$$\begin{aligned} P(a, b, c, d, e) &= P(e \mid a, b, c, d) P(a, b, c, d) && \text{by the product rule} \\ &= P(e \mid c) P(a, b, c, d) && \text{by cond. indep. assumption} \\ &= P(e \mid c) P(d \mid a, b, c) P(a, b, c) \\ &= P(e \mid c) P(d \mid b, c) P(c \mid a, b) P(a, b) \\ &= P(e \mid c) P(d \mid b, c) P(c \mid a) P(b \mid a) P(a) \end{aligned}$$

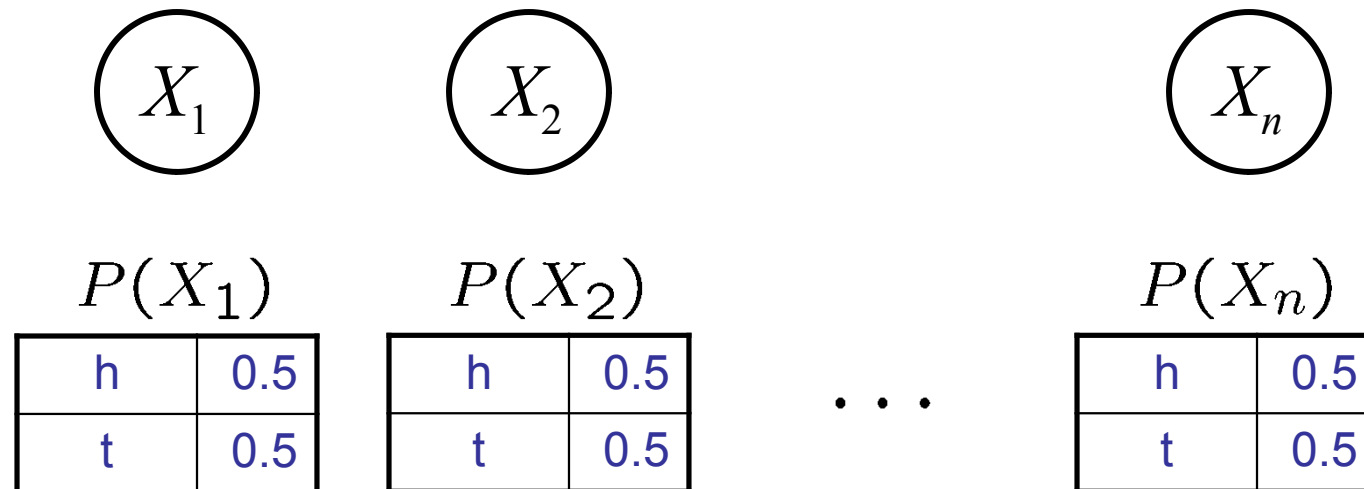
BNs implicitly encode joint distributions

- Belief nets **implicitly** encode joint distributions
 - As a product of local conditional distributions
 - To see what probability a BN gives to a full assignment, multiply all the relevant conditionals together:

$$P(x_1, x_2, \dots, x_n) = \prod_{i=1}^n P(x_i | \text{parents}(X_i))$$

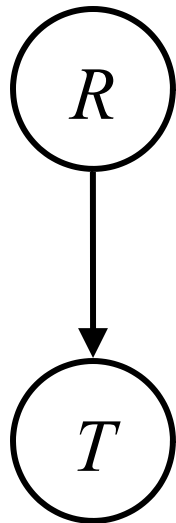
- This lets us reconstruct any entry of the full joint
- Not every BN can represent every full joint
 - The topology enforces certain conditional independencies
- By having the joint distribution, we can answer any query (e.g. by using inference by enumeration)
 - Note: more advanced inference algorithms do not require constructing the entire joint.

Example: Coin Flips



$$P(h, t, h, h, \dots, t) = ???$$

Example: Traffic



$P(R)$

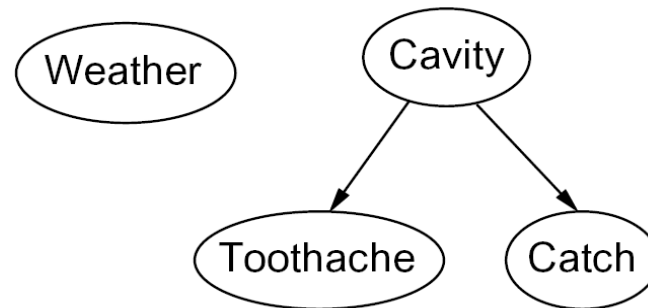
r	$1/4$
$\neg r$	$3/4$

$$P(r, \neg t) =$$

$P(T|R)$

r	t	$3/4$
	$\neg t$	$1/4$
$\neg r$	t	$1/2$
	$\neg t$	$1/2$

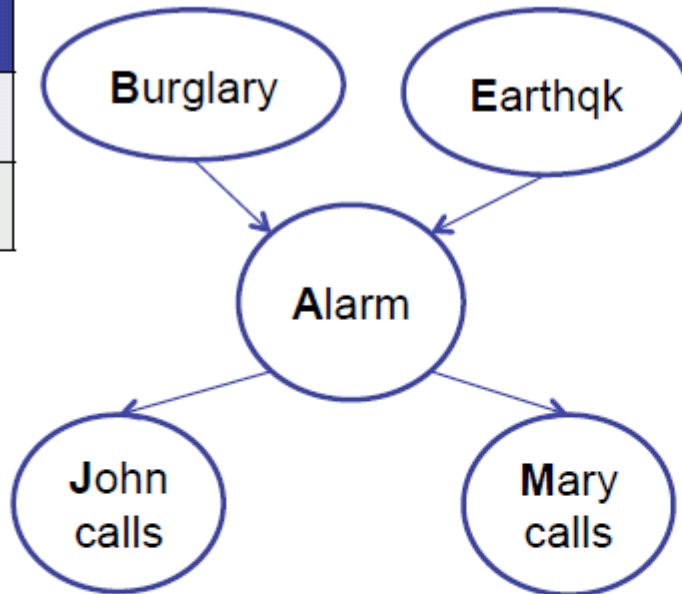
Example: Dental Health



$$P(\neg \text{cavity}, \text{catch}, \neg \text{toothache}, \text{weather}) = \\ P(\text{weather}) P(\neg \text{cavity}) P(\neg \text{toothache} | \neg \text{cavity}) P(\text{catch} | \neg \text{cavity})$$

Example: Alarm Network

B	P(B)
+b	0.001
¬b	0.999



E	P(E)
+e	0.002
¬e	0.998

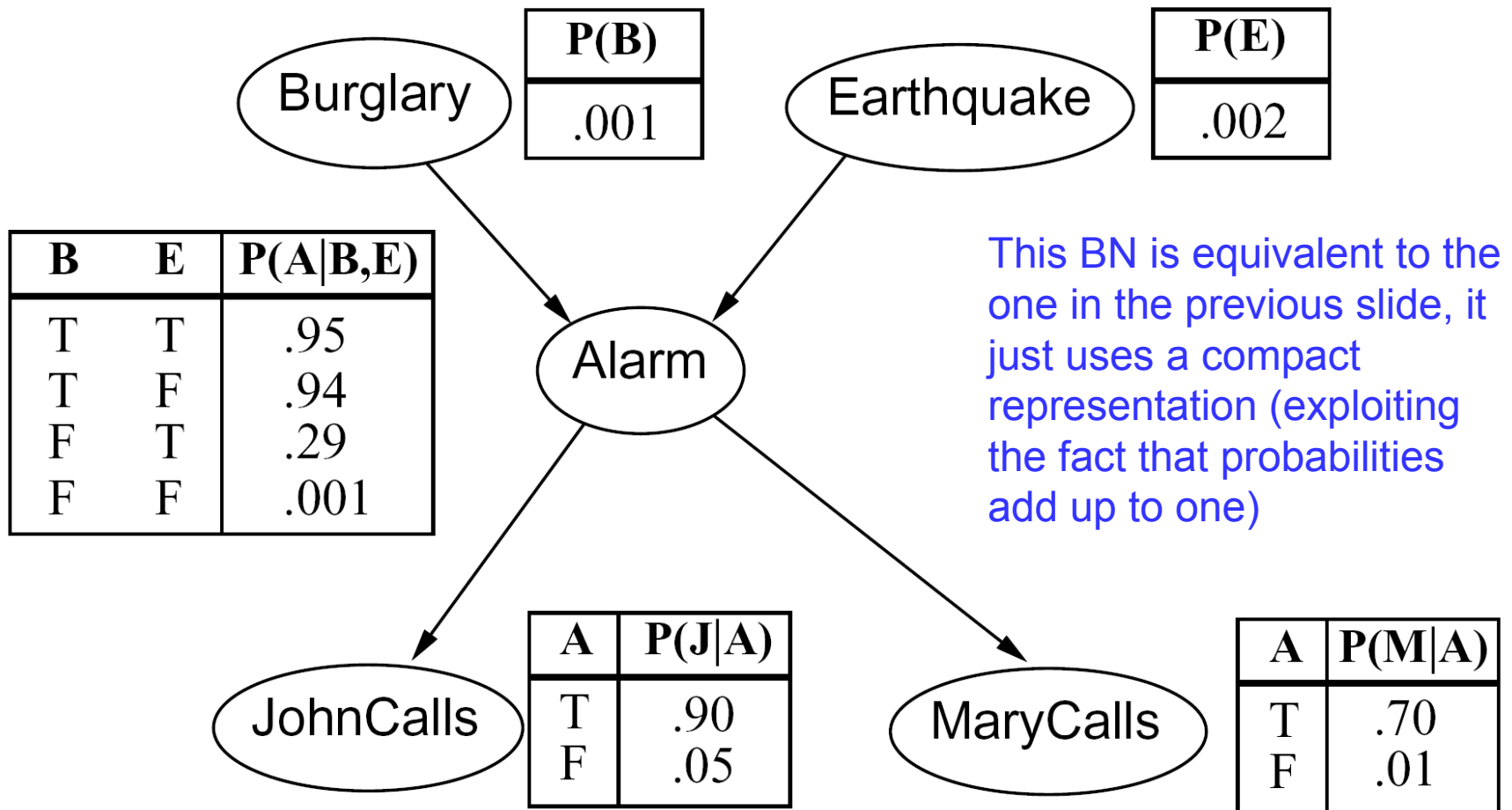
A	J	P(J A)
+a	+j	0.9
+a	¬j	0.1
¬a	+j	0.05
¬a	¬j	0.95

A	M	P(M A)
+a	+m	0.7
+a	¬m	0.3
¬a	+m	0.01
¬a	¬m	0.99

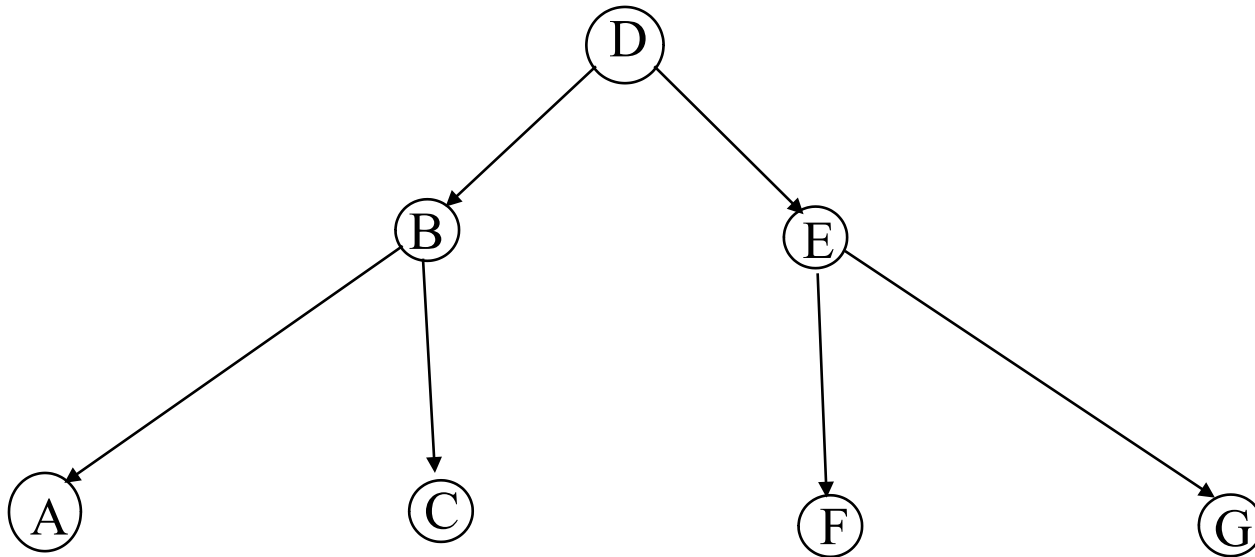
B	E	A	P(A B,E)
+b	+e	+a	0.95
+b	+e	¬a	0.05
+b	¬e	+a	0.94
+b	¬e	¬a	0.06
¬b	+e	+a	0.29
¬b	+e	¬a	0.71
¬b	¬e	+a	0.001
¬b	¬e	¬a	0.999

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

Example: Alarm Network

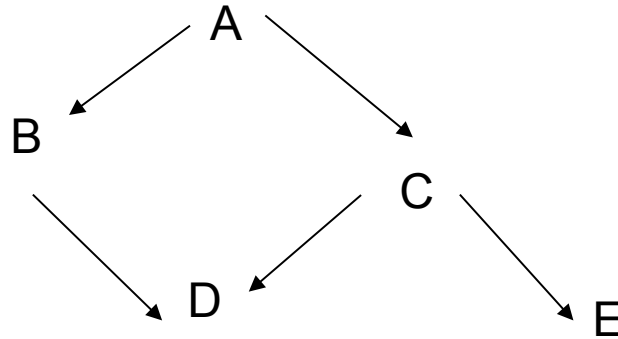


Example: Tree-Structured Bayesian Network



$P(A, B, C, D, E, F, G)$ is modeled as $P(A|B) P(C|B) P(F|E) P(G|E) P(B|D) P(E|D) P(D)$

Computing probabilities using BN



- Any complete joint can be computed:
 - $P(a,b,-c,d,-e) = P(a)P(b|a)P(-c|a)P(d|b,-c)P(-e|-c)$
- Some probabilities are directly available in the CPTs. No calculation is needed:
 - $P(a)$
 - $P(b|-a)$
 - $P(d|b,-c)$
 - $P(-e|c)$
- For other cases, use inference by enumeration.

Inference by enumeration on BNs

- Want to compute $P(Y \mid \mathbf{e})$
 - Y : The query variable
 - \mathbf{e} : Evidence (observed) variables
- Add all of the terms (atomic event probabilities) from the full joint distribution
- If \mathbf{H} are the other unobserved (hidden) variables, then:
$$P(Y|\mathbf{e}) = \alpha P(Y, \mathbf{e}) = \alpha \sum_{\mathbf{H}} P(Y, \mathbf{e}, \mathbf{H})$$
- Y is a variable. We have to compute the above for every value in the domain of Y .
- Each $P(Y, \mathbf{e}, \mathbf{H})$ term can be obtained from the network.
- Computationally expensive!
 - There are more efficient algorithms but this is enough for the scope of COSC367. See chapter 8 for more efficient algorithms and how to compute the joint using *factors*.

Inference: Example

■ Inference by enumeration

$$P(B \mid j, m) = P(B, j, m) / P(j, m)$$

$$= \alpha P(B, j, m)$$

$$= \alpha \sum_A \sum_E P(B, E, A, j, m)$$

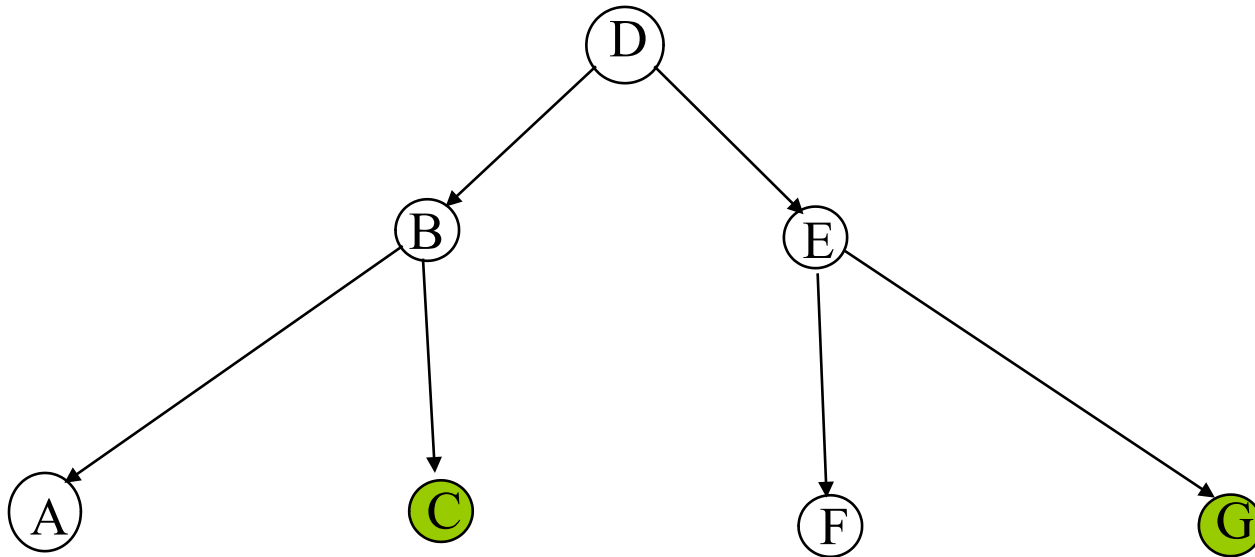
$$= \alpha \sum_A \sum_E P(B) P(E) P(A \mid B, E) P(j \mid A) P(m \mid A)$$

Here, hidden variables are A and E.

Compute it for $B=+b$ and $B=-b$, you will get 0.00059224 and 0.0014919. Find α such that they add up to one:

$$P(B \mid j, m) = \alpha (0.00059224, 0.0014919) \approx (0.284, 0.716) .$$

Example: Using BN and enumeration



$$P(a|c,g) = \alpha \sum_{BDEF} P(a,B,D,E,F,c,g)$$

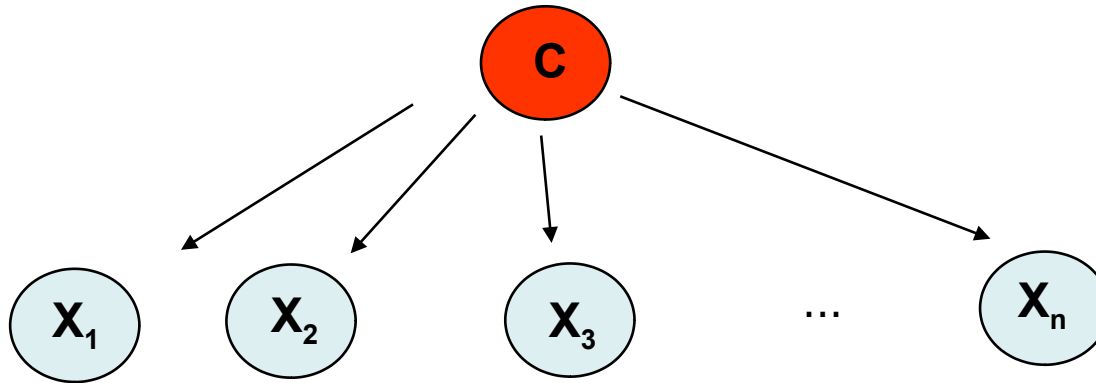
$$\alpha = 1 / \left(\sum_{BDEF} P(a,B,D,E,F,c,g) + \sum_{BDEF} P(-a,B,D,E,F,c,g) \right)$$

Note that
inside the brackets
equals to $p(c,g)$

Inference: Other cases

- How to answer $P(Y)$?
 - There is no evidence therefore all variables except the query variable are hidden (must be summed over)
- How to answer $P(y|\mathbf{e})$?
 - First answer $P(Y|\mathbf{e})$, then pick the result for $Y=y$
- How to answer $P(Y_1=y_1, Y_2=y_2|\mathbf{e})$?
 - It is $P(y_1|y_2, \mathbf{e})P(y_2|\mathbf{e})$

BNs Example: Naïve Bayes Models



$$P(C | X_1, \dots, X_n) = \alpha \prod P(X_i | C) P(C) = \alpha P(X_1 | C) P(X_2 | C) \dots P(X_n | C) P(C)$$

Features X are conditionally independent given the **class** variable C

$P(C)$: Prior distribution of C , the class random variable

$P(X_i | C)$: Likelihood conditional distributions

$P(C | X_1, \dots, X_n)$: Posterior distribution

Widely used in machine learning

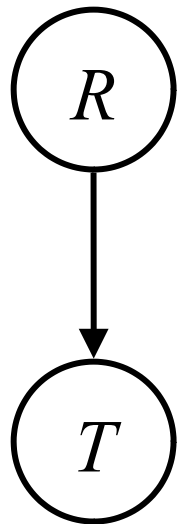
e.g., spam email classification: X 's = counts of words in emails

Conditional probabilities $P(X_i | C)$ can easily be estimated from labeled data.

Further Discussion

Reverse causality?

- Consider the basic traffic net
- Let's multiply out the joint



$P(R)$

r	$1/4$
$\neg r$	$3/4$

r	t	$3/4$
	$\neg t$	$1/4$

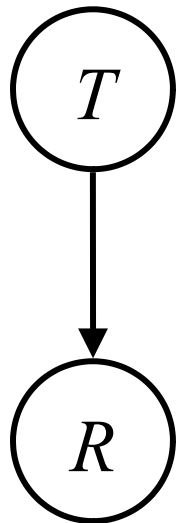
$\neg r$	t	$1/2$
	$\neg t$	$1/2$

$P(T, R)$

r	t	$3/16$
r	$\neg t$	$1/16$
$\neg r$	t	$6/16$
$\neg r$	$\neg t$	$6/16$

Reverse causality? (cont'd)

- Can we express the same joint distribution by a network where arrows no longer mean causality?
 - Yes, but the network is now harder for humans to construct and understand.



t	$9/16$
$\neg t$	$7/16$

$P(R|T)$

t	r	$1/3$
	$\neg r$	$2/3$
$\neg t$	r	$1/7$
	$\neg r$	$6/7$

r	t	$3/16$
r	$\neg t$	$1/16$
$\neg r$	t	$6/16$
$\neg r$	$\neg t$	$6/16$

Causality?

- When Belief nets reflect the true causal patterns:
 - Often simpler (nodes have fewer parents)
 - Often easier to think about
 - Often easier to elicit from experts
- BNs need not actually be causal
 - Sometimes no causal net exists over the domain (especially if variables are missing)
 - E.g. consider the variables *Traffic* and *Drips*
 - End up with arrows that reflect correlation, not causation
- What do the arrows really mean?
 - Topology may happen to encode causal structure
 - **Topology only guaranteed to encode conditional independencies**

Constructing belief networks

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

Example: non-causal modeling (not desired in COSC367)

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

■

MaryCalls

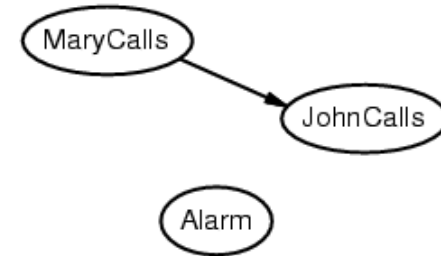
JohnCalls

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

Example

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

■



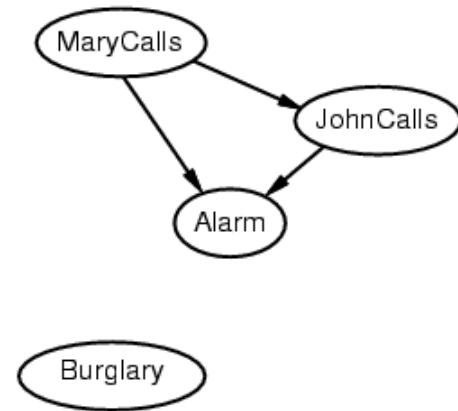
$P(J \mid M) = P(J)$ **No**

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

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**

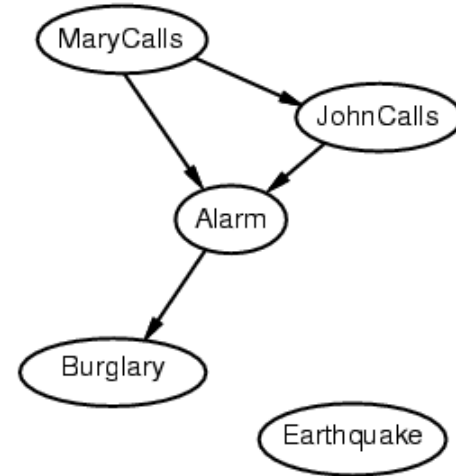
$P(B \mid A, J, M) = P(B \mid A)$?

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

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**

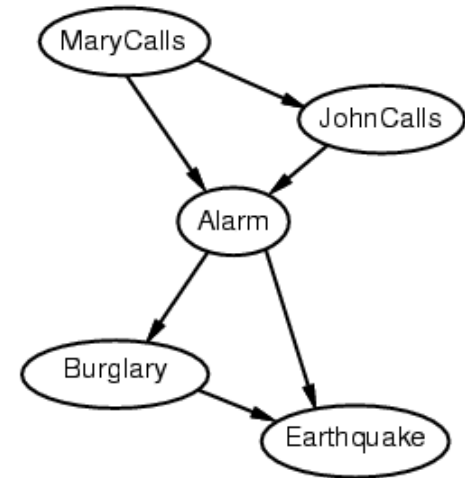
$P(E \mid B, A, J, M) = P(E \mid A)$?

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

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**

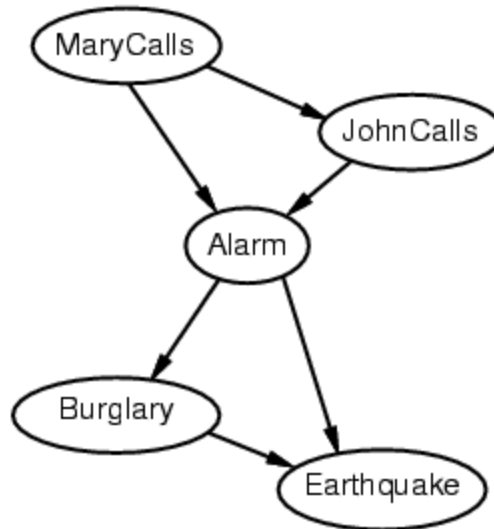
$P(B \mid A, J, M) = P(B \mid A)$? **Yes**

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

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

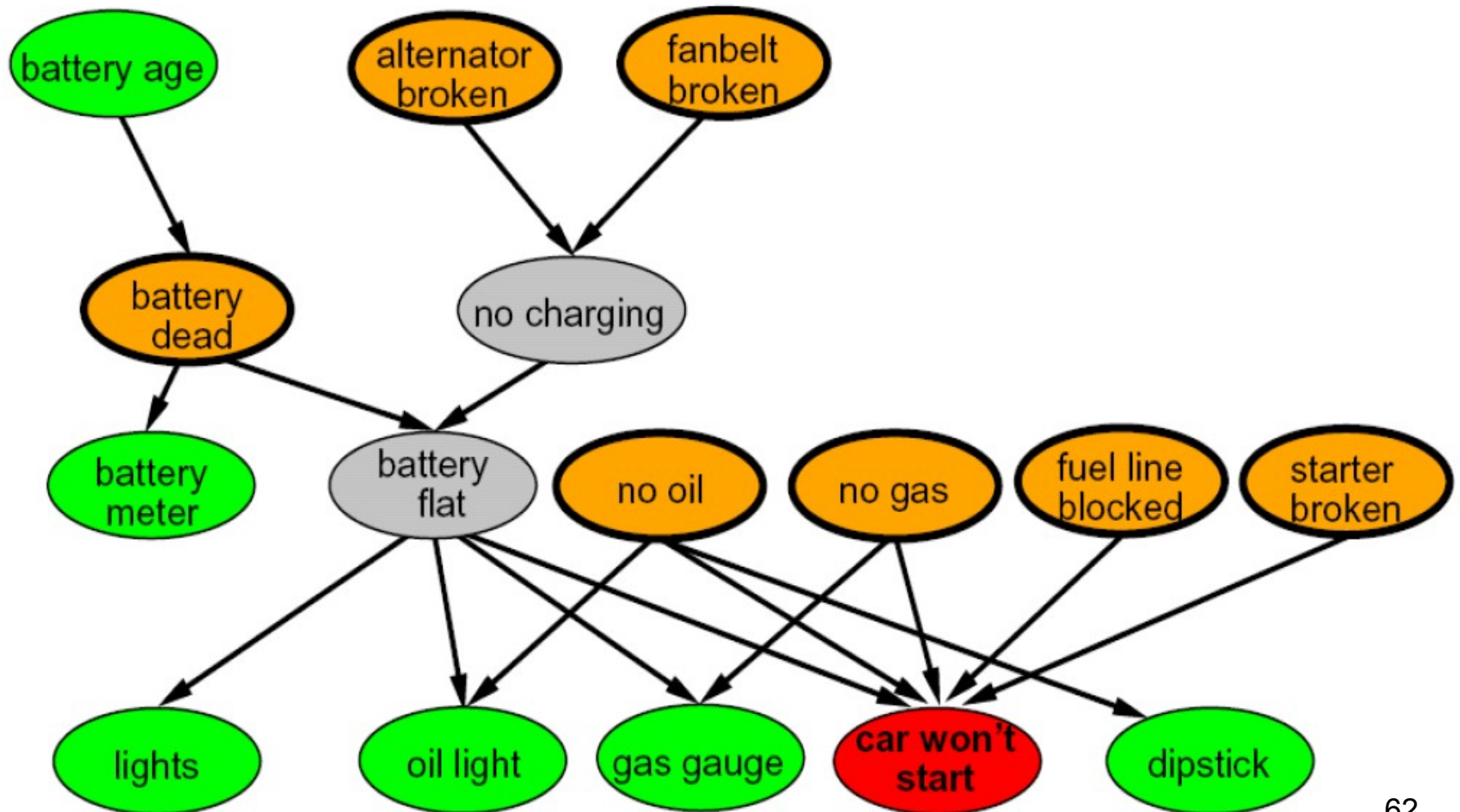
$P(E \mid B, A, J, M) = P(E \mid A, B)$? **Yes**

Example contd.



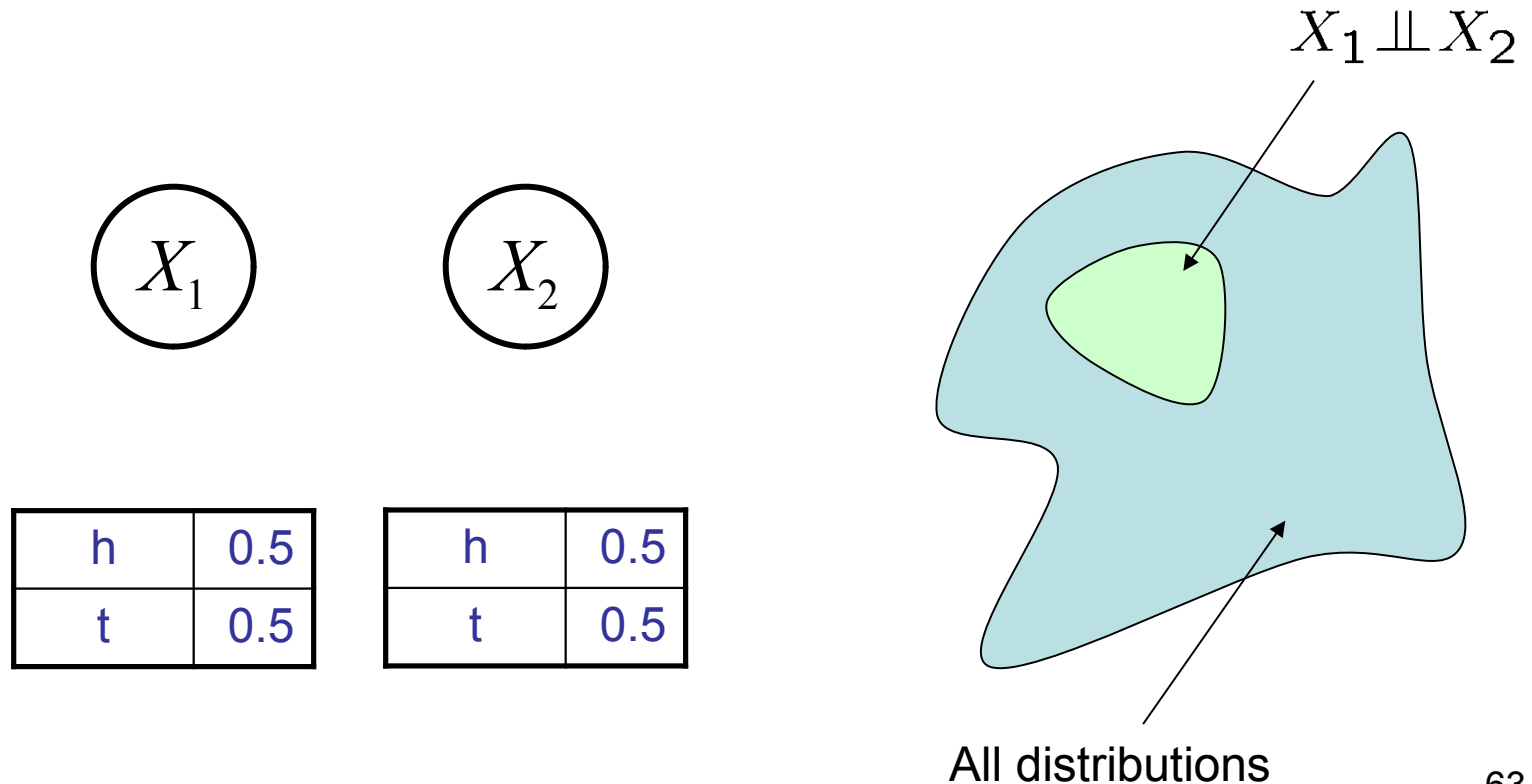
- Deciding conditional independence is hard in noncausal directions
- (Causal models and conditional independence seem hardwired for humans!)
- Network is less compact: $1 + 2 + 4 + 2 + 4 = 13$ numbers needed
- We normally prefer causal and compact networks (not the one in this example)

Example: Car Breakdown



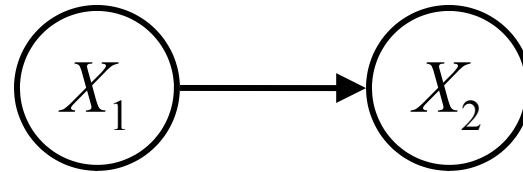
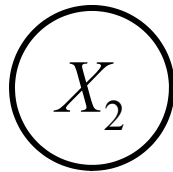
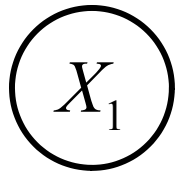
Example: Independence

- For this graph, you can fiddle with the CPTs all you want, but you won't be able to represent any distribution in which the flips are dependent!



Example: Coins

- Arcs don't prevent representing independence, they just allow dependence.



h	0.5
t	0.5

h	0.5
t	0.5

h	0.5
t	0.5

$P(X_2|X_1)$

h h	0.5
t h	0.5

h t	0.5
t t	0.5

Topology Limits Distributions

- Given some graph topology G , only certain joint distributions can be encoded
- The graph structure guarantees certain (conditional) independences
- Adding arcs increases the set of distributions, but has several costs
- Full conditioning can encode any distribution

