

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
- 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

Probabilistic Models

- A probabilistic model is a joint distribution over a set of variables

$$P(X_1, X_2, \dots, X_n)$$

- Inference: given a joint distribution, we can reason about unobserved variables given observations (evidence)
- General form of a query:

$$P(X_q | x_{e_1}, \dots, x_{e_k})$$

Stuff you care about *Stuff you already know*

- This conditional distribution is called a **posterior distribution** or the **belief function** of an agent which uses this model

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*

Conditional Probabilities

- *Conditional probabilities:*
 - E.g., $P(\text{cavity} \mid \text{toothache}) = 0.8$
 - Given that *toothache* is all I know...
- *Notation for conditional distributions:*
 - $P(\text{cavity} \mid \text{toothache})$ = a single number
 - $P(\text{Cavity}, \text{Toothache})$ = 2x2 table summing to 1
 - $P(\text{Cavity} \mid \text{Toothache})$ = Two 2-element distributions over Cavity, each summing to 1
- *If we know more:*
 - $P(\text{cavity} \mid \text{toothache}, \text{catch}) = 0.9$
 - $P(\text{cavity} \mid \text{toothache}, \text{cavity}) = 1$
- *Note: the less specific belief remains **valid** after more evidence arrives, but is not always **useful***
- *New evidence may be irrelevant, allowing simplification:*
 - $P(\text{cavity} \mid \text{toothache}, \text{traffic}) = P(\text{cavity} \mid \text{toothache}) = 0.8$
- *This kind of inference, guided by domain knowledge, is crucial*

The Product Rule

- Sometimes have conditional distributions but want the joint

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

- Example:

$P(W)$		$P(D W)$			$P(D, W)$		
R	P	D	W	P	D	W	P
sun	0.8	wet	sun	0.1	wet	sun	0.08
rain	0.2	dry	sun	0.9	dry	sun	0.72
		wet	rain	0.7	wet	rain	0.14
		dry	rain	0.3	dry	rain	0.06

The Chain Rule

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

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

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

- Why is this always true?

Bayes' Rule

- Two ways to factor a joint distribution over two variables:

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

That's my rule!

- Dividing, we get:

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



- Why is this at all helpful?
 - Lets us build one conditional from its reverse
 - Often one conditional is tricky but the other one is simple
 - Foundation of many systems we'll see later
- In the running for most important AI equation!

Inference with Bayes' Rule

- Example: Diagnostic probability from causal probability:

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

- Example:

- m is meningitis, s is stiff neck
- | | |
|-----------------|---------------------|
| $P(s m) = 0.8$ | } Example
givens |
| $P(m) = 0.0001$ | |
| $P(s) = 0.1$ | |

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

- Note: posterior probability of meningitis still very small
- Note: you should still get stiff necks checked out! Why?

Ghostbusters, Revisited

- Let's say we have two distributions:
 - Prior distribution** over ghost location: $P(G)$
 - Let's say this is uniform
 - Sensor reading model: $P(R | G)$
 - Given: we know what our sensors do
 - R = reading color measured at $(1,1)$
 - E.g. $P(R = \text{yellow} | G=(1,1)) = 0.1$
- We can calculate the **posterior distribution** $P(G|r)$ over ghost locations given a reading using Bayes' rule:

$$P(g|r) \propto P(r|g)P(g)$$

0.11	0.11	0.11
0.11	0.11	0.11
0.11	0.11	0.11

0.17	0.10	0.10
0.09	0.17	0.10
<0.01	0.09	0.17

Model for Ghostbusters

- Reminder: ghost is hidden, sensors are noisy

- T: Top sensor is red
B: Bottom sensor is red
G: Ghost is in the top

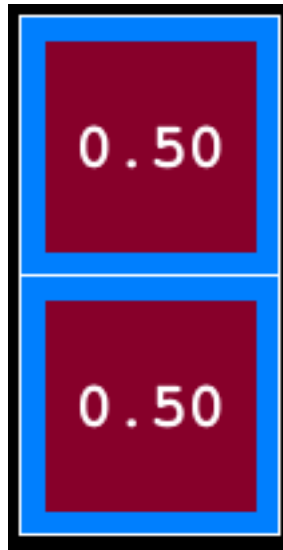
- Queries:

$$P(+g) = ??$$

$$P(+g \mid +t) = ??$$

$$P(+g \mid +t, -b) = ??$$

- Problem: joint distribution too large / complex



Joint Distribution

T	B	G	P(T,B,G)
+t	+b	+g	0.16
+t	+b	¬g	0.16
+t	¬b	+g	0.24
+t	¬b	¬g	0.04
¬t	+b	+g	0.04
¬t	+b	¬g	0.24
¬t	¬b	+g	0.06
¬t	¬b	¬g	0.06

Independence

- Two variables are *independent* if:

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

- This says that their joint distribution *factors* into a product two simpler distributions
- Another form:

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

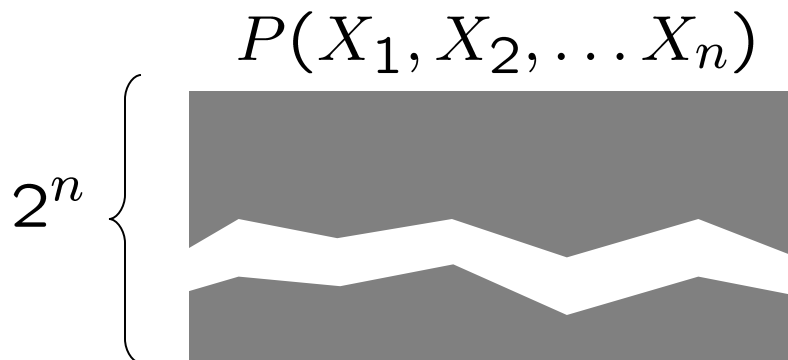
- We write: $X \perp\!\!\!\perp Y$

- Independence is a simplifying *modeling assumption*
 - Empirical* joint distributions: at best “close” to independent
 - What could we assume for {Weather, Traffic, Cavity, Toothache}?

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?

$P_1(T, W)$

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

$P(T)$

T	P
warm	0.5
cold	0.5

$P_2(T, W)$

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

$P(W)$

W	P
sun	0.6
rain	0.4

Conditional Independence

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

Conditional Independence

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

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

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

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

- What about this domain:
 - Traffic
 - Umbrella
 - Raining
- What about fire, smoke, alarm?

The Chain Rule

$$P(X_1, X_2, \dots, X_n) = P(X_1)P(X_2|X_1)P(X_3|X_1, X_2) \dots$$

- Trivial decomposition:

$$P(\text{Traffic}, \text{Rain}, \text{Umbrella}) =$$

$$P(\text{Rain})P(\text{Traffic}|\text{Rain})P(\text{Umbrella}|\text{Rain}, \text{Traffic})$$

- With assumption of conditional independence:

$$P(\text{Traffic}, \text{Rain}, \text{Umbrella}) =$$

$$P(\text{Rain})P(\text{Traffic}|\text{Rain})P(\text{Umbrella}|\text{Rain})$$

- Bayes' nets / graphical models help us express conditional independence assumptions

Ghostbusters Chain Rule

- Each sensor depends only on where the ghost is
- That means, the two sensors are conditionally independent, given the ghost position
- T: Top square is red
B: Bottom square is red
G: Ghost is in the top

$$P(T,B,G) = P(G) P(T|G) P(B|G)$$

T	B	G	P(T,B,G)
+t	+b	+g	0.16
+t	+b	¬g	0.16
+t	¬b	+g	0.24
+t	¬b	¬g	0.04
¬t	+b	+g	0.04
¬t	+b	¬g	0.24
¬t	¬b	+g	0.06
¬t	¬b	¬g	0.06

- Givens:

$$P(+g) = 0.5$$

$$P(+t \mid +g) = 0.8$$

$$P(+t \mid \neg g) = 0.4$$

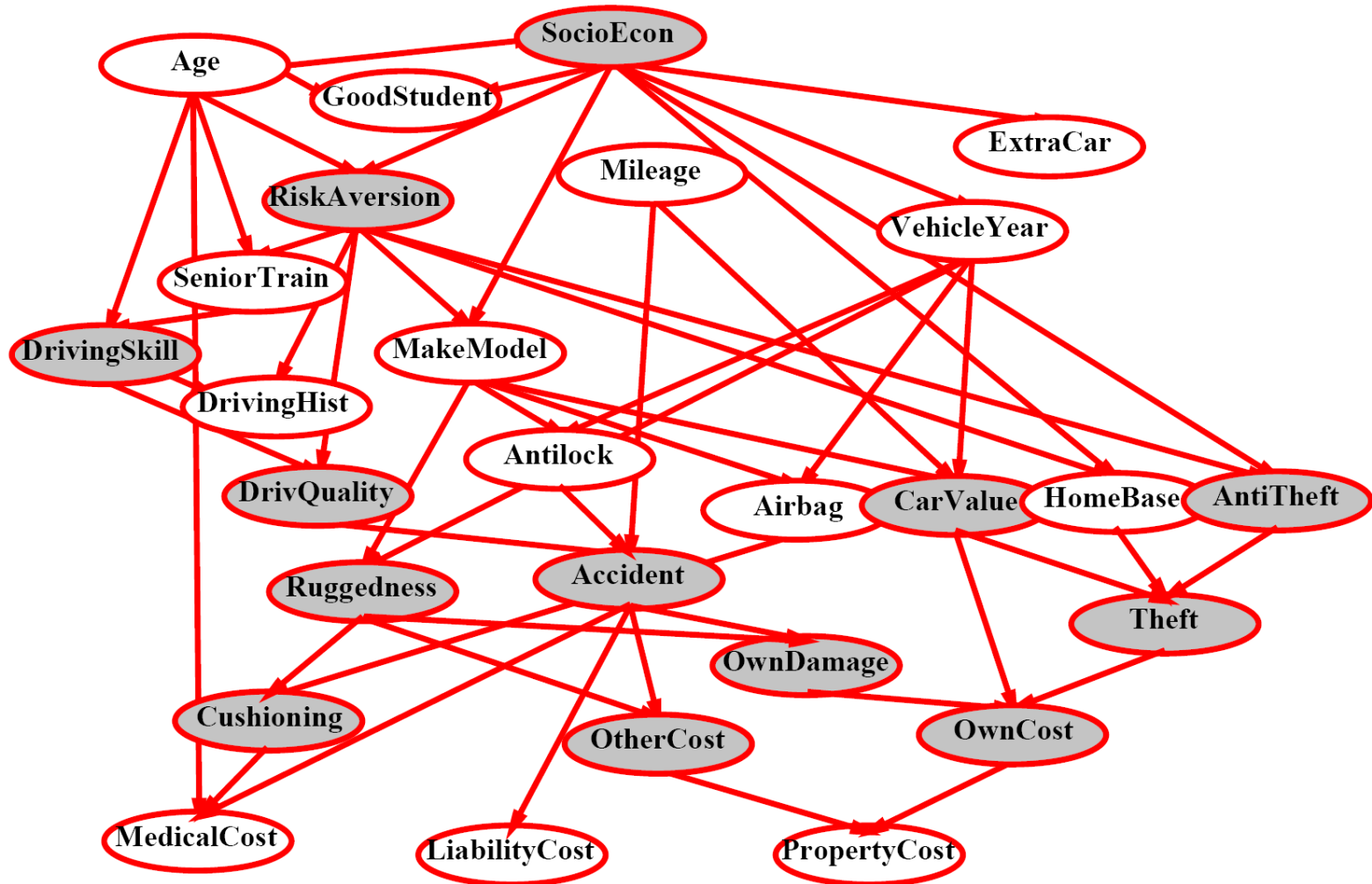
$$P(+b \mid +g) = 0.4$$

$$P(+b \mid \neg g) = 0.8$$

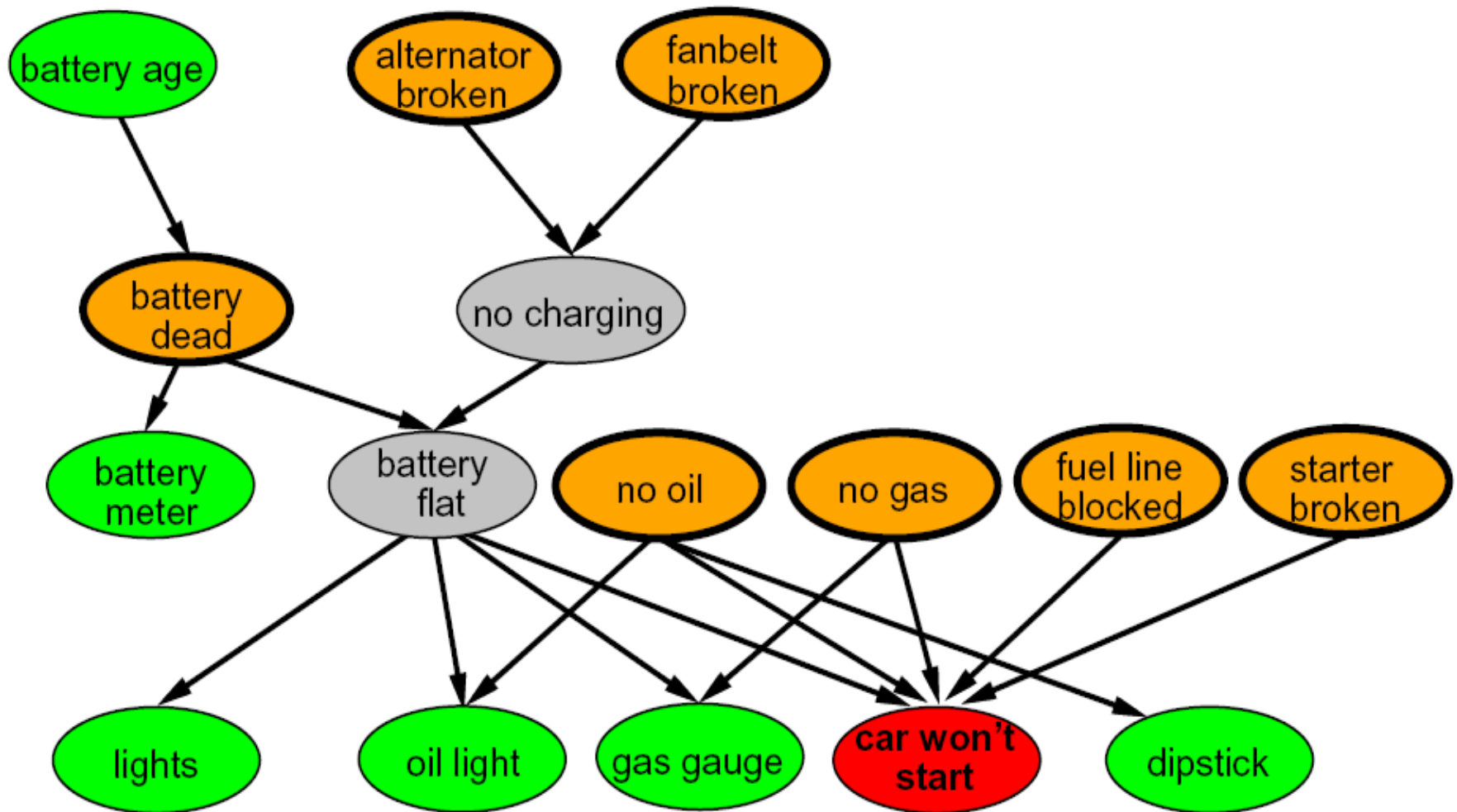
Bayes' Nets: Big Picture

- Two problems with using full joint distribution tables as our 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
- **Bayes' 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
 - For about 10 min, we'll be vague about how these interactions are specified

Example Bayes' Net: Insurance



Example Bayes' Net: Car



Graphical Model Notation

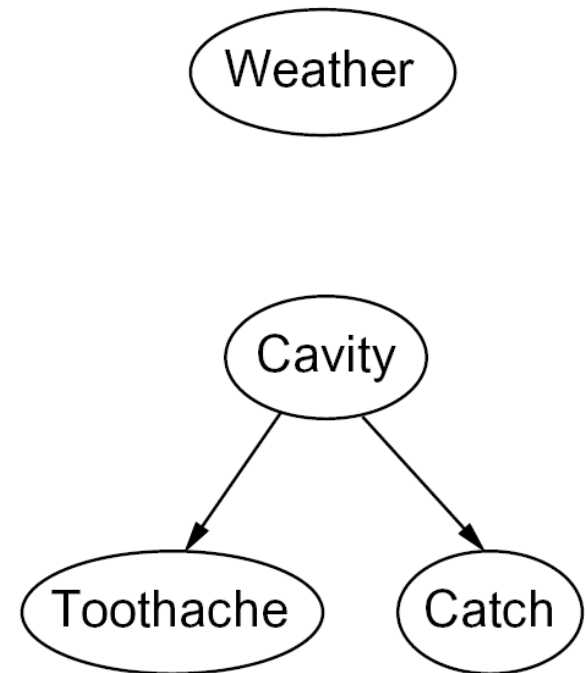
- **Nodes: variables (with domains)**

- Can be assigned (observed) or unassigned (unobserved)

- **Arcs: interactions**

- Similar to CSP constraints
- Indicate “direct influence” between variables
- Formally: encode conditional independence (more later)

- For now: imagine that arrows mean direct causation (in general, they don't!)

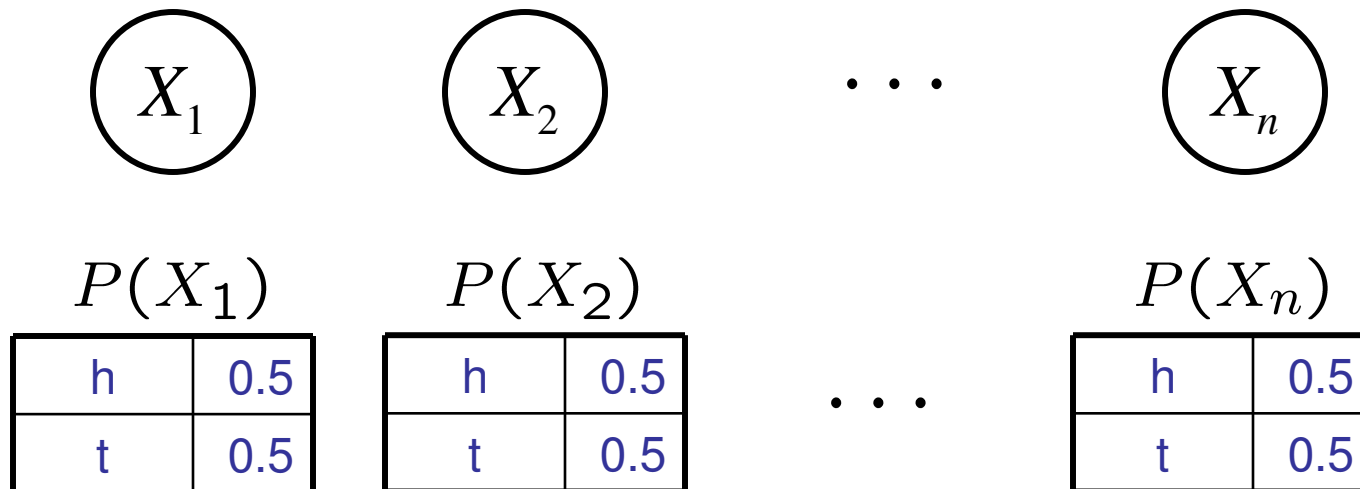


- N independent coin flips

Example: Coin Flips

- No interactions between variables:
absolute independence

Example: Coin Flips

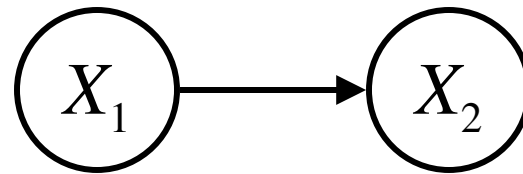
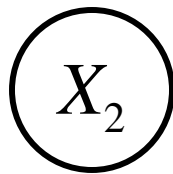
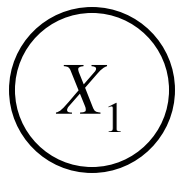


$$P(h, h, t, h) =$$

Only distributions whose variables are absolutely independent can be represented by a Bayes' net with no arcs.

Example: Coins

- Extra arcs don't prevent representing independence, just allow non-independence



$P(X_1)$

h	0.5
t	0.5

$P(X_2)$

h	0.5
t	0.5

$P(X_1)$

h	0.5
t	0.5

$P(X_2|X_1)$

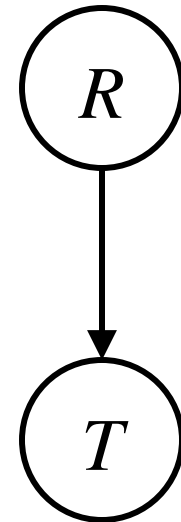
h h	0.5
t h	0.5

h t	0.5
t t	0.5

- Adding unneeded arcs isn't wrong, it's just inefficient

Example: Traffic

- Variables:
 - R: It rains
 - T: There is traffic
- Model 1: independence
- Model 2: rain causes traffic
- Why is an agent using model 2 better?



Example: Traffic II

- Let's build a causal graphical model
- Variables
 - T: Traffic
 - R: It rains
 - L: Low pressure
 - D: Roof drips
 - B: Ballgame
 - C: Cavity

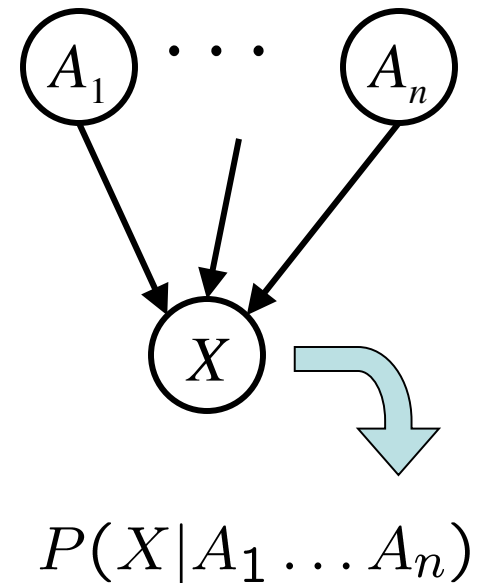
Example: Alarm Network

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

Bayes' Net Semantics

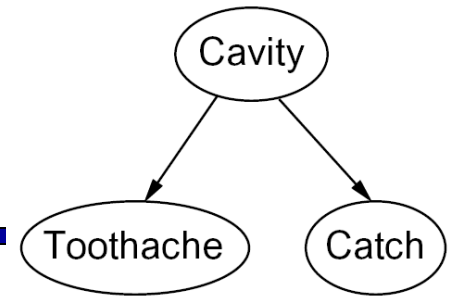
- Let's formalize the semantics of a Bayes' net
 - A set of nodes, one per variable X
 - A directed, acyclic graph
 - A conditional distribution for each node
 - A collection of distributions over X , one for each combination of parents' values
- CPT: conditional probability table
- Description of a noisy “causal” process

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



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

Probabilities in BNs

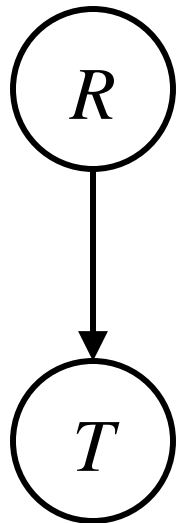


- Bayes' nets **implicitly** encode joint distributions
 - As a product of local conditional distributions
 - To see what probability a BN gives to a full assignment, multiply all the relevant conditionals together:

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

- Example:
 $P(+cavity, +catch, \neg toothache)$
- This lets us reconstruct any entry of the full joint
- Not every BN can represent every joint distribution
 - The topology enforces certain conditional independencies

Example: Traffic



$P(R)$

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

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

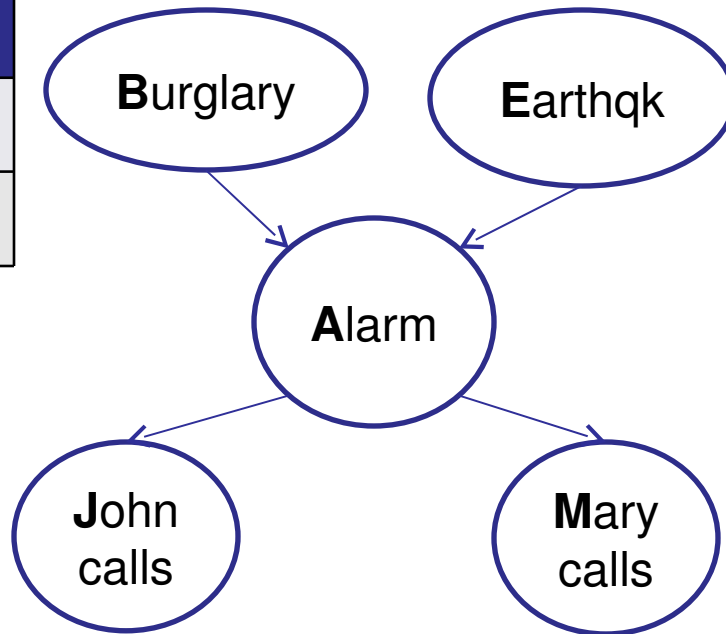
$P(T|R)$

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

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

Example: Alarm Network

B	P(B)
+b	0.001
¬b	0.999



E	P(E)
+e	0.002
¬e	0.998

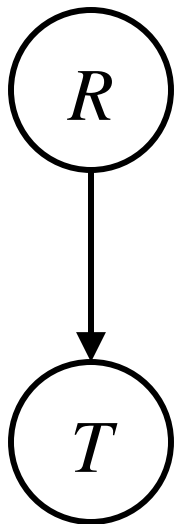
A	J	P(J A)
+a	+j	0.9
+a	¬j	0.1
¬a	+j	0.05
¬a	¬j	0.95

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

Example: Traffic

- Causal direction



$P(R)$

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

$P(T|R)$

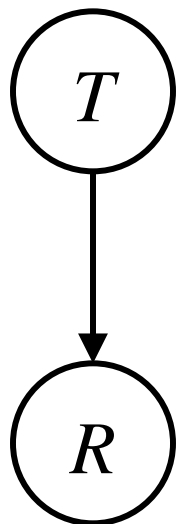
r	t	$3/4$
	$\neg t$	$1/4$
$\neg r$	t	$1/2$
	$\neg t$	$1/2$

$P(T, R)$

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

Example: Reverse Traffic

- Reverse causality?



$P(T)$

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

$P(T, R)$

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

Causality?

- When Bayes' 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 really encodes conditional independence**