# 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

# 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 = yellow | G=(1,1)) = 0.1

| | | |
|---|---|---|
| 0.11 | 0.11 | 0.11 |
| 0.11 | 0.11 | 0.11 |
| 0.11 | 0.11 | 0.11 |

- We can calculate the posterior distribution P(G|r) over ghost locations given a reading using Bayes' rule:

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

| | | |
|---|---|---|
| 0.17 | 0.10 | 0.10 |
| 0.09 | 0.17 | 0.10 |
| <0.01 | 0.09 | 0.17 |

# The Chain Rule

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

- Trivial decomposition:

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

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

- With assumption of conditional independence:

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

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

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

# 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
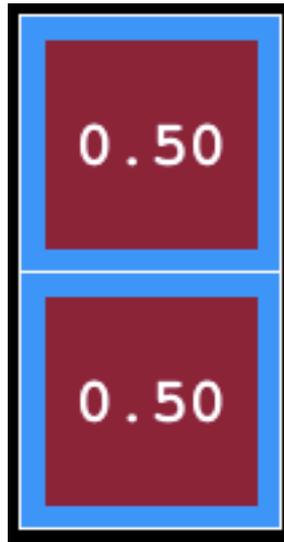
- Queries:
  P( +g) = ??
  P( +g | +t) = ??
  P( +g | +t, -b) = ??

- 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 |

0.50

0.50

# 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

- Givens:
  P( +g ) = 0.5
  P( +t | +g ) = 0.8
  P( +t | ¬g ) = 0.4
  P( +b | +g ) = 0.4
  P( +b | ¬g ) = 0.8

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

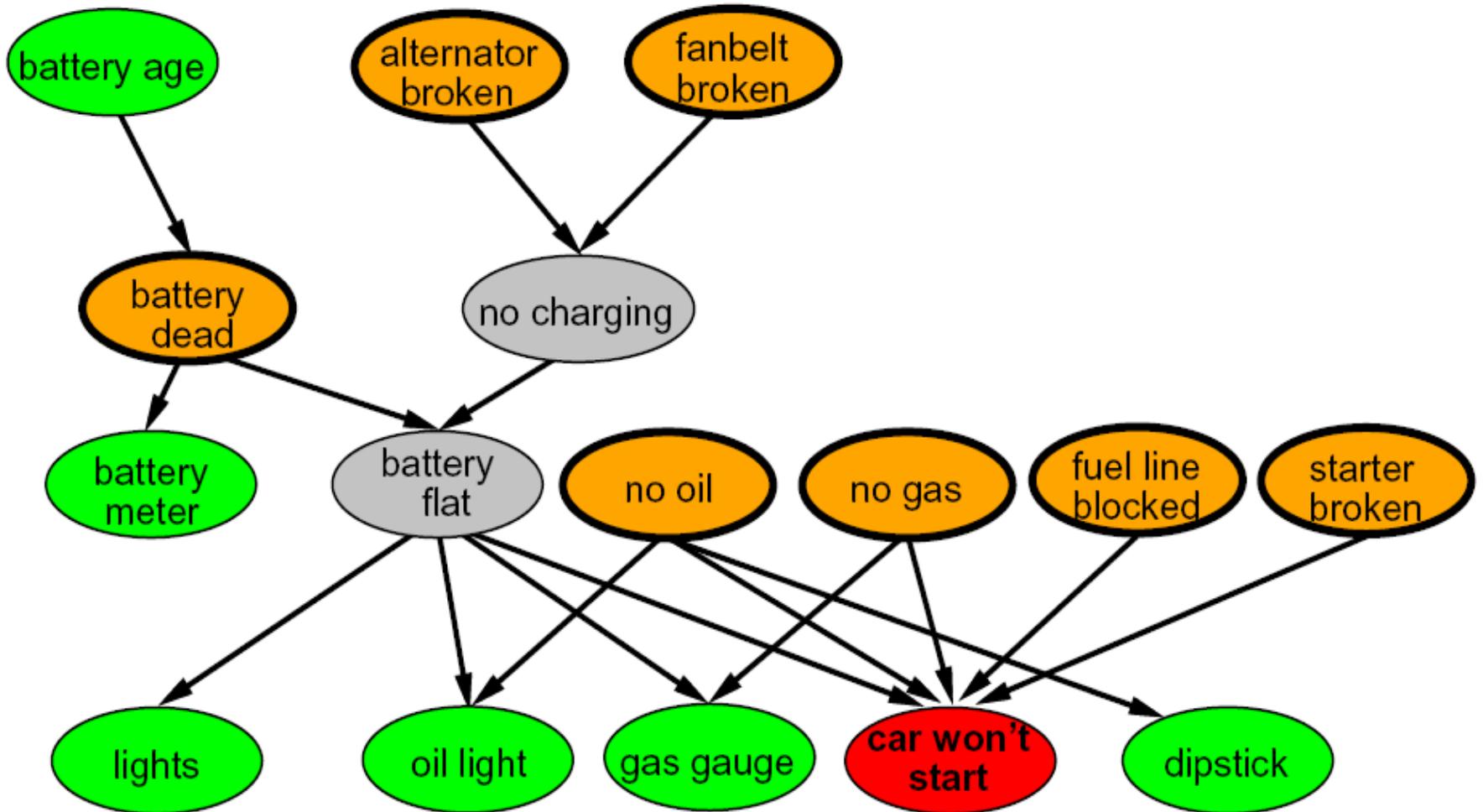| T | B | G | P(T,B,G) |
|---|---|---|---|
| +t | +b | +g | 0.16 |
| +t | +b | ¬g | 0.16 |
| +t | ¬b | +g | 0.24 |
| +t | ¬b | ¬g | 0.04 |
| ¬t | +b | +g | 0.04 |
| ¬t | +b | ¬g | 0.24 |
| ¬t | ¬b | +g | 0.06 |
| ¬t | ¬b | ¬g | 0.06 |

# 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 now, we'll be vague about how these interactions are specified

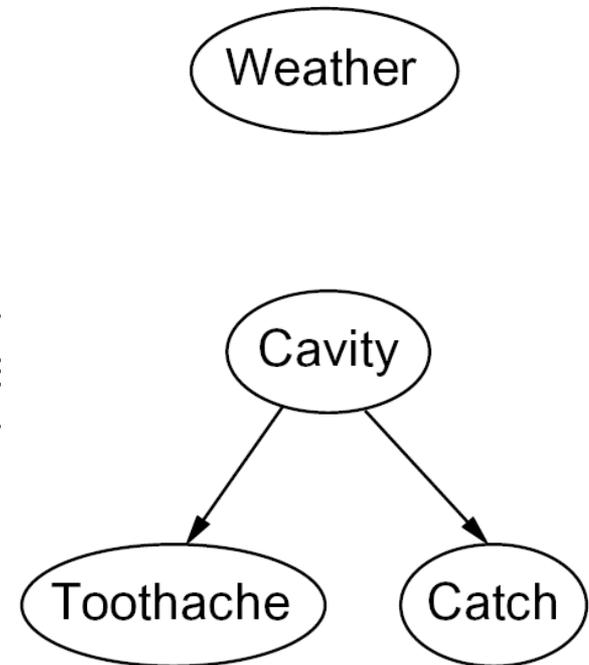# Example Bayes' Net: Insurance

# Example Bayes' Net: Car

# Graphical Model Notation

- Nodes: variables (with domains)
  - Can be assigned (observed) or unassigned (unobserved)

- Arcs: interactions
  - Indicate "direct influence" between variables
  - Formally: encode conditional independence (more later)

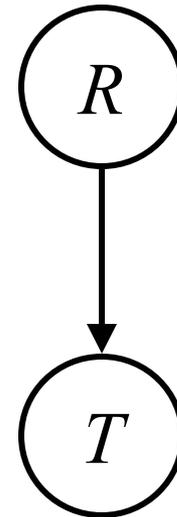- For now: imagine that arrows mean direct causation (in general, they don't!)

Weather

Cavity

Toothache    Catch

# Example: Coin Flips

- N independent coin flips

$$X_1 \qquad X_2 \qquad \cdots \qquad X_n$$

- No interactions between variables: absolute independence

# Example: Traffic

- Variables:
  - R: It rains
  - T: There is traffic

- Model 1: independence

- Model 2: rain causes traffic

- Would an agent using model 2 better?
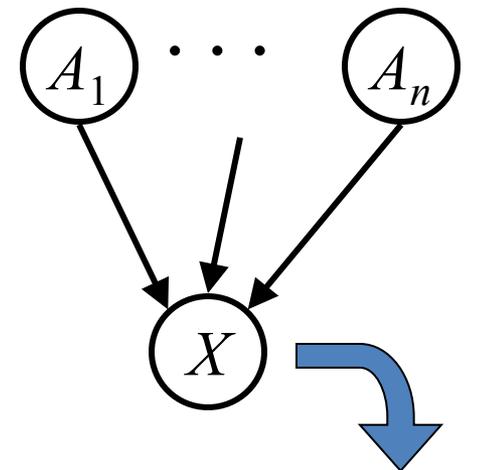
$R$

$T$

# 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

# 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
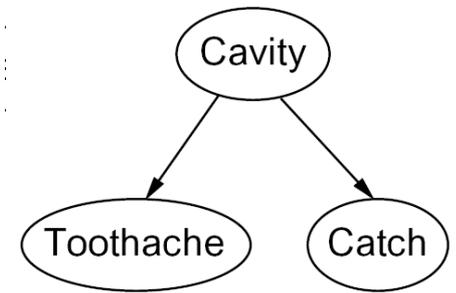
    $$P(X|a_1 \ldots a_n)$$

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



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

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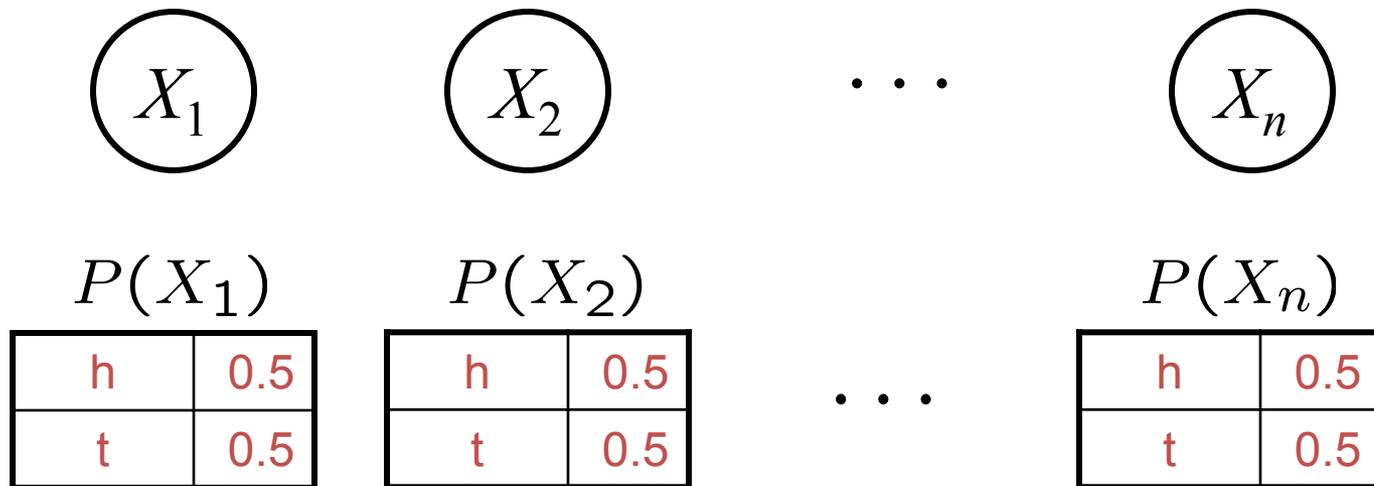
$$P(x_1, x_2, \ldots x_n) = \prod_{i=1}^{n} P(x_i | parents(X_i))$$

  - Example:

$$P(+cavity, +catch, \neg toothache)$$

- This lets us reconstruct any entry of the full joint
- Not every BN can represent every joint distribution
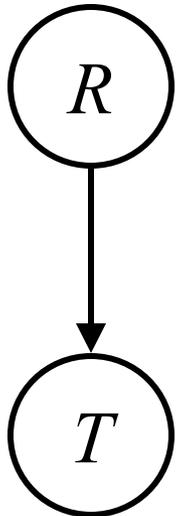  - The topology enforces certain conditional independencies

# Example: 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: Traffic

$P(R)$

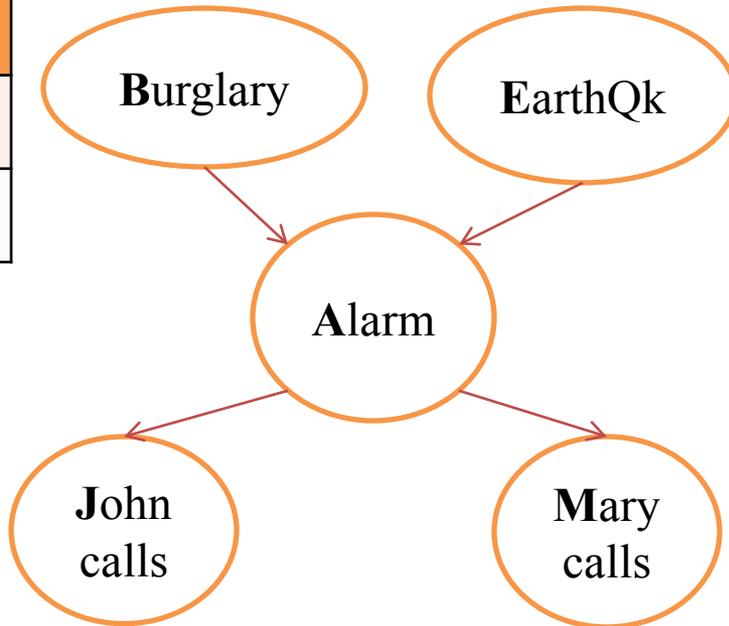| | |
|---|---|
| +r | 1/4 |
| ¬r | 3/4 |

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

$P(T|R)$

+r →

| | |
|---|---|
| +t | 3/4 |
| ¬t | 1/4 |

¬r →

| | |
|---|---|
| +t | 1/2 |
| ¬t | 1/2 |

# Example: Alarm Network

| B | P(B) |
|---|---|
| +b | 0.001 |
| ¬b | 0.999 |

| E | P(E) |
|---|---|
| +e | 0.002 |
| ¬e | 0.998 |

**B**urglary

**E**arthQk

**A**larm

**J**ohn calls

**M**ary calls

| 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 |

| 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 |

# Example: Alarm Network

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

**B**urglary   **E**arthQk

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

**A**larm

**J**ohn calls   **M**ary calls

| B | E | P(A\|B,E) |
|-----|-----|-----|
| +b | +e | 0.95 |
| +b | ¬e | 0.94 |
| ¬b | +e | 0.29 |
| ¬b | ¬e | 0.001 |

| A | P(J\|A) |
|-----|-----|
| +a | 0.9 |
| ¬a | 0.05 |

| A | P(M\|A) |
|-----|-----|
| +a | 0.7 |
| ¬a | 0.01 |

# Bayes' Nets

- So far: how a Bayes' net encodes a joint distribution

- Next: how to answer queries about that distribution
  - Key idea: conditional independence
  - Main goal: answer queries about conditional independence and influence

- After that: how to answer numerical queries (inference)

# 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

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

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



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

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

# Example: Alarm Network

| B | P(B) |
|---|---|
| +b | 0.001 |
| ¬b | 0.999 |

| E | P(E) |
|---|---|
| +e | 0.002 |
| ¬e | 0.998 |

**B**urglary

**E**arthqk

**A**larm

**J**ohn calls

**M**ary calls

| 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 |

# Building the (Entire) Joint

- We can take a Bayes' net and build any entry from the full joint distribution it encodes

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

  - Typically, there's no reason to build ALL of it
  - We build what we need on the fly

- To emphasize: every BN over a domain <span style="color:red">implicitly defines a joint distribution</span> over that domain, specified by local probabilities and graph structure

# Size of a Bayes' Net

- How big is a joint distribution over N Boolean variables?
  $$2^N$$

- How big is an N-node net if nodes have up to k parents?
  $$O(N * 2^{k+1})$$

- Both give you the power to calculate $P(X_1, X_2, \ldots X_n)$
- BNs: Huge space savings!
- Also easier to elicit local CPTs
- Also turns out to be faster to answer queries (coming)

# Bayes' Nets So Far

- We now know:
  - What is a Bayes' net?
  - What joint distribution does a Bayes' net encode?

- Now: properties of that joint distribution (independence)
  - Key idea: conditional independence
  - Last class: assembled BNs using an intuitive notion of conditional independence as causality
  - Today: formalize these ideas
  - Main goal: answer queries about conditional independence and influence

- Next: how to compute posteriors quickly (inference)

# Inference by Enumeration

- Given unlimited time, inference in BNs is easy
- Recipe:
  - State the marginal probabilities you need
  - Figure out ALL the atomic probabilities you need
  - Calculate and combine them
- Example:

$$P(+b| + j, +m) =$$

$$\frac{P(+b, +j, +m)}{P(+j, +m)}$$

# Example: Enumeration

- In this simple method, we only need the BN to synthesize the joint entries

$$P(+b, +j, +m) =$$

$$P(+b)P(+e)P(+a|+b, +e)P(+j|+a)P(+m|+a)+$$

$$P(+b)P(+e)P(-a|+b, +e)P(+j|-a)P(+m|-a)+$$

$$P(+b)P(-e)P(+a|+b, -e)P(+j|+a)P(+m|+a)+$$

$$P(+b)P(-e)P(-a|+b, -e)P(+j|-a)P(+m|-a)$$

P(+m | +b, +e)?

- P(+m | +b, +e)?
- P(+m, +b, +e) / P(+b, +e)

P(+m, +b, +e) =

P(+b)P(+e)P(+a|+b,+e)P(+m|+a) +
P(+b)P(+e)P(-a|+b,+e)P(+m|-a)

Find  P(-m, +b, +e)
   Or
Find  P(+b, +e)

# Assume a= true. What is P(B,E)?

- P(B,E|+a) =?



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

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

| A | P(J|A) |
|---|---|
| +a | 0.9 |
| ¬a | 0.05 |

| A | P(M|A) |
|---|---|
| +a | 0.7 |
| ¬a | 0.01 |

| B | E | P(A|B,E) |
|---|---|---|
| +b | +e | 0.95 |
| +b | ¬e | 0.94 |
| ¬b | +e | 0.29 |
| ¬b | ¬e | 0.001 |

# Inference by Enumeration?

# Variable Elimination

- Why is inference by enumeration so slow?
  - You join up the whole joint distribution before you sum out the hidden variables
  - You end up repeating a lot of work!

- Idea: interleave joining and marginalizing!
  - Called "Variable Elimination"
  - Still NP-hard, but usually much faster than inference by enumeration

- We'll need some new notation to define VE

# The Chain Rule

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

- Trivial decomposition:

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

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

- With assumption of conditional independence:

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

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

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

# Conditional Independence

- Reminder: independence
  - X and Y are <span style="color:red">independent</span> if

$$\forall x, y \ \ P(x,y) = P(x)P(y) \ \ \text{-- -- -->} \ \ X \perp\!\!\!\perp Y$$

  - X and Y are <span style="color:red">conditionally independent</span> given Z

$$\forall x, y, z \ \ P(x,y|z) = P(x|z)P(y|z) \text{-- -- -->} X \perp\!\!\!\perp Y | Z$$

  - (Conditional) independence is a property of a distribution

# 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

# Independence in a BN

- Important question about a BN:
  - Are two nodes independent given certain evidence?
  - If yes, can prove using algebra (tedious in general)
  - If no, can prove with a counter example
  - Example:

```
   X  ------>  Y  ------>  Z
```

  - Question: are X and Z necessarily independent?
    - Answer: no.  Example: low pressure causes rain, which causes traffic.
    - X can influence Z, Z can influence X (via Y)
    - Addendum: they *could* be independent: how?

# Causal Chains

- This configuration is a "causal chain"

$$X \rightarrow Y \rightarrow Z$$

X: Low pressure

Y: Rain

Z: Traffic

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

- Is X independent of Z given Y?

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

$$= P(z|y) \qquad \textit{Yes!}$$

- Evidence along the chain "blocks" the influence

# Common Cause

- Another basic configuration: two effects of the same cause
  - Are X and Z independent?

  - Are X and Z independent given Y?

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

$$= P(z|y) \quad \textit{Yes!}$$

  - Observing the cause blocks influence between effects.

Y: Project due

X: Newsgroup busy

Z: Lab full

6

# Common Effect

- Last configuration: two causes of one effect (v-structures)
  - Are X and Z independent?
    - Yes: the ballgame and the rain cause traffic, but they are not correlated
    - Still need to prove they must be (try it!)
  - Are X and Z independent given Y?
    - No: seeing traffic puts the rain and the ballgame in competition as explanation?
  - This is backwards from the other cases
    - Observing an effect activates influence between possible causes.

X: Raining

Z: Ballgame

Y: Traffic

# The General Case

- Any complex example can be analyzed using these three canonical cases

- General question: in a given BN, are two variables independent (given evidence)?
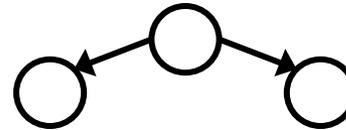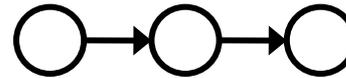
- Solution: analyze the graph

# Reachability

- Recipe: shade evidence nodes

- Attempt 1: if two nodes are connected by an undirected path not blocked by a shaded node, they are conditionally independent

- Almost works, but not quite
  - Where does it break?
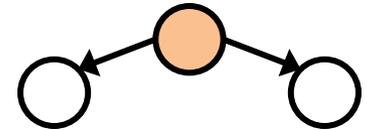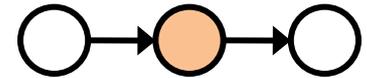  - Answer: the v-structure at T doesn't count as a link in a path unless "active"

# Reachability (D-Separation)

- Question: Are X and Y conditionally independent given evidence vars {Z}?
  - Yes, if X and Y "separated" by Z
  - Look for active paths from X to Y
  - No active paths = independence!

- A path is active if each triple is active:
  - Causal chain $A \rightarrow B \rightarrow C$ where B is unobserved (either direction)
  - Common cause $A \leftarrow B \rightarrow C$ where B is unobserved
  - Common effect (aka v-structure) $A \rightarrow B \leftarrow C$ where B *or one of its descendents* is observed

- All it takes to block a path is a single inactive segment
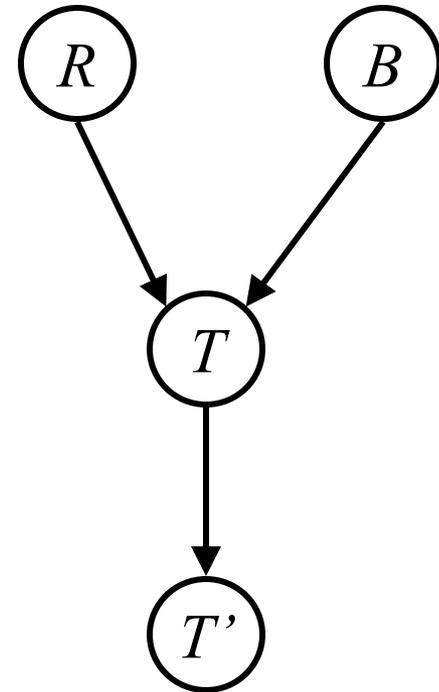
Active Triples | Inactive Triples

# Example

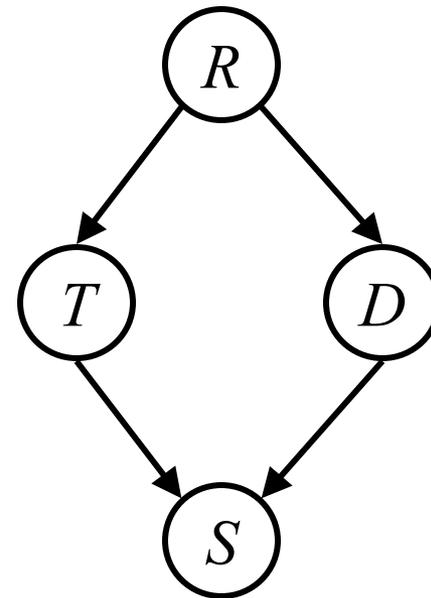$R \perp\!\!\!\perp B$      *Yes*

$R \perp\!\!\!\perp B | T$

$R \perp\!\!\!\perp B | T'$

# Example

- Variables:
  - R: Raining
  - T: Traffic
  - D: Roof drips
  - S: I'm sad

- Questions:

$$T \perp\!\!\!\perp D$$
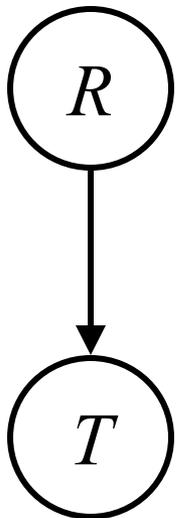
$$T \perp\!\!\!\perp D | R \qquad \textit{Yes}$$

$$T \perp\!\!\!\perp D | R, S$$

# 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
  - 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 independence

# Example: Traffic

- Basic traffic net
- Let's multiply out the joint

$P(R)$

| | |
|---|---|
| r | 1/4 |
| ¬r | 3/4 |

$P(T,R)$

| | | |
|---|---|---|
| r | t | 3/16 |
| r | ¬t | 1/16 |
| ¬r | t | 6/16 |
| ¬r | ¬t | 6/16 |

$P(T|R)$

| r | t | 3/4 |
|---|---|---|
| | ¬t | 1/4 |

| ¬r | t | 1/2 |
|---|---|---|
| | ¬t | 1/2 |

# Example: Reverse Traffic

- Reverse causality?



$P(T)$

| | |
|---|---|
| t | 9/16 |
| ¬t | 7/16 |

$P(R|T)$

| | | |
|---|---|---|
| t | r | 1/3 |
| | ¬r | 2/3 |

| | | |
|---|---|---|
| ¬t | r | 1/7 |
| | ¬r | 6/7 |

$P(T, R)$

| | | |
|---|---|---|
| r | t | 3/16 |
| r | ¬t | 1/16 |
| ¬r | t | 6/16 |
| ¬r | ¬t | 6/16 |

# 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 |

$P(X_2|X_1)$

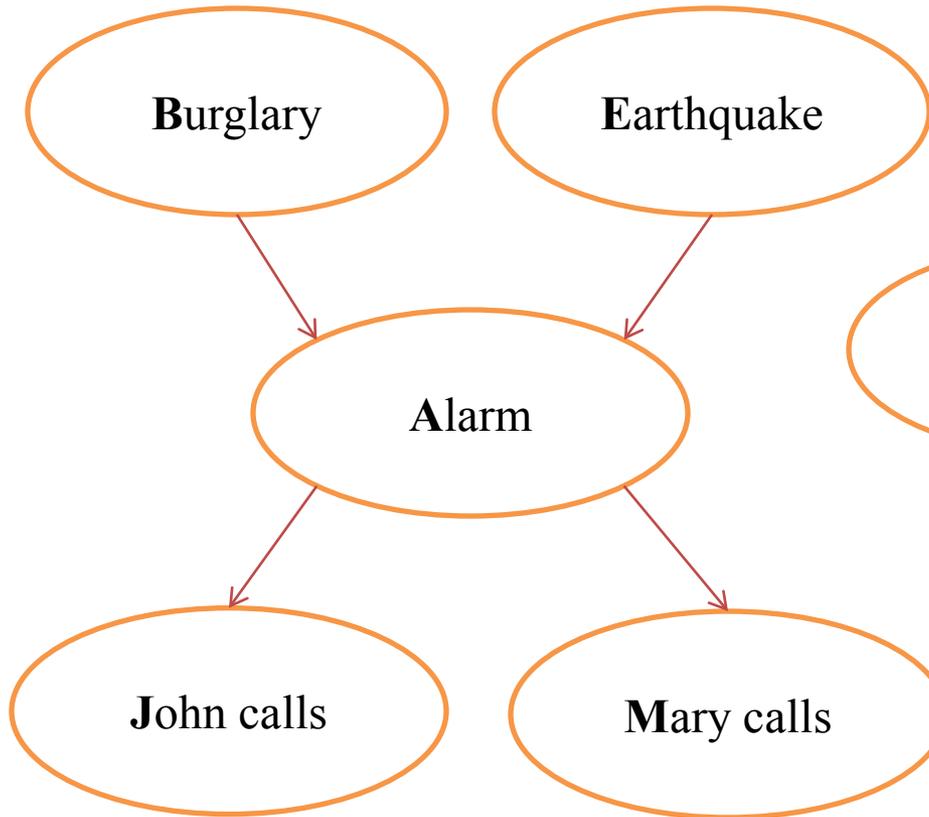| h \| h | 0.5 |
|--------|-----|
| t \| h | 0.5 |

| h \| t | 0.5 |
|--------|-----|
| t \| t | 0.5 |

- Adding unneeded arcs isn't wrong, it's just inefficient
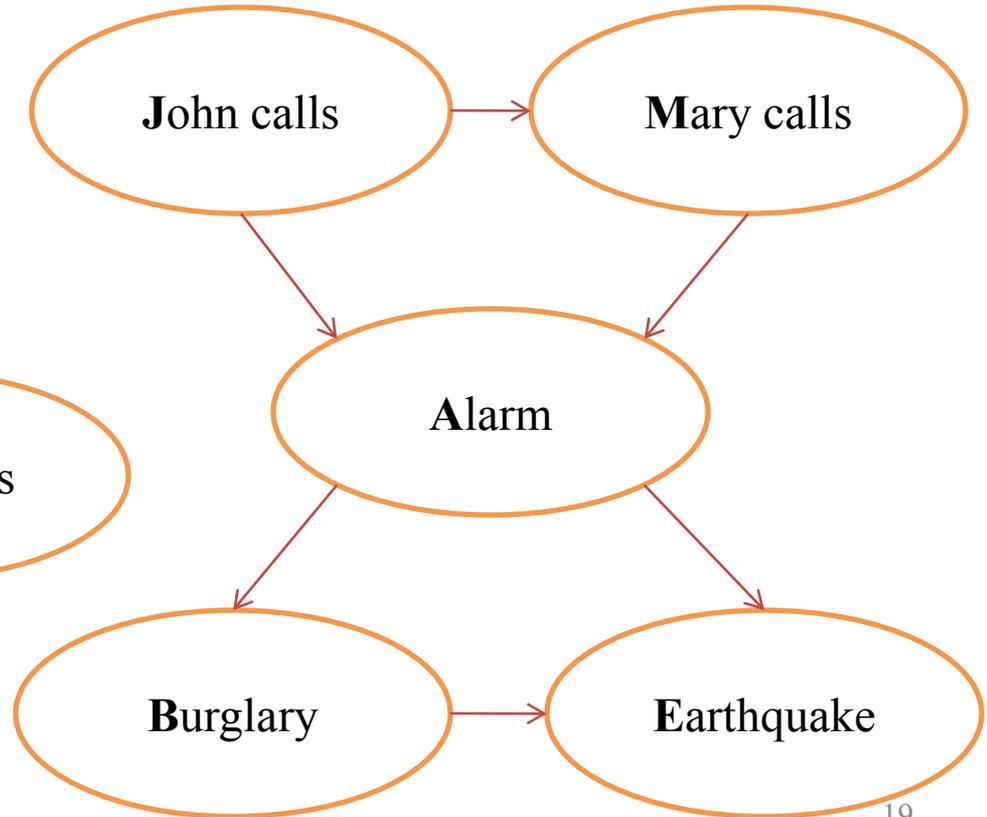
# Changing Bayes' Net Structure

- The same joint distribution can be encoded in many different Bayes' nets
  - Causal structure tends to be the simplest

- Analysis question: given some edges, what other edges do you need to add?
  - One answer: fully connect the graph
  - Better answer: don't make any false conditional independence assumptions

# Example: Alternate Alarm

**B**urglary

**E**arthquake

**A**larm

**J**ohn calls

**M**ary calls

If we reverse the edges, we make different conditional independence assumptions

**J**ohn calls → **M**ary calls

**A**larm

**B**urglary → **E**arthquake

To capture the same joint distribution, we have to add more edges to the graph

19

# Summary

- Bayes nets compactly encode joint distributions

- Guaranteed independencies of distributions can be deduced from BN graph structure

- D-separation gives precise conditional independence guarantees from graph alone

- A Bayes' net's joint distribution may have further (conditional) independence that is not detectable until you inspect its specific distribution