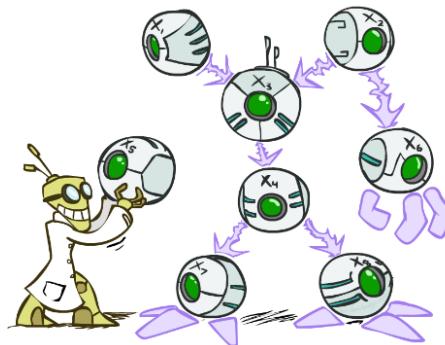


# CS 188: Artificial Intelligence

## Bayes' Nets



Instructors: Dan Klein and Pieter Abbeel --- University of California, Berkeley

[These slides were created by Dan Klein and Pieter Abbeel for CS188 Intro to AI at UC Berkeley. All CS188 materials are available at <http://ai.berkeley.edu>.]

## Probabilistic Models

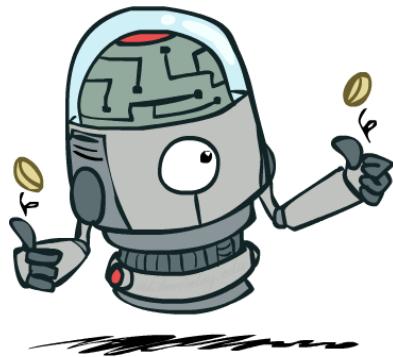
---

- 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



# Independence

---



# Independence

---

- Two variables are *independent* if:

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

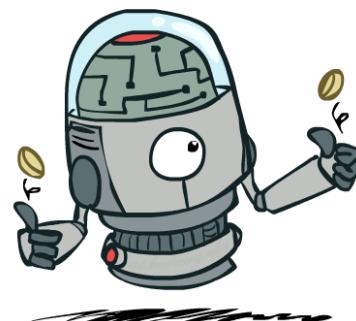
- This says that their joint distribution *factors* into a product two simpler distributions
- Another form:

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

- 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?

$P_1(T, W)$			$P(T)$	$P_2(T, W)$
T	W	P	T	T
hot	sun	0.4	hot	0.3
hot	rain	0.1	hot	0.2
cold	sun	0.2	cold	0.3
cold	rain	0.3	cold	0.2

$P(W)$	
W	P
sun	0.6
rain	0.4

## 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

{



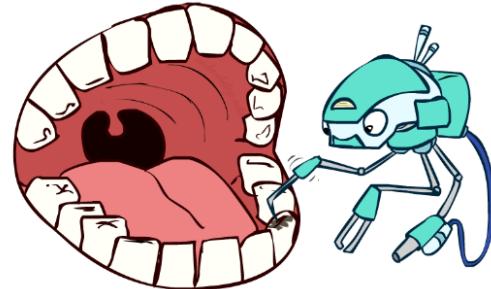
$$2^n \left\{ \begin{array}{c} P(X_1, X_2, \dots, X_n) \\ \text{[Diagram: A wavy grey bar]} \end{array} \right.$$



# 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:
  - $P(+\text{catch} | +\text{toothache}, +\text{cavity}) = P(+\text{catch} | +\text{cavity})$
- The same independence holds if I don't have a cavity:
  - $P(+\text{catch} | +\text{toothache}, -\text{cavity}) = P(+\text{catch} | -\text{cavity})$
- Catch is conditionally independent of Toothache given Cavity:
  - $P(\text{Catch} | \text{Toothache}, \text{Cavity}) = P(\text{Catch} | \text{Cavity})$
- Equivalent statements:
  - $P(\text{Toothache} | \text{Catch}, \text{Cavity}) = P(\text{Toothache} | \text{Cavity})$
  - $P(\text{Toothache}, \text{Catch} | \text{Cavity}) = P(\text{Toothache} | \text{Cavity}) P(\text{Catch} | \text{Cavity})$
  - 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.
- $X$  is conditionally independent of  $Y$  given  $Z$  
$$X \perp\!\!\!\perp Y | Z$$

if and only if:

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

or, equivalently, if and only if

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

# Conditional Independence

---

- What about this domain:

- Traffic
- Umbrella
- Raining

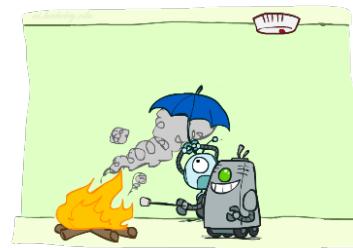
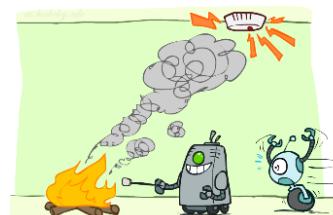


# Conditional Independence

---

- What about this domain:

- Fire
- Smoke
- Alarm



# Conditional Independence and the Chain Rule

---

- 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:

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



- With assumption of conditional independence:

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

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

## Ghostbusters Chain Rule

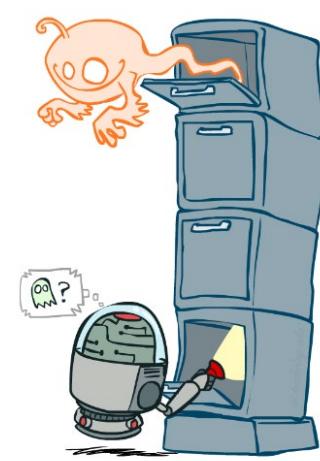
---

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

		0.50	
		0.50	

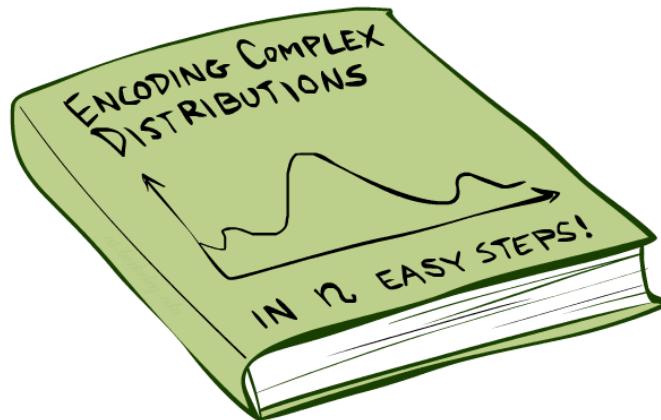
$$P(T, B, G) = P(G) P(T|G) P(B|G)$$

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

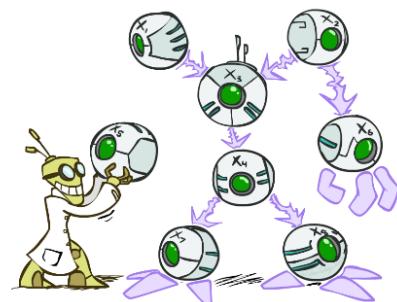
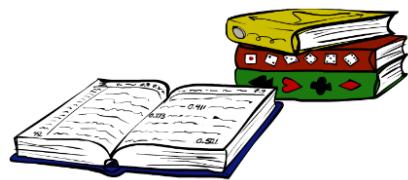
---



## Bayes' Nets: Big Picture

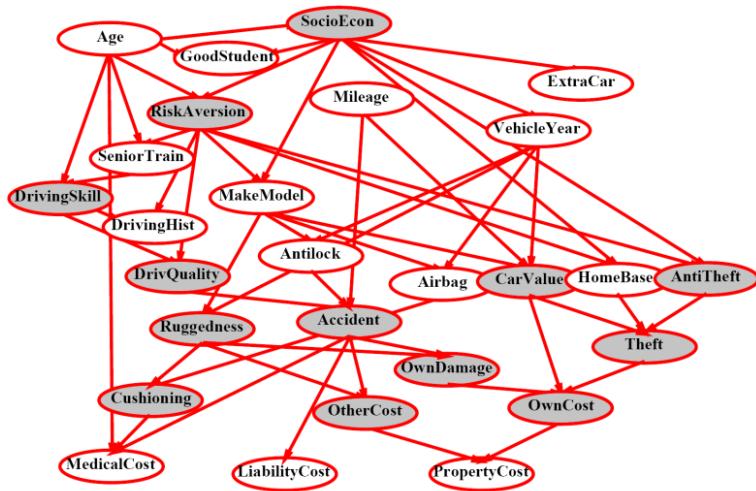
---

- 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



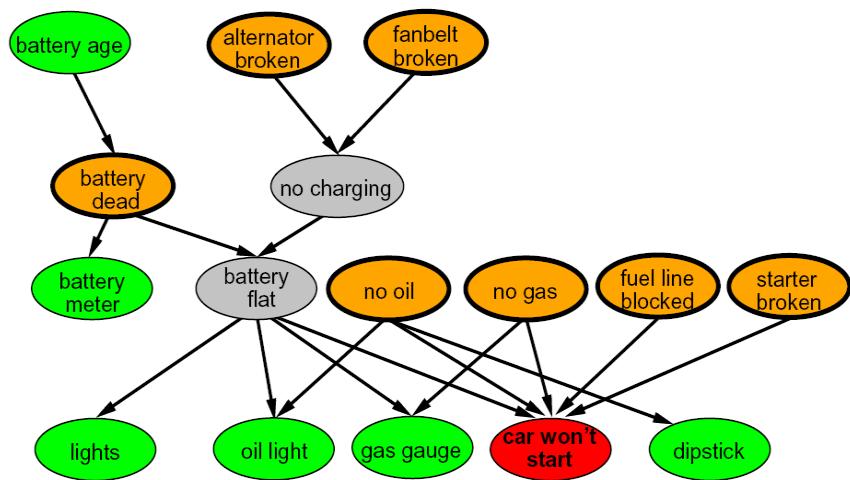
## Example Bayes' Net: Insurance

---



## Example Bayes' Net: Car

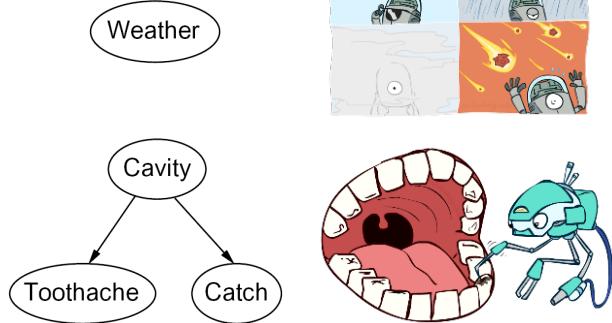
---



# Graphical Model Notation

---

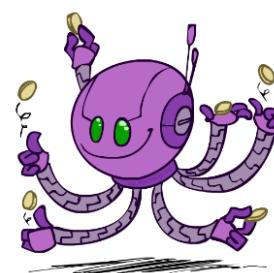
- **Nodes: variables (with domains)**
  - Can be assigned (observed) or unassigned (unobserved)
- **Arcs: interactions**
  - Similar to CSP constraints
  - Indicate “direct influence” between variables
  - Formally: encode conditional independence (more later)
- For now: imagine that arrows mean direct causation (in general, they don’t!)



## Example: Coin Flips

---

- N independent coin flips



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

## Example: Traffic

---

- **Variables:**

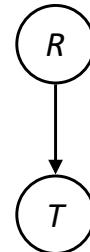
- R: It rains
- T: There is traffic



- Model 1: independence



- 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

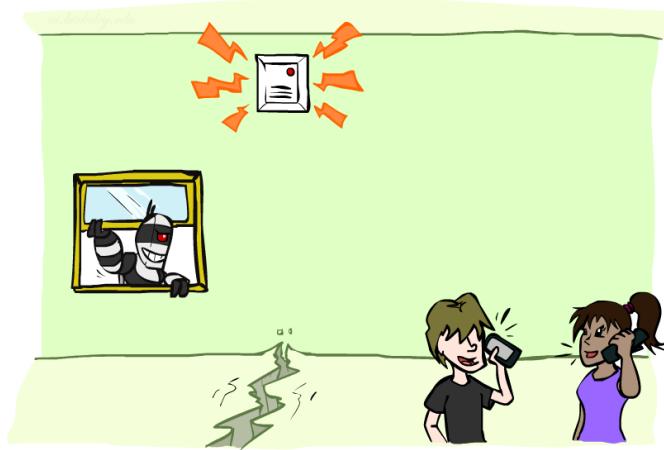


## Example: Alarm Network

---

- **Variables**

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



## Bayes' Net Semantics

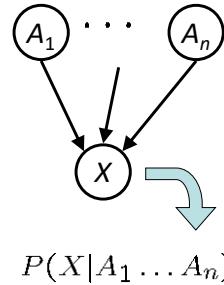
---



# Bayes' Net Semantics



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

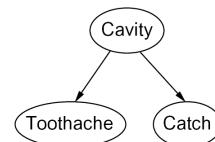
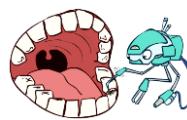
## Probabilities in BNs



- 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:



$$P(+\text{cavity}, +\text{catch}, -\text{toothache})$$

# Probabilities in BNs



- Why are we guaranteed that setting

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

results in a proper joint distribution?

- Chain rule (valid for all distributions):  $P(x_1, x_2, \dots, x_n) = \prod_{i=1}^n P(x_i | x_1 \dots x_{i-1})$

- Assume conditional independences:  $P(x_i | x_1, \dots, x_{i-1}) = P(x_i | \text{parents}(X_i))$

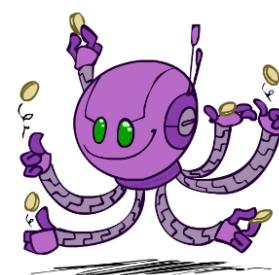
→ Consequence:  $P(x_1, x_2, \dots, x_n) = \prod_{i=1}^n P(x_i | \text{parents}(X_i))$

- Not every BN can represent every joint distribution

- The topology enforces certain conditional independencies

## Example: Coin Flips

$X_1$	$X_2$	...	$X_n$												
$P(X_1)$	$P(X_2)$	...	$P(X_n)$												
<table border="1"><tr><td>h</td><td>0.5</td></tr><tr><td>t</td><td>0.5</td></tr></table>	h	0.5	t	0.5	<table border="1"><tr><td>h</td><td>0.5</td></tr><tr><td>t</td><td>0.5</td></tr></table>	h	0.5	t	0.5	...	<table border="1"><tr><td>h</td><td>0.5</td></tr><tr><td>t</td><td>0.5</td></tr></table>	h	0.5	t	0.5
h	0.5														
t	0.5														
h	0.5														
t	0.5														
h	0.5														
t	0.5														

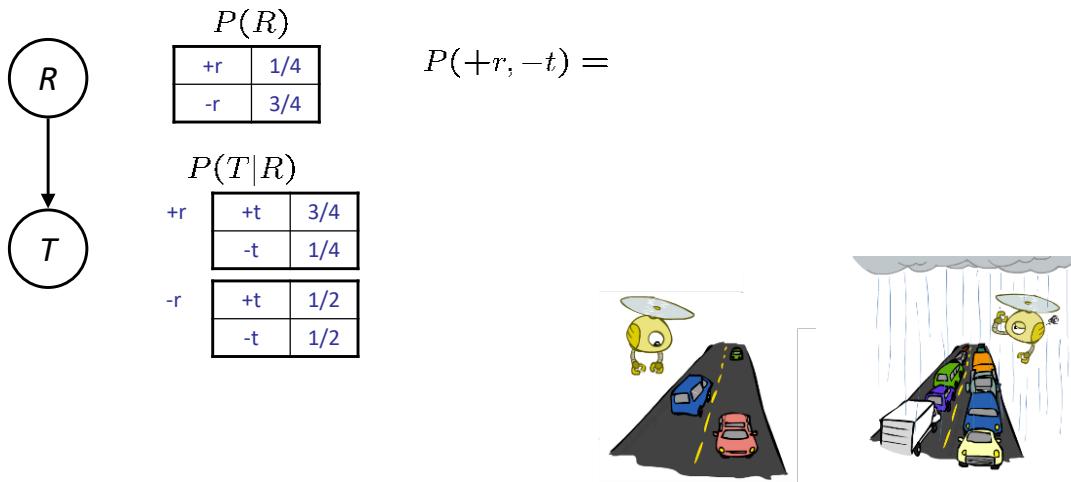


$$P(h, h, t, h) =$$

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

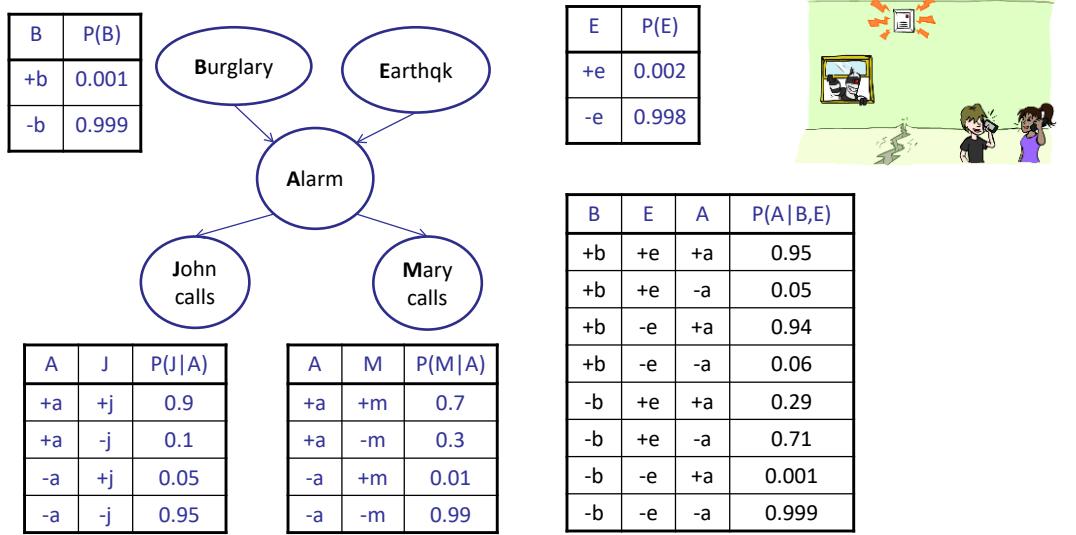
## Example: Traffic

---



## Example: Alarm Network

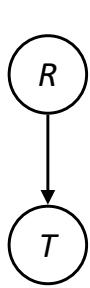
---



## Example: Traffic

---

- Causal direction



$$P(R)$$

+r	1/4
-r	3/4

$$P(T|R)$$

+r	+t	3/4
-r	+t	1/2
-r	-t	1/2
-r	-t	1/2

 $P(T, R)$ 

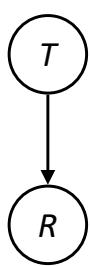
+r	+t	3/16
+r	-t	1/16
-r	+t	6/16
-r	-t	6/16



## Example: Reverse Traffic

---

- Reverse causality?



$$P(T)$$

+t	9/16
-t	7/16

$$P(R|T)$$

+t	+r	1/3
-t	+r	1/7
-t	-r	6/7
-t	-r	6/7

 $P(T, R)$ 

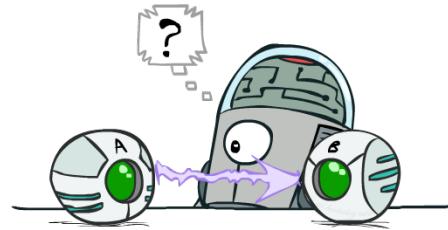
+r	+t	3/16
+r	-t	1/16
-r	+t	6/16
-r	-t	6/16



# 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**

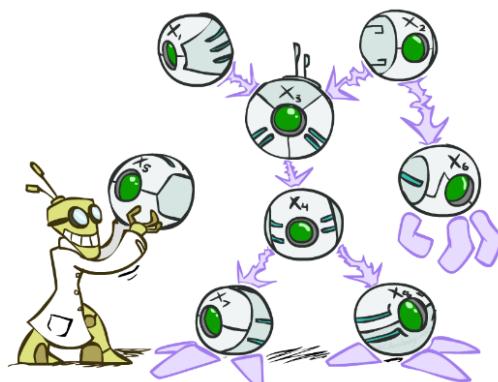


$$P(x_i|x_1, \dots, x_{i-1}) = P(x_i|\text{parents}(X_i))$$

# Bayes' Nets

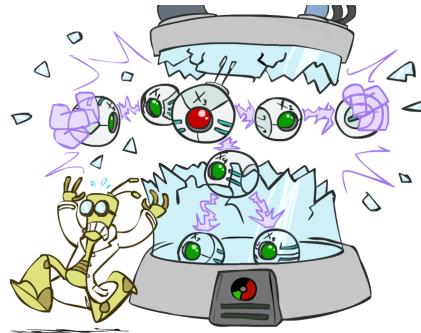
---

- So far: how a Bayes' net encodes a joint distribution
- Next: how to answer queries about that distribution
  - Today:
    - First assembled BNs using an intuitive notion of conditional independence as causality
    - Then saw that key property is conditional independence
  - Main goal: answer queries about conditional independence and influence
- After that: how to answer numerical queries (inference)



# CS 188: Artificial Intelligence

## Bayes' Nets: Independence



Instructors: Pieter Abbeel & Dan Klein --- University of California, Berkeley

[These slides were created by Dan Klein and Pieter Abbeel for CS188 Intro to AI at UC Berkeley. All CS188 materials are available at <http://ai.berkeley.edu>.]

## Probability Recap

---

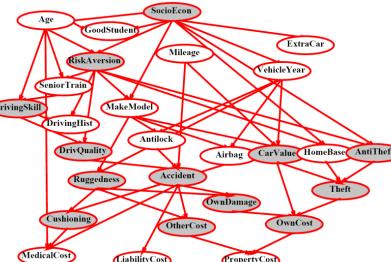
- Conditional probability      
$$P(x|y) = \frac{P(x,y)}{P(y)}$$
- Product rule      
$$P(x,y) = P(x|y)P(y)$$
- Chain rule      
$$\begin{aligned} P(X_1, X_2, \dots, X_n) &= P(X_1)P(X_2|X_1)P(X_3|X_1, X_2)\dots \\ &= \prod_{i=1}^n P(X_i|X_1, \dots, X_{i-1}) \end{aligned}$$
- X, Y independent if and only if:       $\forall x, y : P(x, y) = P(x)P(y)$
- X and Y are conditionally independent given Z if and only if:

$$\forall x, y, z : P(x, y|z) = P(x|z)P(y|z) \quad X \perp\!\!\!\perp Y | Z$$

# Bayes' Nets

---

- A Bayes' net is an efficient encoding of a probabilistic model of a domain



- Questions we can ask:

- Inference: given a fixed BN, what is  $P(X | e)$ ?
- Representation: given a BN graph, what kinds of distributions can it encode?
- Modeling: what BN is most appropriate for a given domain?

# Bayes' Net Semantics

---

- A directed, acyclic graph, one node per random variable
- A conditional probability table (CPT) for each node

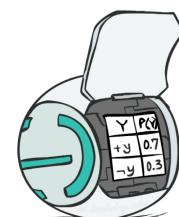
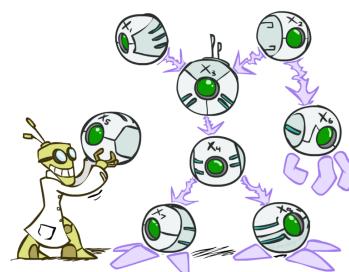
- A collection of distributions over  $X$ , one for each combination of parents' values

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

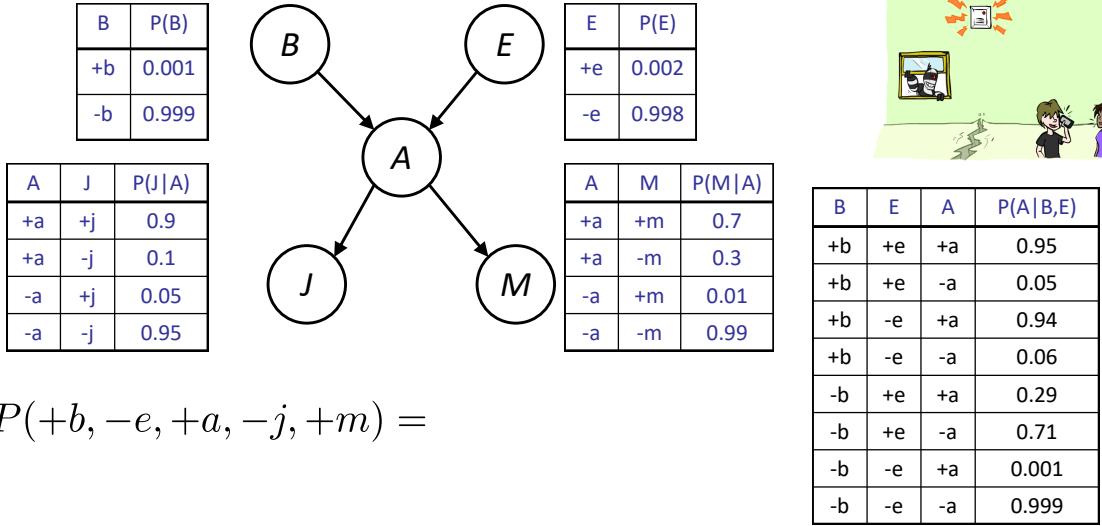
- Bayes' nets implicitly encode joint distributions

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

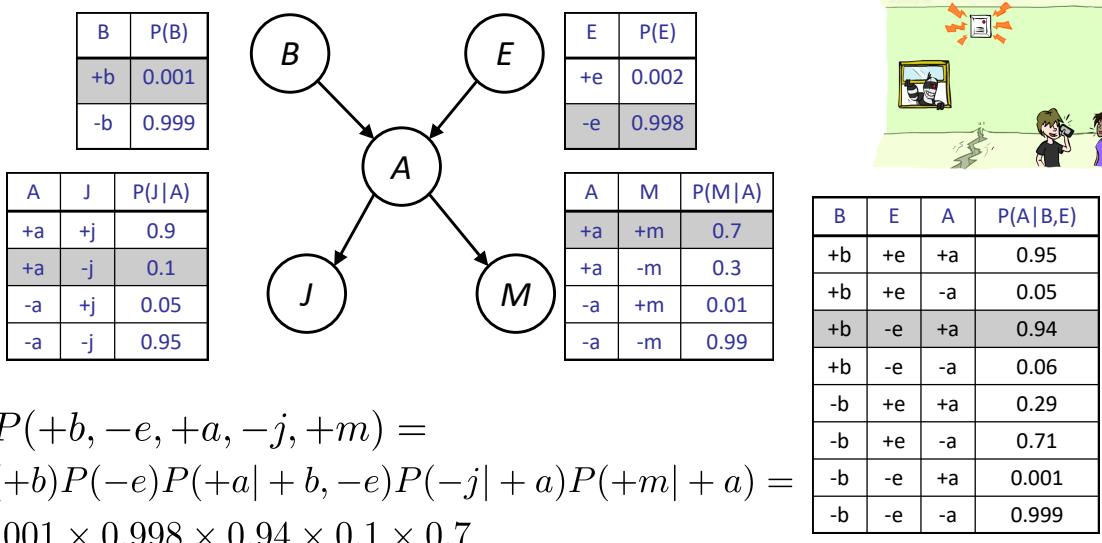
$$P(x_1, x_2, \dots, x_n) = \prod_{i=1}^n P(x_i | \text{parents}(X_i))$$



## Example: Alarm Network



## Example: Alarm Network



# Size of a Bayes' Net

---

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

$$2^N$$

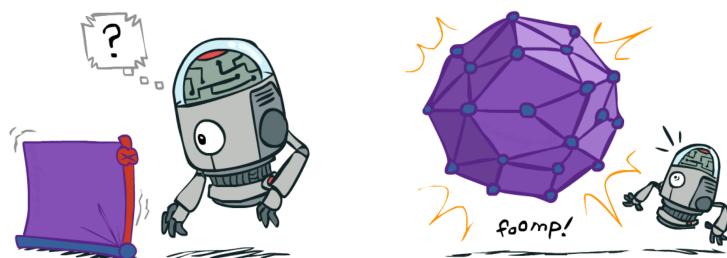
- How big is an N-node net if nodes have up to k parents?

$$O(N * 2^{k+1})$$

- Both give you the power to calculate

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

- BNs: Huge space savings!
- Also easier to elicit local CPTs
- Also faster to answer queries (coming)



# Bayes' Nets

---

## ✓ Representation

- Conditional Independences
- Probabilistic Inference
- Learning Bayes' Nets from Data

# Conditional Independence

---

- X and Y are independent if

$$\forall x, y \ P(x, y) = P(x)P(y) \dashrightarrow X \perp\!\!\!\perp Y$$

- X and Y are conditionally independent given Z

$$\forall x, y, z \ P(x, y|z) = P(x|z)P(y|z) \dashrightarrow X \perp\!\!\!\perp Y|Z$$

- (Conditional) independence is a property of a distribution

- Example:  $\text{Alarm} \perp\!\!\!\perp \text{Fire}|\text{Smoke}$



## Bayes Nets: Assumptions

---

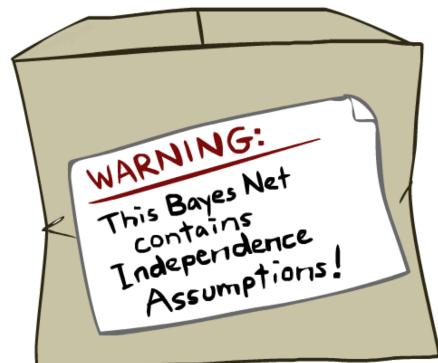
- Assumptions we are required to make to define the Bayes net when given the graph:

$$P(x_i|x_1 \cdots x_{i-1}) = P(x_i|\text{parents}(X_i))$$

- Beyond above “chain rule  $\rightarrow$  Bayes net” conditional independence assumptions

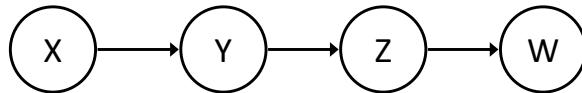
- Often additional conditional independences
  - They can be read off the graph

- Important for modeling: understand assumptions made when choosing a Bayes net graph



## Example

---



- Conditional independence assumptions directly from simplifications in chain rule:
  
  
  
- Additional implied conditional independence assumptions?

## Independence in a BN

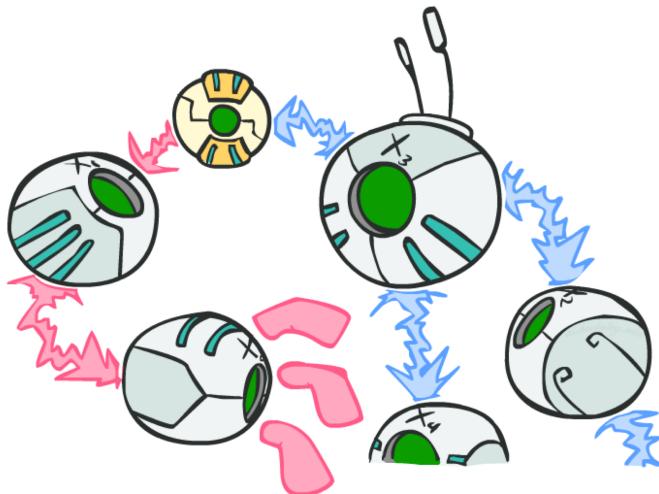
---

- 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:  

```
graph LR; X((X)) --> Y((Y)); Y --> Z((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?

## D-separation: Outline

---



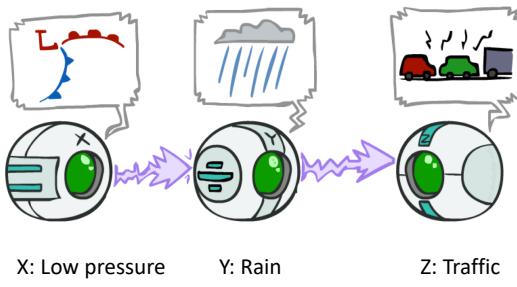
## D-separation: Outline

---

- Study independence properties for triples
- Analyze complex cases in terms of member triples
- D-separation: a condition / algorithm for answering such queries

# Causal Chains

- This configuration is a “causal chain”



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

- Guaranteed X independent of Z? **No!**

▪ One example set of CPTs for which X is not independent of Z is sufficient to show this independence is not guaranteed.

- Example:

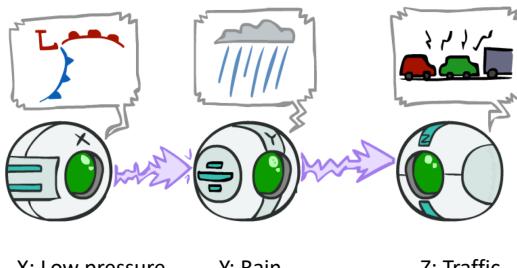
▪ Low pressure causes rain causes traffic, high pressure causes no rain causes no traffic

▪ In numbers:

$$\begin{aligned} P(+y | +x) &= 1, P(-y | -x) = 1, \\ P(+z | +y) &= 1, P(-z | -y) = 1 \end{aligned}$$

# Causal Chains

- This configuration is a “causal chain”



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

- Guaranteed X independent of Z given Y?

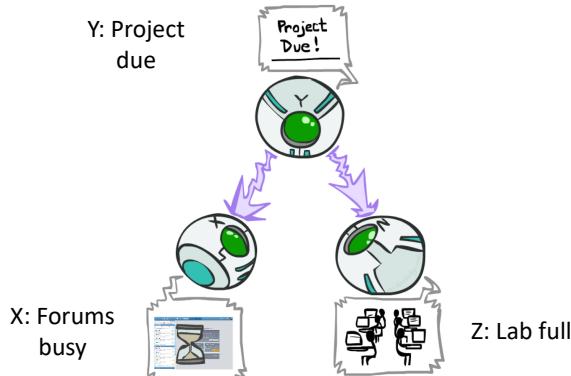
$$\begin{aligned} P(z|x, y) &= \frac{P(x, y, z)}{P(x, y)} \\ &= \frac{P(x)P(y|x)P(z|y)}{P(x)P(y|x)} \\ &= P(z|y) \end{aligned}$$

**Yes!**

- Evidence along the chain “blocks” the influence

# Common Cause

- This configuration is a “common cause”
- Guaranteed X independent of Z? **No!**



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

- One example set of CPTs for which X is not independent of Z is sufficient to show this independence is not guaranteed.

- Example:

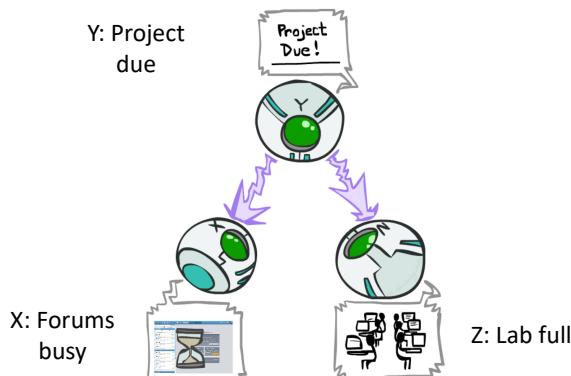
- Project due causes both forums busy and lab full

- In numbers:

$$\begin{aligned} P(+x | +y) &= 1, P(-x | -y) = 1, \\ P(+z | +y) &= 1, P(-z | -y) = 1 \end{aligned}$$

# Common Cause

- This configuration is a “common cause”
- Guaranteed X and Z independent given Y?



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

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

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

$$= P(z|y)$$

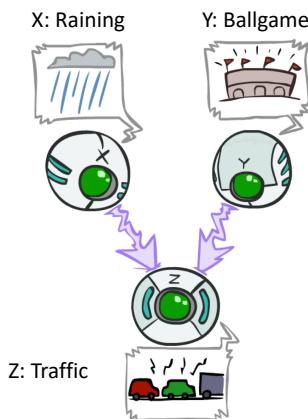
**Yes!**

- Observing the cause blocks influence between effects.

# Common Effect

---

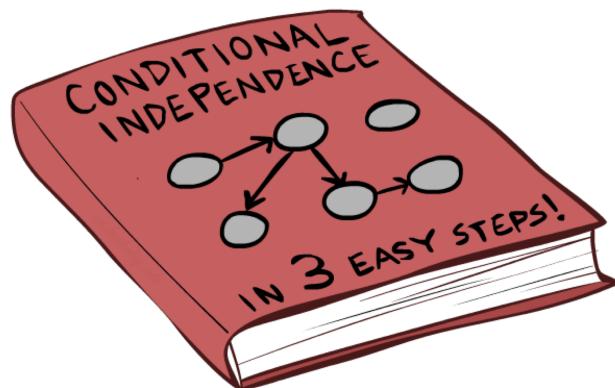
- Last configuration: two causes of one effect (v-structures)



- Are X and Y 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 Y independent given Z?
  - **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.

## The General Case

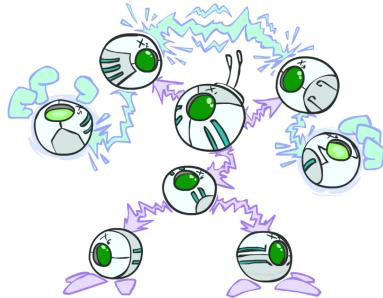
---



# The General Case

---

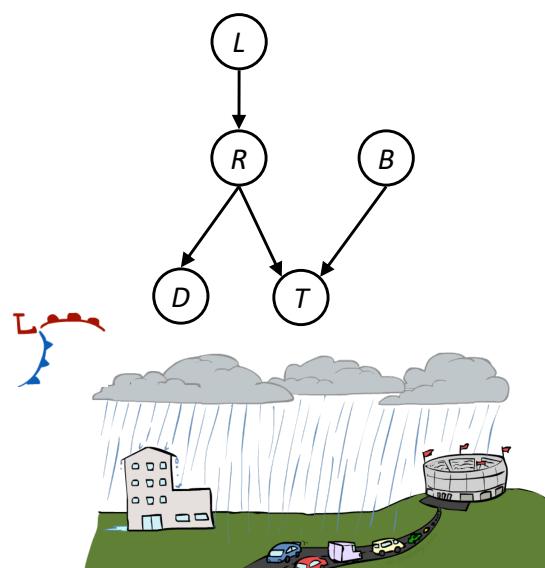
- General question: in a given BN, are two variables independent (given evidence)?
- Solution: analyze the graph
- Any complex example can be broken into repetitions of the three canonical cases



## Reachability

---

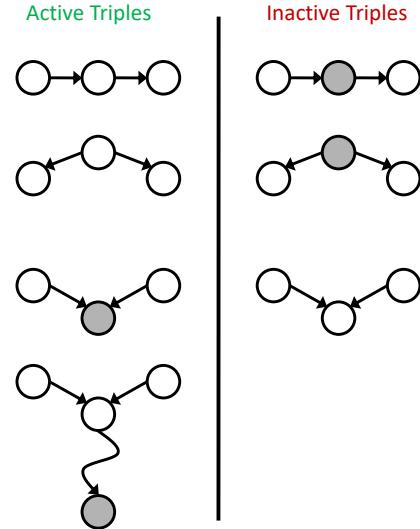
- Recipe: shade evidence nodes, look for paths in the resulting graph
- Attempt 1: if two nodes are connected by an undirected path not blocked by a shaded node, they are conditionally independent
- Almost works, but not quite
  - Where does it break?
  - Answer: the v-structure at T doesn't count as a link in a path unless "active"



# Active / Inactive Paths

---

- Question: Are X and Y conditionally independent given evidence variables  $\{Z\}$ ?
  - Yes, if X and Y “d-separated” by Z
  - Consider all (undirected) paths from X to Y
  - No active paths = independence!
- A path is active if each triple is active:
  - Causal chain A → B → C where B is unobserved (either direction)
  - Common cause A ← B → C where B is unobserved
  - Common effect (aka v-structure)  
 $A \rightarrow B \leftarrow C$  where B or one of its descendants is observed
- All it takes to block a path is a single inactive segment



# D-Separation

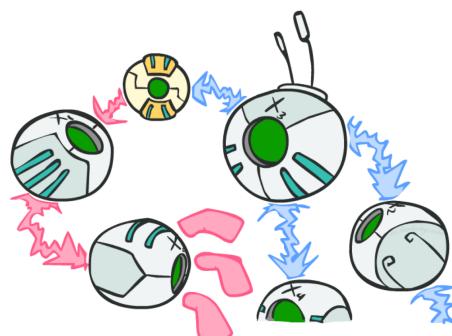
---

- Query:  $X_i \perp\!\!\!\perp X_j | \{X_{k_1}, \dots, X_{k_n}\}$  ?
- Check all (undirected!) paths between  $X_i$  and  $X_j$ 
  - If one or more active, then independence not guaranteed

$$X_i \not\perp\!\!\!\perp X_j | \{X_{k_1}, \dots, X_{k_n}\}$$

- Otherwise (i.e. if all paths are inactive), then independence is guaranteed

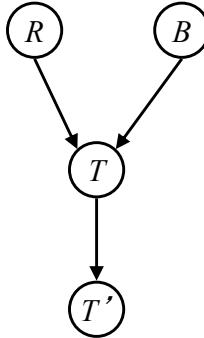
$$X_i \perp\!\!\!\perp X_j | \{X_{k_1}, \dots, X_{k_n}\}$$



## Example

---

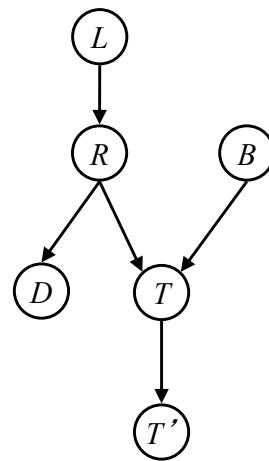
$R \perp\!\!\!\perp B$       Yes  
 $R \perp\!\!\!\perp B|T$   
 $R \perp\!\!\!\perp B|T'$



## Example

---

$L \perp\!\!\!\perp T'|T$       Yes  
 $L \perp\!\!\!\perp B$       Yes  
 $L \perp\!\!\!\perp B|T$   
 $L \perp\!\!\!\perp B|T'$   
 $L \perp\!\!\!\perp B|T, R$       Yes



## Example

---

- **Variables:**

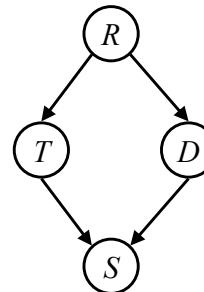
- R: Raining
- T: Traffic
- D: Roof drips
- S: I'm sad

- **Questions:**

$$T \perp\!\!\!\perp D$$

$$T \perp\!\!\!\perp D | R \quad \text{Yes}$$

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



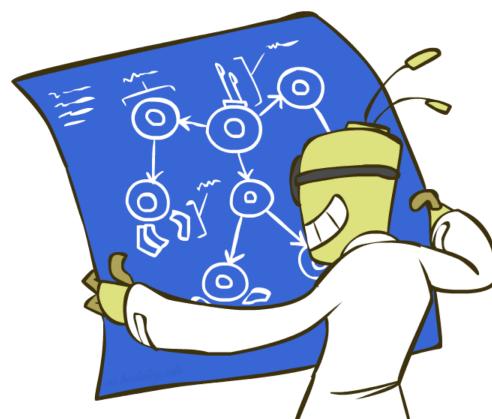
## Structure Implications

---

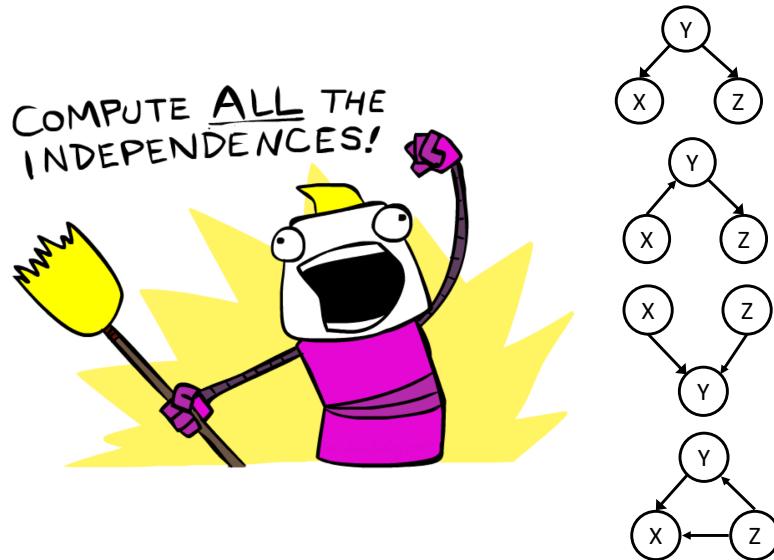
- Given a Bayes net structure, can run d-separation algorithm to build a complete list of conditional independences that are necessarily true of the form

$$X_i \perp\!\!\!\perp X_j | \{X_{k_1}, \dots, X_{k_n}\}$$

- This list determines the set of probability distributions that can be represented

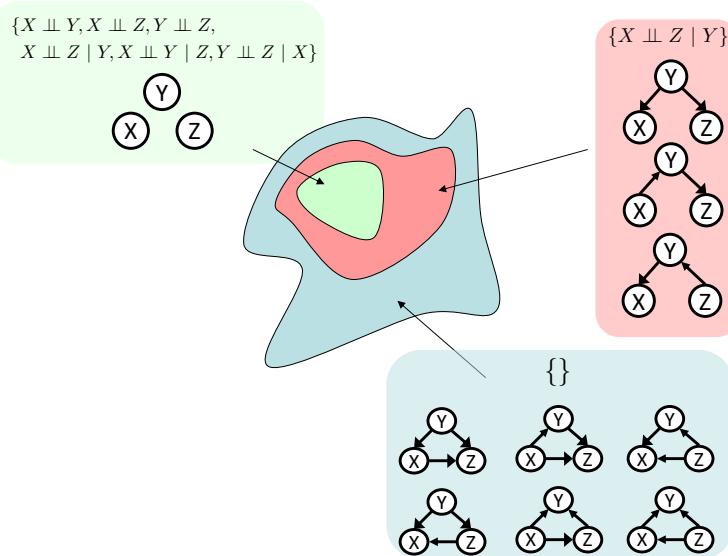


# Computing All Independences



# Topology Limits Distributions

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



## Bayes Nets Representation 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

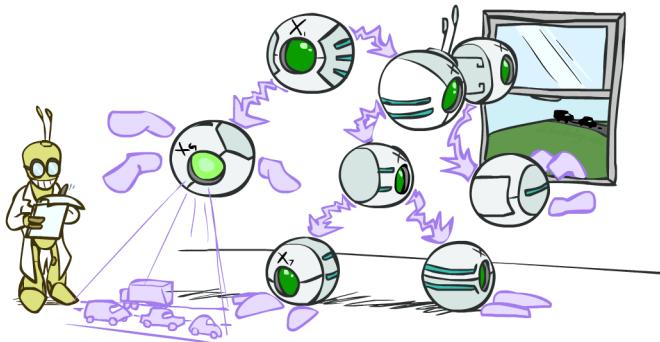
## Bayes' Nets

---

- ✓ Representation
- ✓ Conditional Independences
- Probabilistic Inference
  - Enumeration (exact, exponential complexity)
  - Variable elimination (exact, worst-case exponential complexity, often better)
  - Probabilistic inference is NP-complete
  - Sampling (approximate)
- Learning Bayes' Nets from Data

# CS 188: Artificial Intelligence

## Bayes' Nets: Inference



Instructors: Dan Klein and Pieter Abbeel --- University of California, Berkeley

[These slides were created by Dan Klein and Pieter Abbeel for CS188 Intro to AI at UC Berkeley. All CS188 materials are available at <http://ai.berkeley.edu>.]

## Bayes' Net Representation

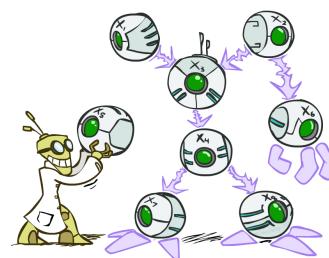
- A directed, acyclic graph, one node per random variable
- A conditional probability table (CPT) for each node
  - A collection of distributions over X, one for each combination of parents' values

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

- Bayes' nets implicitly encode joint distributions

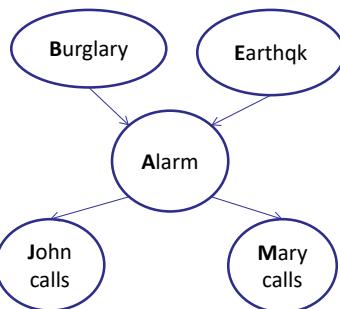
- As a product of local conditional distributions
- To see what probability a BN gives to a full assignment, multiply all the relevant conditionals together:

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



## Example: 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

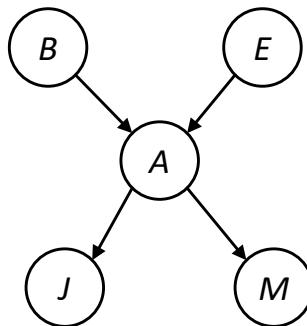
A	M	P(M A)
+a	+m	0.7
+a	-m	0.3
-a	+m	0.01
-a	-m	0.99

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

[Demo: BN Applet]

## Example: Alarm Network

B	P(B)
+b	0.001
-b	0.999



E	P(E)
+e	0.002
-e	0.998



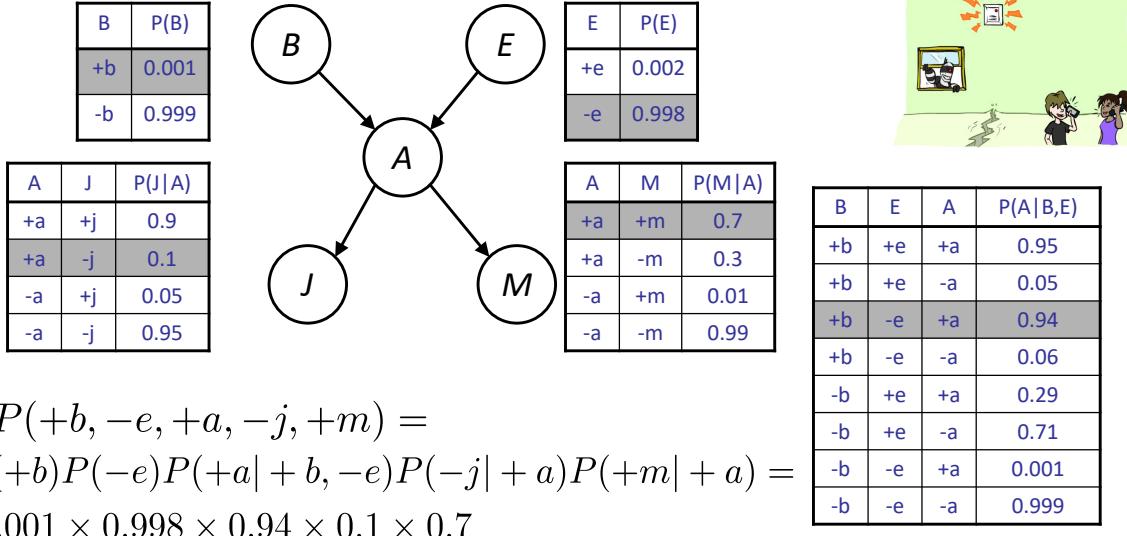
A	J	P(J A)
+a	+j	0.9
+a	-j	0.1
-a	+j	0.05
-a	-j	0.95

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

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

## Example: Alarm Network



## Bayes' Nets

- ✓ Representation
- ✓ Conditional Independences
- Probabilistic Inference
  - Enumeration (exact, exponential complexity)
  - Variable elimination (exact, worst-case exponential complexity, often better)
  - Inference is NP-complete
  - Sampling (approximate)
- Learning Bayes' Nets from Data

# Inference

- Inference: calculating some useful quantity from a joint probability distribution

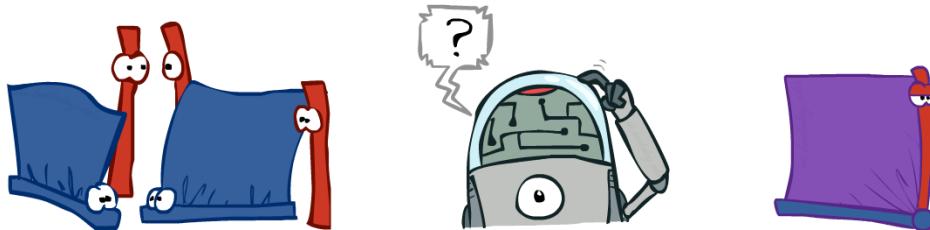
- Examples:

- Posterior probability

$$P(Q|E_1 = e_1, \dots, E_k = e_k)$$

- Most likely explanation:

$$\operatorname{argmax}_q P(Q = q|E_1 = e_1, \dots)$$



## Inference by Enumeration

- General case:

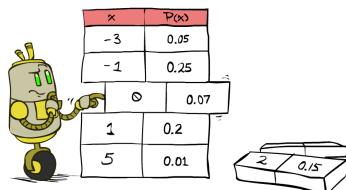
- Evidence variables:  $E_1 \dots E_k = e_1 \dots e_k$
- Query\* variable:  $Q$
- Hidden variables:  $H_1 \dots H_r$

- We want:

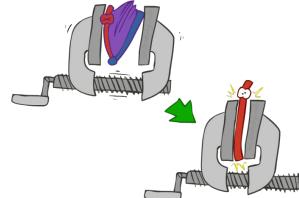
\* Works fine with multiple query variables, too

$$P(Q|e_1 \dots e_k)$$

- Step 1: Select the entries consistent with the evidence



- Step 2: Sum out H to get joint of Query and evidence



- Step 3: Normalize

$$\times \frac{1}{Z}$$

$$Z = \sum_q P(Q, e_1 \dots e_k)$$

$$P(Q|e_1 \dots e_k) = \frac{1}{Z} P(Q, e_1 \dots e_k)$$

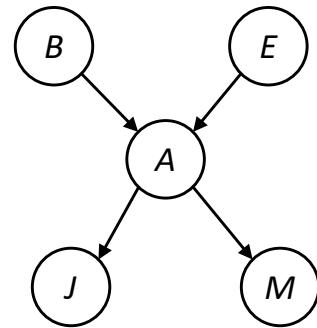
$$P(Q, e_1 \dots e_k) = \sum_{h_1 \dots h_r} \underbrace{P(Q, h_1 \dots h_r, e_1 \dots e_k)}_{X_1, X_2, \dots, X_n}$$

# Inference by Enumeration in Bayes' Net

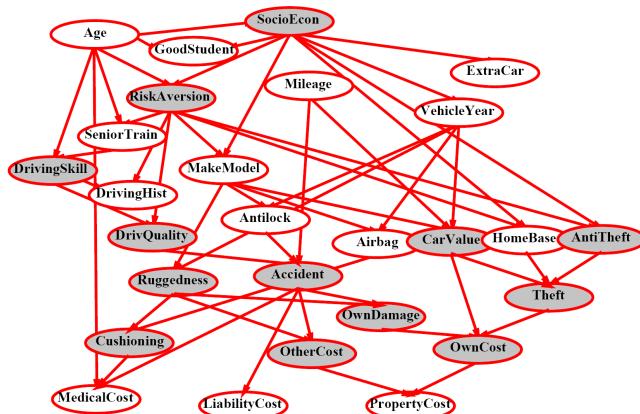
- Given unlimited time, inference in BNs is easy
- Reminder of inference by enumeration by example:

$$P(B \mid +j, +m) \propto_B P(B, +j, +m)$$

$$\begin{aligned} &= \sum_{e,a} P(B, e, a, +j, +m) \\ &= \sum_{e,a} P(B)P(e)P(a|B, e)P(+j|a)P(+m|a) \\ &= P(B)P(+e)P(+a|B, +e)P(+j|+a)P(+m|+a) + P(B)P(+e)P(-a|B, +e)P(+j|-a)P(+m|-a) \\ &\quad P(B)P(-e)P(+a|B, -e)P(+j|+a)P(+m|+a) + P(B)P(-e)P(-a|B, -e)P(+j|-a)P(+m|-a) \end{aligned}$$



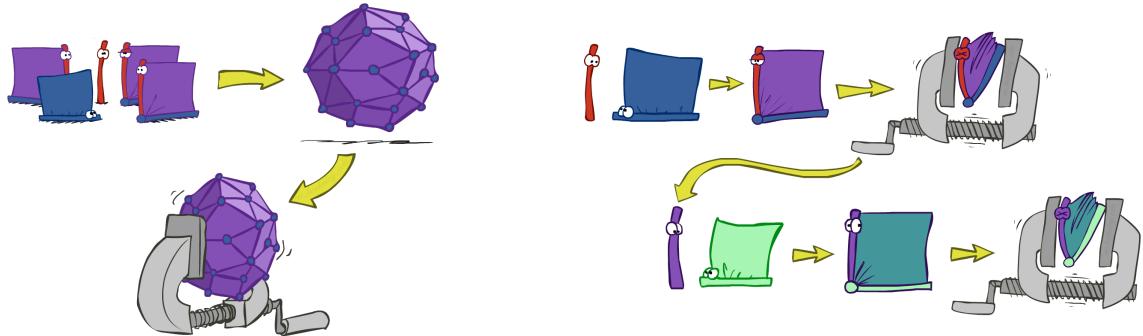
## Inference by Enumeration?



# Inference by Enumeration vs. Variable Elimination

---

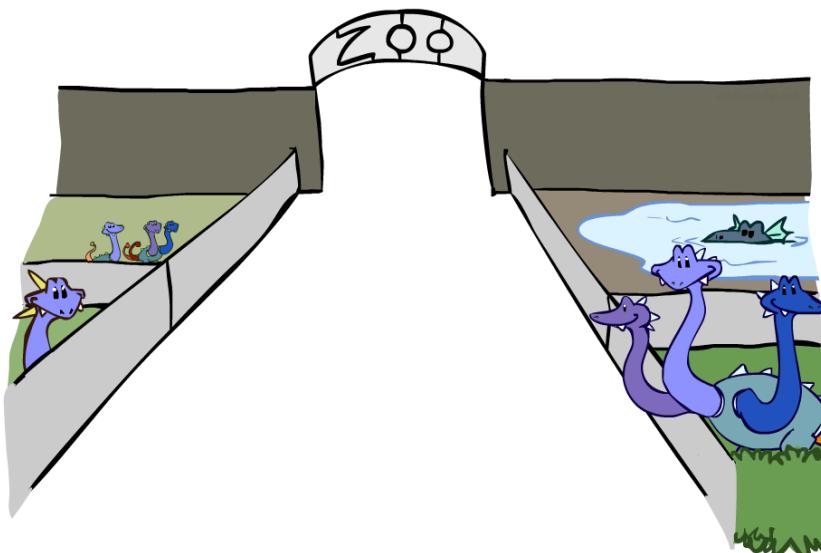
- Why is inference by enumeration so slow?
  - You join up the whole joint distribution before you sum out the hidden variables
- Idea: interleave joining and marginalizing!
  - Called “Variable Elimination”
  - Still NP-hard, but usually much faster than inference by enumeration



- First we'll need some new notation: factors

## Factor Zoo

---



# Factor Zoo I

- Joint distribution:  $P(X, Y)$

- Entries  $P(x, y)$  for all  $x, y$
- Sums to 1

$P(T, W)$

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

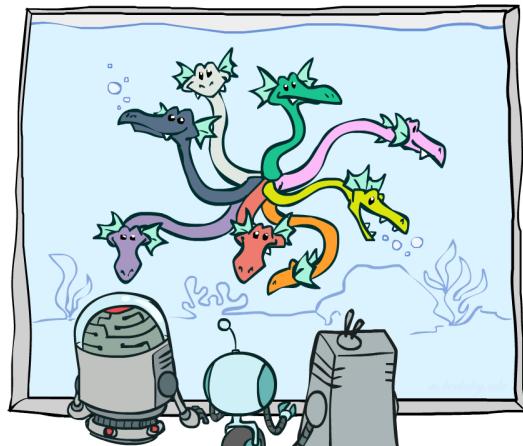
- Selected joint:  $P(x, Y)$

- A slice of the joint distribution
- Entries  $P(x, y)$  for fixed  $x$ , all  $y$
- Sums to  $P(x)$

$P(cold, W)$

T	W	P
cold	sun	0.2
cold	rain	0.3

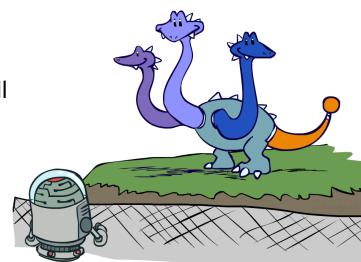
- Number of capitals = dimensionality of the table



# Factor Zoo II

- Single conditional:  $P(Y | x)$

- Entries  $P(y | x)$  for fixed  $x$ , all  $y$
- Sums to 1



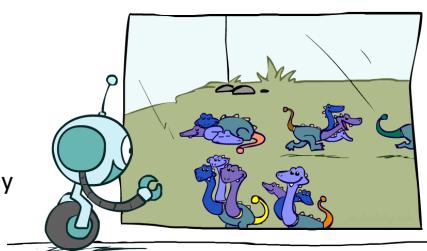
$P(W|cold)$

T	W	P
cold	sun	0.4
cold	rain	0.6

- Family of conditionals:

$P(Y | X)$

- Multiple conditionals
- Entries  $P(y | x)$  for all  $x, y$
- Sums to  $|X|$



$P(W|T)$

T	W	P
hot	sun	0.8
hot	rain	0.2
cold	sun	0.4
cold	rain	0.6

$\} P(W|hot)$

$\} P(W|cold)$

## Factor Zoo III

---

- Specified family:  $P(y | X)$

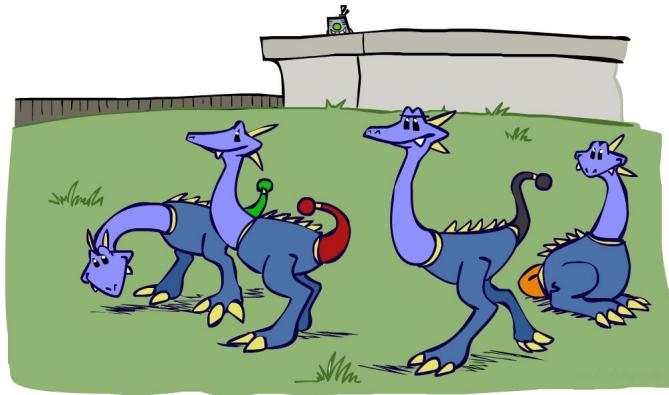
- Entries  $P(y | x)$  for fixed  $y$ ,  
but for all  $x$
- Sums to ... who knows!

$$P(\text{rain}|T)$$

T	W	P
hot	rain	0.2
cold	rain	0.6

$P(\text{rain}|\text{hot})$

$P(\text{rain}|\text{cold})$

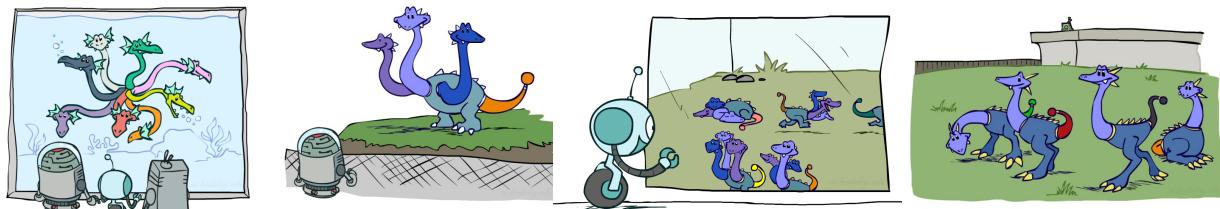


## Factor Zoo Summary

---

- In general, when we write  $P(Y_1 \dots Y_N | X_1 \dots X_M)$

- It is a “factor,” a multi-dimensional array
- Its values are  $P(y_1 \dots y_N | x_1 \dots x_M)$
- Any assigned (=lower-case) X or Y is a dimension missing (selected) from the array



## Example: Traffic Domain

---

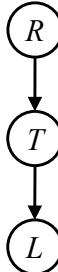
- Random Variables

- R: Raining
- T: Traffic
- L: Late for class!

$$P(L) = ?$$

$$= \sum_{r,t} P(r,t,L)$$

$$= \sum_{r,t} P(r)P(t|r)P(L|t)$$



$$P(R)$$

+r	0.1
-r	0.9

$$P(T|R)$$

+r	+t	0.8
+r	-t	0.2
-r	+t	0.1
-r	-t	0.9

$$P(L|T)$$

+t	+l	0.3
+t	-l	0.7
-t	+l	0.1
-t	-l	0.9

## Inference by Enumeration: Procedural Outline

---

- Track objects called **factors**
- Initial factors are local CPTs (one per node)

$$P(R)$$

+r	0.1
-r	0.9

$$P(T|R)$$

+r	+t	0.8
+r	-t	0.2
-r	+t	0.1
-r	-t	0.9

$$P(L|T)$$

+t	+l	0.3
+t	-l	0.7
-t	+l	0.1
-t	-l	0.9

- Any known values are selected

- E.g. if we know  $L = +\ell$ , the initial factors are

$$P(R)$$

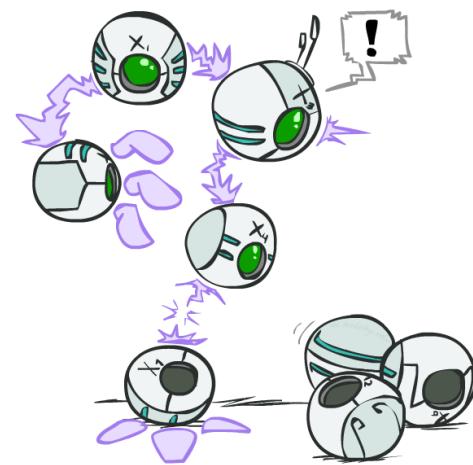
+r	0.1
-r	0.9

$$P(T|R)$$

+r	+t	0.8
+r	-t	0.2
-r	+t	0.1
-r	-t	0.9

$$P(+\ell|T)$$

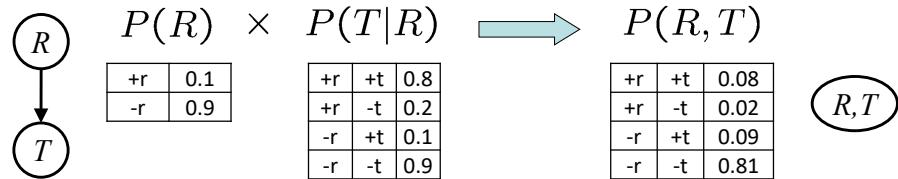
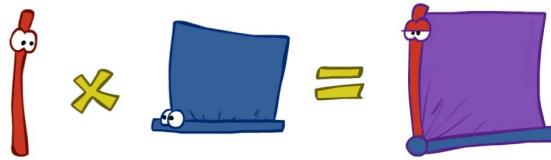
+t	+l	0.3
+t	-l	0.1
-t	+l	0.1
-t	-l	0.9



- Procedure: Join all factors, eliminate all hidden variables, normalize

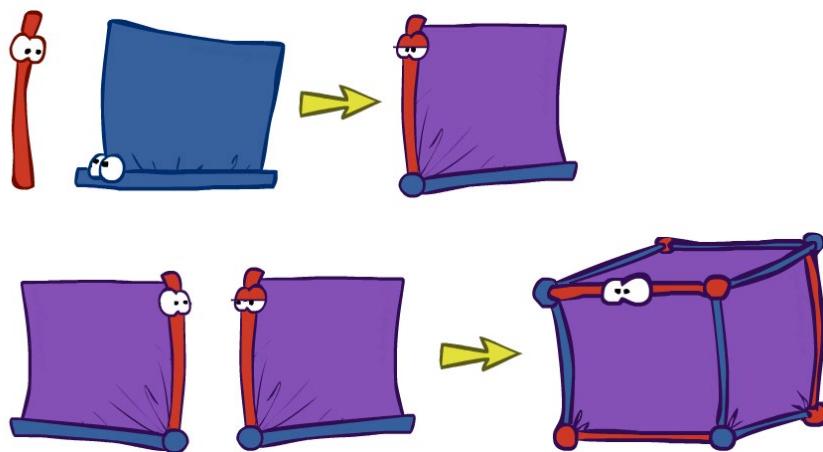
# Operation 1: Join Factors

- First basic operation: **joining factors**
- Combining factors:
  - Just like a database join
  - Get all factors over the joining variable
  - Build a new factor over the union of the variables involved
- Example: Join on R

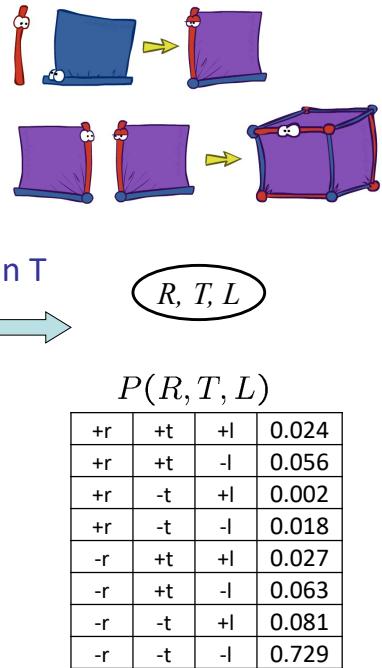


- Computation for each entry: pointwise products     $\forall r, t : P(r, t) = P(r) \cdot P(t|r)$

## Example: Multiple Joins



## Example: Multiple Joins



## Operation 2: Eliminate

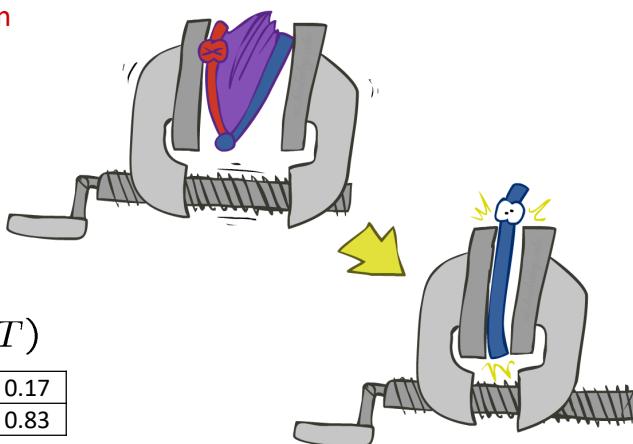
- Second basic operation: **marginalization**
- Take a factor and sum out a variable
  - Shrinks a factor to a smaller one
  - A **projection** operation
- Example:

$P(R, T)$

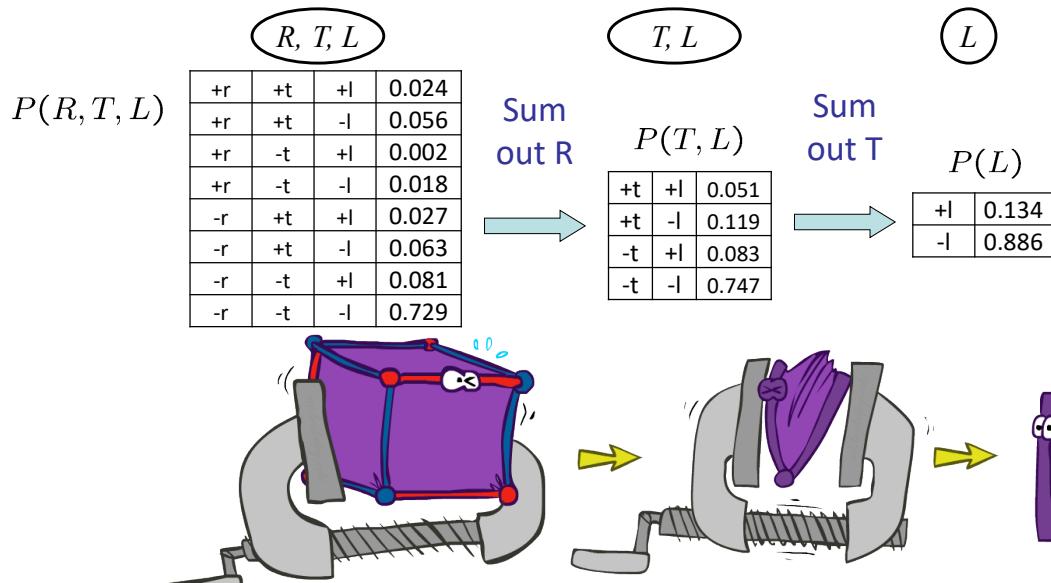
+r	+t	0.08
+r	-t	0.02
-r	+t	0.09
-r	-t	0.81

sum  $R$

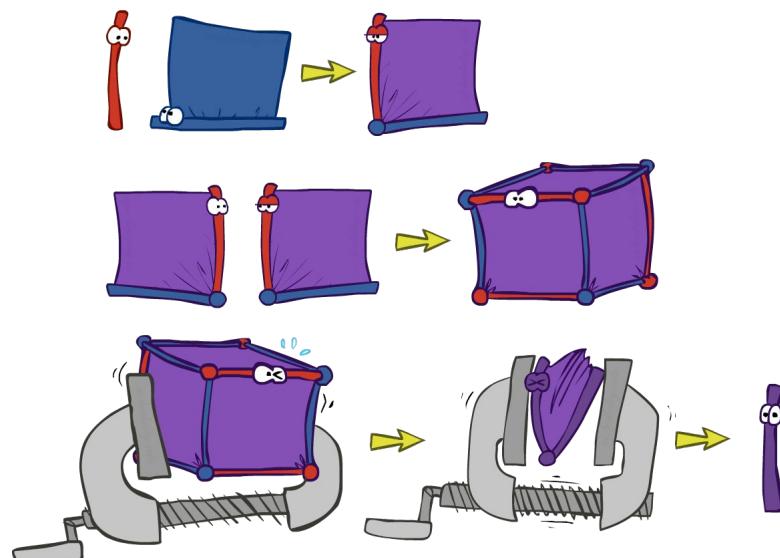
+t	0.17
-t	0.83



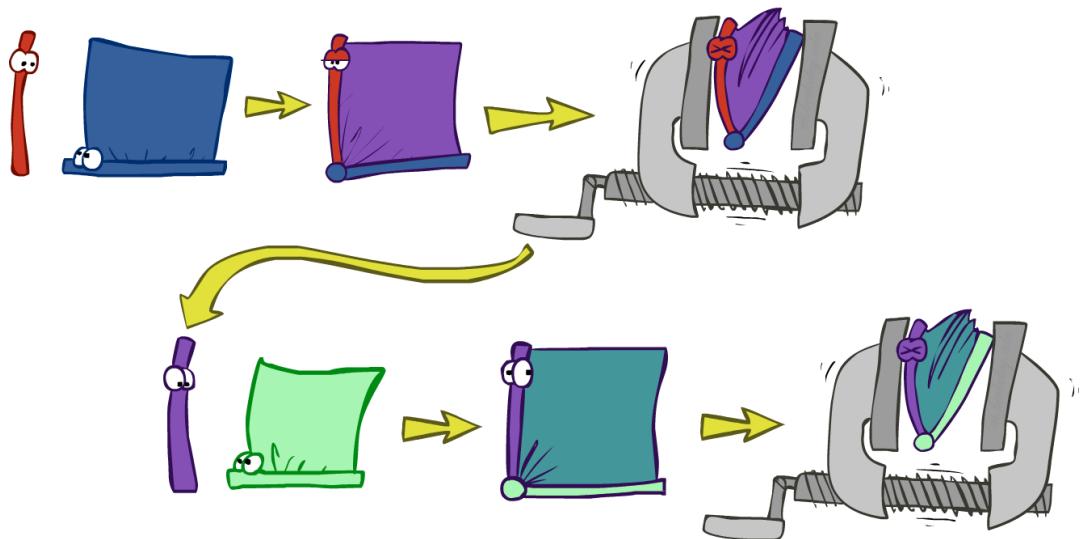
## Multiple Elimination



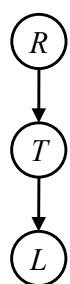
Thus Far: Multiple Join, Multiple Eliminate (= Inference by Enumeration)



## Marginalizing Early (= Variable Elimination)



## Traffic Domain



$$P(L) = ?$$

### Inference by Enumeration

$$= \sum_t \sum_r P(L|t) P(r) P(t|r)$$

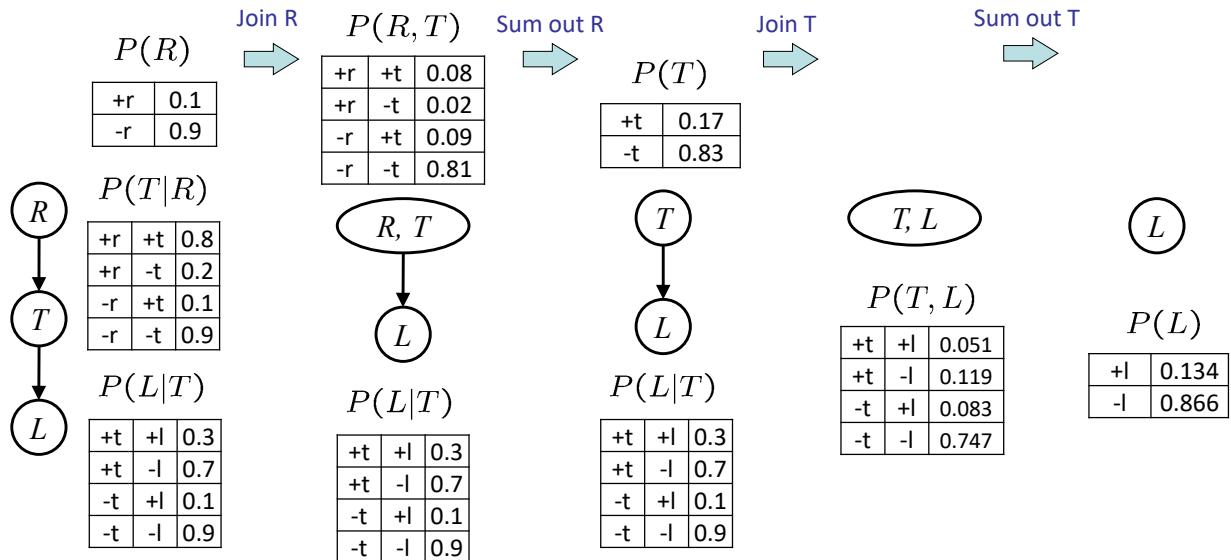
Join on r  
Join on t  
Eliminate r  
Eliminate t

### Variable Elimination

$$= \sum_t P(L|t) \sum_r P(r) P(t|r)$$

Join on r  
Eliminate r  
Join on t  
Eliminate t

# Marginalizing Early! (aka VE)



## Evidence

- If evidence, start with factors that select that evidence

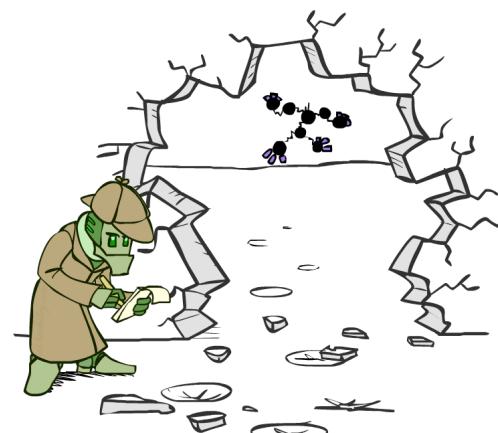
- No evidence uses these initial factors:

$P(R)$	$P(T R)$	$P(L T)$
+r   0.1	+r   +t   0.8 +r   -t   0.2 -r   +t   0.1 -r   -t   0.9	+t   +l   0.3 +t   -l   0.7 -t   +l   0.1 -t   -l   0.9
-r   0.9		

- Computing  $P(L|r)$  the initial factors become:

$P(+r)$	$P(T +r)$	$P(L T)$
+r   0.1	+r   +t   0.8 +r   -t   0.2	+t   +l   0.3 +t   -l   0.7 -t   +l   0.1 -t   -l   0.9
-r   0.9		

- We eliminate all vars other than query + evidence



# Evidence II

---

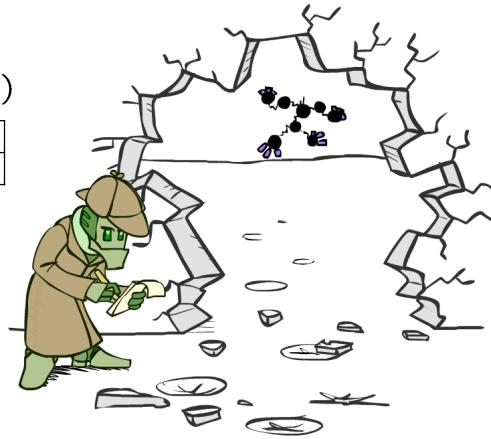
- Result will be a selected joint of query and evidence

- E.g. for  $P(L | +r)$ , we would end up with:

$P(+r, L)$		
+r	+l	0.026
+r	-l	0.074

Normalize 

$P(L   +r)$	
+l	0.26
-l	0.74

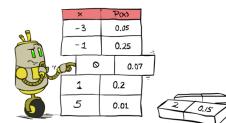


- To get our answer, just normalize this!
- That's it!

# General Variable Elimination

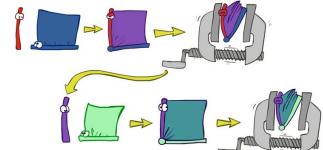
---

- Query:  $P(Q | E_1 = e_1, \dots, E_k = e_k)$



- Start with initial factors:

- Local CPTs (but instantiated by evidence)



- While there are still hidden variables (not Q or evidence):

- Pick a hidden variable H
- Join all factors mentioning H
- Eliminate (sum out) H

- Join all remaining factors and normalize

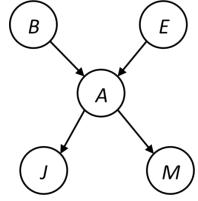
$$\{ * \text{ [blue]} = \text{ [purple]} \} \times \frac{1}{Z}$$

## Example

---

$$P(B|j, m) \propto P(B, j, m)$$

$P(B)$	$P(E)$	$P(A B, E)$	$P(j A)$	$P(m A)$
--------	--------	-------------	----------	----------



Choose A

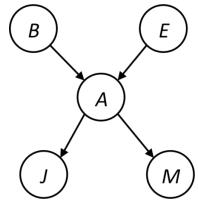
$$\begin{array}{c} P(A|B, E) \\ P(j|A) \\ P(m|A) \end{array} \quad \times \quad P(j, m, A|B, E) \quad \sum \quad P(j, m|B, E)$$

$P(B)$	$P(E)$	$P(j, m B, E)$
--------	--------	----------------

## Example

---

$P(B)$	$P(E)$	$P(j, m B, E)$
--------	--------	----------------



Choose E

$$\begin{array}{c} P(E) \\ P(j, m|B, E) \end{array} \quad \times \quad P(j, m, E|B) \quad \sum \quad P(j, m|B)$$

$P(B)$	$P(j, m B)$
--------	-------------

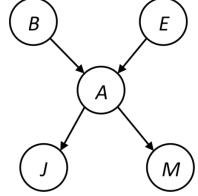
Finish with B

$$\begin{array}{c} P(B) \\ P(j, m|B) \end{array} \quad \times \quad P(j, m, B) \quad \xrightarrow{\text{Normalize}} \quad P(B|j, m)$$

## Same Example in Equations

$$P(B|j, m) \propto P(B, j, m)$$

$P(B)$	$P(E)$	$P(A B, E)$	$P(j A)$	$P(m A)$
--------	--------	-------------	----------	----------



$$\begin{aligned}
 P(B|j, m) &\propto P(B, j, m) \\
 &= \sum_{e,a} P(B, j, m, e, a) && \text{marginal obtained from joint by summing out} \\
 &= \sum_{e,a} P(B)P(e)P(a|B, e)P(j|a)P(m|a) && \text{use Bayes' net joint distribution expression} \\
 &= \sum_e P(B)P(e) \sum_a P(a|B, e)P(j|a)P(m|a) && \text{use } x^*(y+z) = xy + xz \\
 &= \sum_e P(B)P(e)f_1(B, e, j, m) && \text{joining on } a, \text{ and then summing out gives } f_1 \\
 &= P(B) \sum_e P(e)f_1(B, e, j, m) && \text{use } x^*(y+z) = xy + xz \\
 &= P(B)f_2(B, j, m) && \text{joining on } e, \text{ and then summing out gives } f_2
 \end{aligned}$$

All we are doing is exploiting  $uwv + uwz + uxy + uxz + vwy + vwz + vxy + vxz = (u+v)(w+x)(y+z)$  to improve computational efficiency!

## Another Variable Elimination Example

Query:  $P(X_3|Y_1 = y_1, Y_2 = y_2, Y_3 = y_3)$

Start by inserting evidence, which gives the following initial factors:

$$p(Z)p(X_1|Z)p(X_2|Z)p(X_3|Z)p(y_1|X_1)p(y_2|X_2)p(y_3|X_3)$$

Eliminate  $X_1$ , this introduces the factor  $f_1(Z, y_1) = \sum_{x_1} p(x_1|Z)p(y_1|x_1)$ , and we are left with:

$$p(Z)f_1(Z, y_1)p(X_2|Z)p(X_3|Z)p(y_2|X_2)p(y_3|X_3)$$

Eliminate  $X_2$ , this introduces the factor  $f_2(Z, y_2) = \sum_{x_2} p(x_2|Z)p(y_2|x_2)$ , and we are left with:

$$p(Z)f_1(Z, y_1)f_2(Z, y_2)p(X_3|Z)p(y_3|X_3)$$

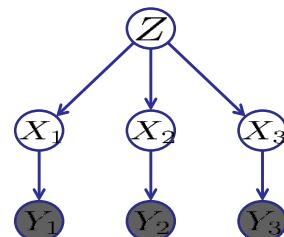
Eliminate  $Z$ , this introduces the factor  $f_3(y_1, y_2, X_3) = \sum_z p(z)f_1(z, y_1)f_2(z, y_2)p(X_3|z)$ , and we are left:

$$p(y_3|X_3), f_3(y_1, y_2, X_3)$$

No hidden variables left. Join the remaining factors to get:

$$f_4(y_1, y_2, y_3, X_3) = P(y_3|X_3)f_3(y_1, y_2, X_3).$$

Normalizing over  $X_3$  gives  $P(X_3|y_1, y_2, y_3)$ .

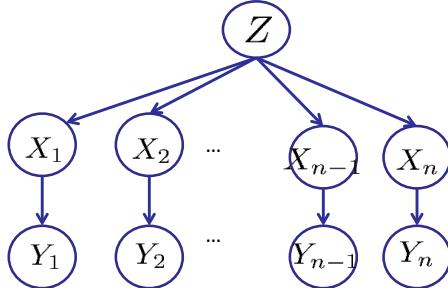


Computational complexity critically depends on the largest factor being generated in this process. Size of factor = number of entries in table. In example above (assuming binary) all factors generated are of size 2 --- as they all only have one variable ( $Z$ ,  $Z$ , and  $X_3$  respectively).

# Variable Elimination Ordering

---

- For the query  $P(X_n | y_1, \dots, y_n)$  work through the following two different orderings as done in previous slide:  $Z, X_1, \dots, X_{n-1}$  and  $X_1, \dots, X_{n-1}, Z$ . What is the size of the maximum factor generated for each of the orderings?



- Answer:  $2^{n+1}$  versus  $2^2$  (assuming binary)
- In general: the ordering can greatly affect efficiency.

## VE: Computational and Space Complexity

---

- The computational and space complexity of variable elimination is determined by the largest factor
- The elimination ordering can greatly affect the size of the largest factor.
  - E.g., previous slide's example  $2^n$  vs. 2
- Does there always exist an ordering that only results in small factors?
  - No!

# Worst Case Complexity?

---

- **CSP:**

$$(x_1 \vee x_2 \vee \neg x_3) \wedge (\neg x_1 \vee x_3 \vee \neg x_4) \wedge (x_2 \vee \neg x_2 \vee x_4) \wedge (\neg x_3 \vee \neg x_4 \vee \neg x_5) \wedge (x_2 \vee x_5 \vee x_7) \wedge (x_4 \vee x_5 \vee x_6) \wedge (\neg x_5 \vee x_6 \vee \neg x_7) \wedge (\neg x_5 \vee \neg x_6 \vee x_7)$$

$$P(X_i = 0) = P(X_i = 1) = 0.5$$

$$Y_1 = X_1 \vee X_2 \vee \neg X_3$$

$$\dots$$

$$Y_8 = \neg X_5 \vee X_6 \vee X_7$$

$$Y_{1,2} = Y_1 \wedge Y_2$$

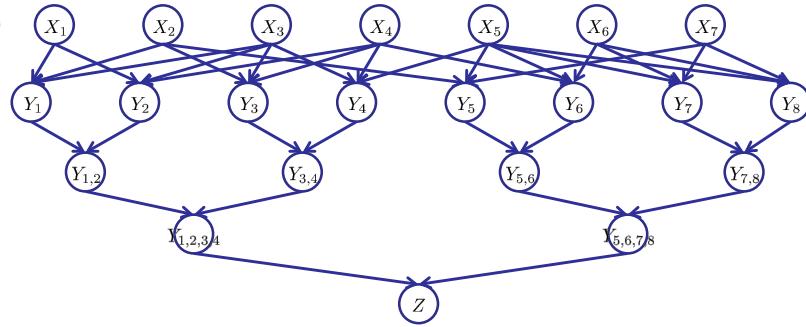
$$\dots$$

$$Y_{7,8} = Y_7 \wedge Y_8$$

$$Y_{1,2,3,4} = Y_{1,2} \wedge Y_{3,4}$$

$$Y_{5,6,7,8} = Y_{5,6} \wedge Y_{7,8}$$

$$Z = Y_{1,2,3,4} \wedge Y_{5,6,7,8}$$



- If we can answer  $P(z)$  equal to zero or not, we answered whether the 3-SAT problem has a solution.
- Hence inference in Bayes' nets is NP-hard. No known efficient probabilistic inference in general.

# Polytrees

---

- A polytree is a directed graph with no undirected cycles
- For poly-trees you can always find an ordering that is efficient
  - Try it!!
- Cut-set conditioning for Bayes' net inference
  - Choose set of variables such that if removed only a polytree remains
  - Exercise: Think about how the specifics would work out!

# Bayes' Nets

---

- ✓ Representation
- ✓ Conditional Independences
  - Probabilistic Inference
    - ✓ Enumeration (exact, exponential complexity)
    - ✓ Variable elimination (exact, worst-case exponential complexity, often better)
    - ✓ Inference is NP-complete
      - Sampling (approximate)
  - Learning Bayes' Nets from Data