

Bayes Nets

With slides from Dan Klein and Stuart Russell

Today

- Bayes Nets
 - Bayesian reasoning (recap)
 - Representation
 - Inference

The Product Rule

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

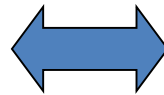
- Example:

$P(W)$

R	P
sun	0.8
rain	0.2

$P(D|W)$

D	W	P
wet	sun	0.1
dry	sun	0.9
wet	rain	0.7
dry	rain	0.3



$P(D, W)$

D	W	P
wet	sun	
dry	sun	
wet	rain	
dry	rain	

The Chain Rule: Inverse of Product 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 true?
 - Recursive decomposition using product rule
 - $P(x_1, x_2, x_3) = p(x_3 | x_1, x_2) * p(x_1, x_2)$
 $= p(x_3 | x_1, x_2) * p(x_2 | x_1) * p(x_1)$

Bayes' Rule

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

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

- Dividing, we get:

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

- Why is this 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 (e.g. ASR, MT)
- In the running for most important AI, ML, DM equation!**

That's my rule!



Inference with Bayes' Rule

- Example: Diagnostic probability from causal probability:

- Example: $P(\text{cause}|\text{effect}) = \frac{P(\text{effect}|\text{cause})P(\text{cause})}{P(\text{effect})}$

- M: meningitis, S: stiff neck

$$\left. \begin{array}{l} P(+m) = 0.0001 \\ P(+s|+m) = 0.8 \\ P(+s|-m) = 0.01 \end{array} \right\} \begin{array}{l} \text{Example} \\ \text{givens} \end{array}$$

$$P(+m|+s) =$$

Inference with Bayes' Rule

- Example: Diagnostic probability from causal probability:

- Example: $P(\text{cause}|\text{effect}) = \frac{P(\text{effect}|\text{cause})P(\text{cause})}{P(\text{effect})}$

- M: meningitis, S: stiff neck

$$\left. \begin{array}{l} P(+m) = 0.0001 \\ P(+s|+m) = 0.8 \\ P(+s|-m) = 0.01 \end{array} \right\} \begin{array}{l} \text{Example} \\ \text{givens} \end{array}$$

$$P(+m|+s) =$$

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

Exercise: Inference with Bayes' Rule

- Given:

$$P(W)$$

W	P
sun	0.8
rain	0.2

$$P(D|W)$$

D	W	P
wet	sun	0.1
dry	sun	0.9
wet	rain	0.7
dry	rain	0.3

- What is $P(W \mid \text{dry})$?

Solution: on board

W	P
sun	
rain	

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

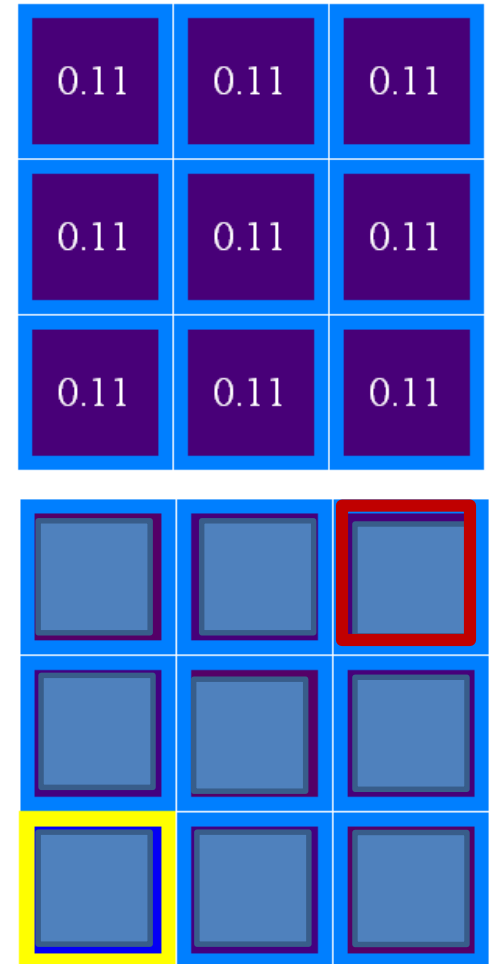
Exercise: Bayes Ghost Localization

- Setup:

- Prior distribution over ghost location: $P(G)$ = uniform (on right)
- R = reading color measured at $(1,1)$ = Yellow
- Sensor reading model: $P(R | G)$

$P(\text{red} 4)$	$P(\text{orange} 4)$	$P(\text{yellow} 4)$	$P(\text{green} 4)$
0.05	0.15	0.5	0.3

- What is probability of ghost at $(3,3)$?



Hands-on Example: Ghost Localization

- Setup:

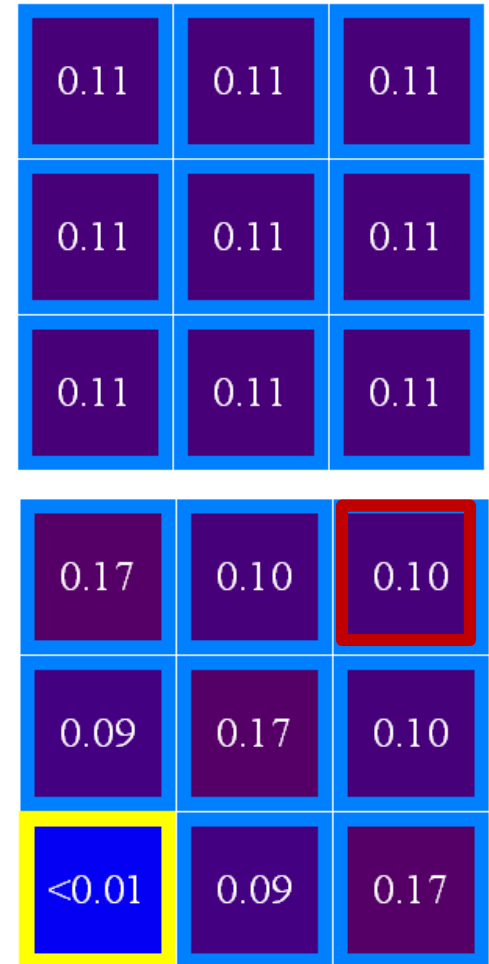
- Prior distribution over ghost location: $P(G)$ = uniform (on right)
- R = reading color measured at (1,1) = Yellow
- Sensor reading model: $P(R | G)$

$P(\text{red} 4)$	$P(\text{orange} 4)$	$P(\text{yellow} 4)$	$P(\text{green} 4)$
0.05	0.15	0.5	0.3

- What is probability of ghost at (3,3)?

- Answer:

$$\begin{aligned} p(3,3 | R) &= p(R=\text{yel} | g=3,3) * p(g=3,3) \\ &\sim 0.5 * 0.11 \text{ (before normalization)} \\ &0.055 \end{aligned}$$



Independence

- Two variables are *independent* if:

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

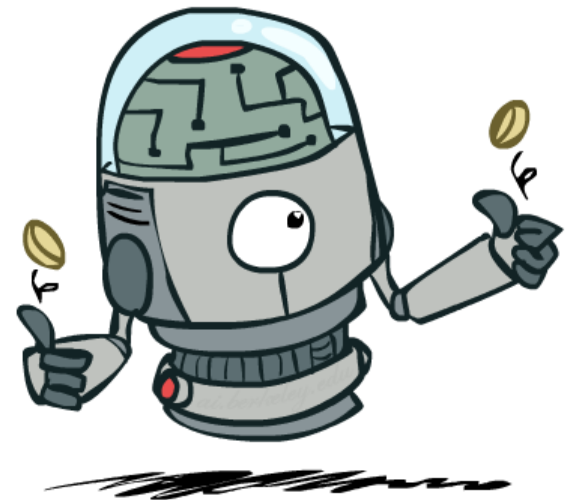
- This says that their joint distribution *factors* into a product two simpler distributions
- Another form:

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

- 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:

$$P(X_1)$$

H	0.5
T	0.5

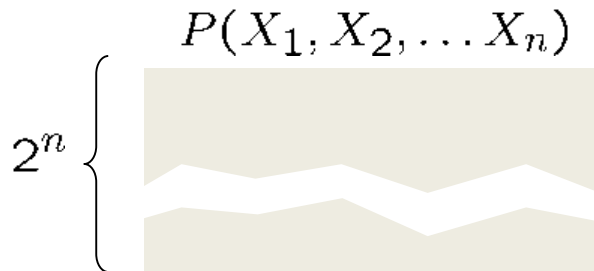
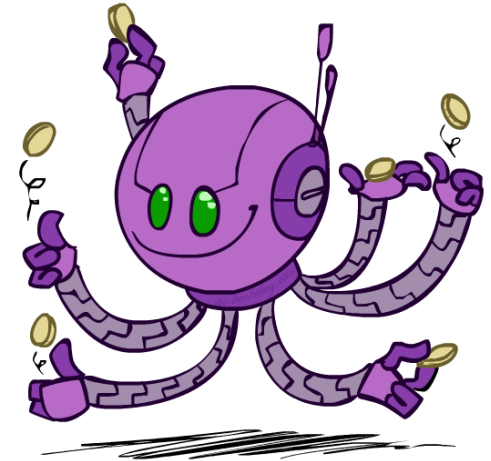
$$P(X_2)$$

H	0.5
T	0.5

...

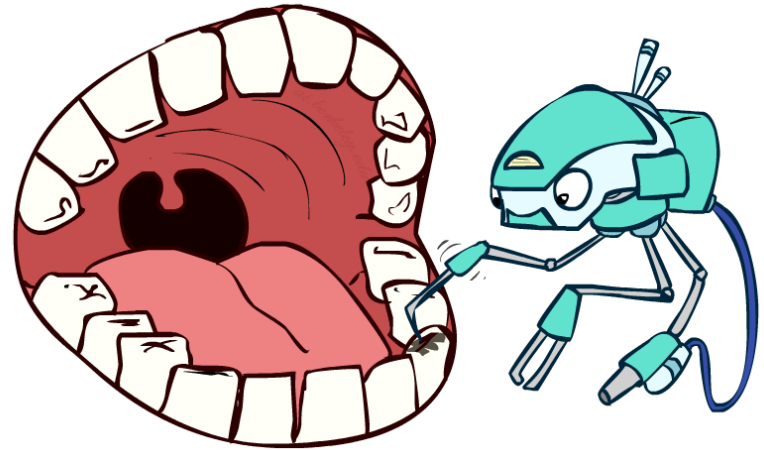
$$P(X_n)$$

H	0.5
T	0.5



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:
 - $P(+\text{catch} \mid +\text{toothache}, -\text{cavity}) = P(+\text{catch} \mid -\text{cavity})$
- 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.

- X is conditionally independent of Y given Z

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

if and only if:

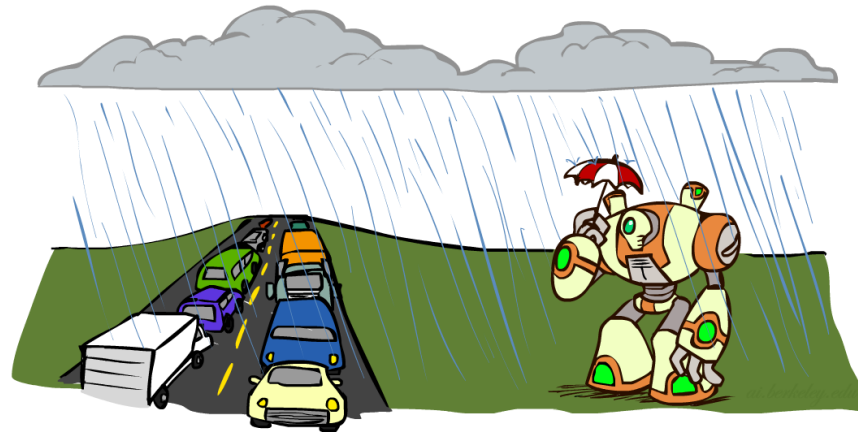
$$\forall x, y, z : P(x, y | z) = P(x | z)P(y | z)$$

or, equivalently, if and only if

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

Which variables conditionally independent? T, U, R?

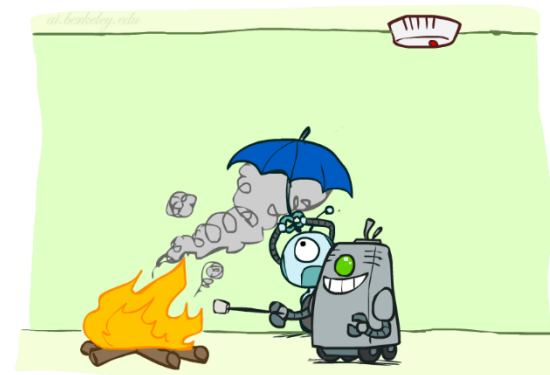
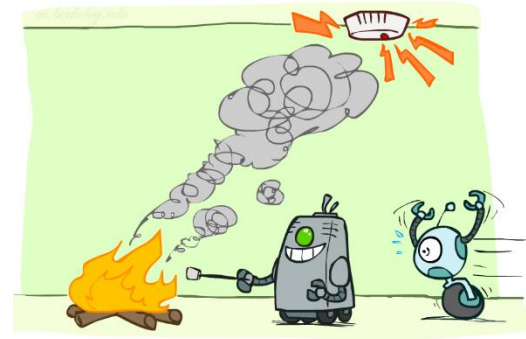
- What about this domain:
 - Traffic
 - Umbrella
 - Raining



Which variables conditionally independent? F, S, A?

- 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:

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

- With assumption of conditional independence:

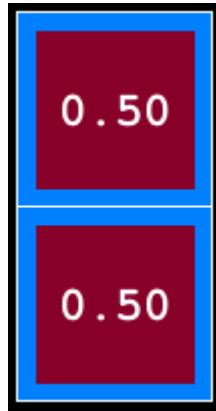
$$P(\text{Traffic, Rain, Umbrella}) = P(\text{Rain})P(\text{Traffic}|\text{Rain})P(\text{Umbrella}|\text{Rain})$$

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



Ghostbusters Chain Rule

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



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

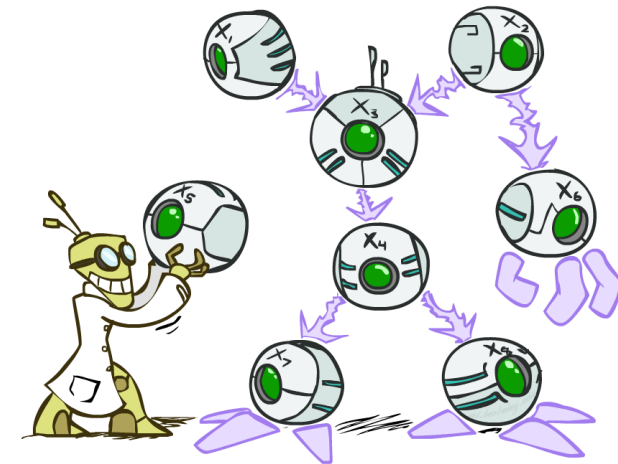
T	B	G	P(T,B,G)
+t	+b	+g	0.16
+t	+b	-g	0.16
+t	-b	+g	0.24
+t	-b	-g	0.04
-t	+b	+g	0.04
-t	+b	-g	0.24
-t	-b	+g	0.06
-t	-b	-g	0.06

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



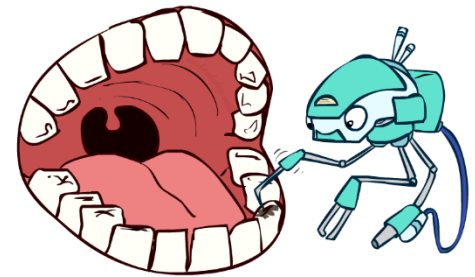
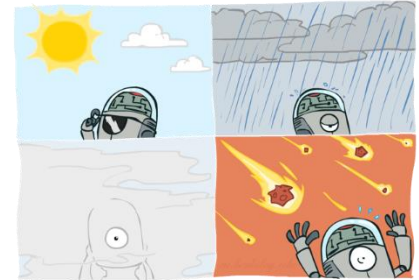
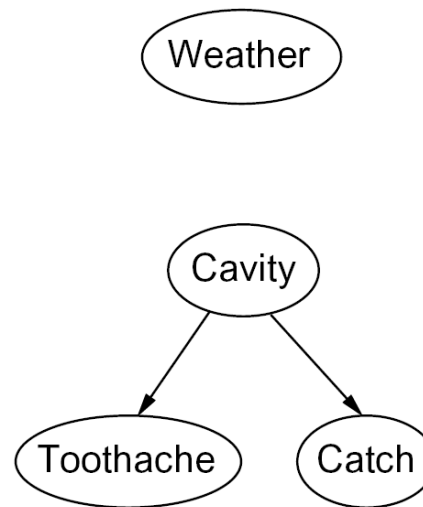
Bayes' Nets: Big Picture

- Two problems with using full joint distribution tables as our probabilistic models:
 - Unless there are only a few variables, joint probability table is MUCH 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 only local distributions (conditional probabilities)
 - More generally: kind of a **graphical model**
 - We describe how variables locally interact
 - Local interactions chain together to give global, indirect interactions



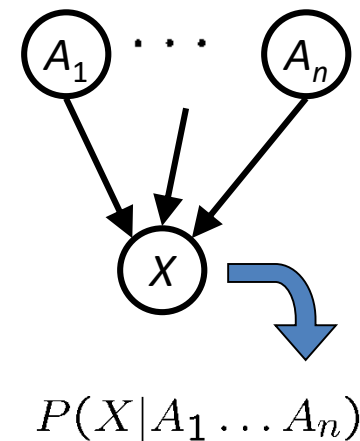
Graphical Model Notation

- 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!)



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



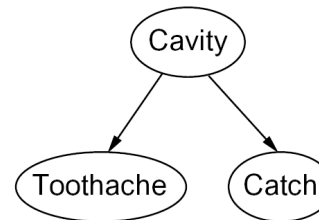
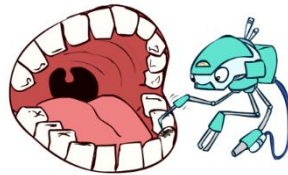
A Bayes net = Topology (graph) + Local Conditional Probabilities

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:

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

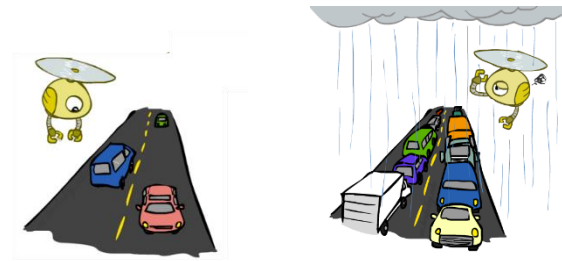
- Example:



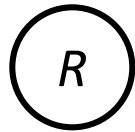
$$P(+cavity, +catch, -toothache)$$

Example: Traffic

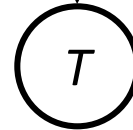
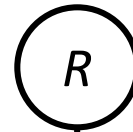
- Variables:
 - R: It rains
 - T: There is traffic



- Model 1: independence



- Model 2: rain causes traffic



- Which model is better for a driving agent?

Probabilities in BNs (2)

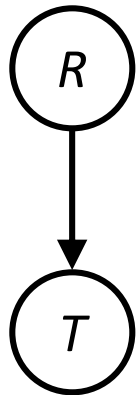
- Why are we guaranteed that setting

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

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: Traffic

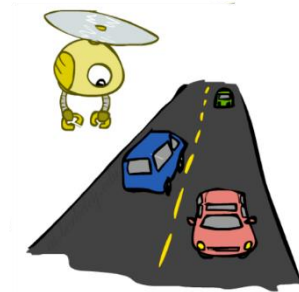

$$P(R)$$

$+r$	$1/4$
$-r$	$3/4$

$$P(T|R)$$

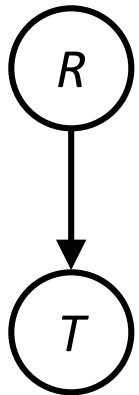
$+r$	$+t$	$3/4$
	$-t$	$1/4$
$-r$	$+t$	$1/2$
	$-t$	$1/2$

$$P(+r, -t) =$$



Example: Traffic

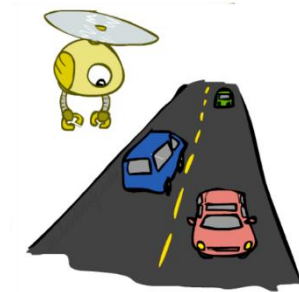
$$P(x_1, x_2, \dots, x_n) = \prod_{i=1}^n P(x_i | \text{parents}(X_i))$$



$P(R)$	
$+r$	$1/4$
$-r$	$3/4$

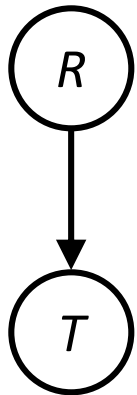
$P(T R)$		
$+r$	$+t$	$3/4$
	$-t$	$1/4$
$-r$	$+t$	$1/2$
	$-t$	$1/2$

$$P(+r, -t) =$$



Example: Traffic

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



$P(R)$

$+r$	$1/4$
$-r$	$3/4$

$P(T|R)$

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


$$P(+r, -t) = 1/4 * 1/4 = 1/16 (0.06)$$



CPT Comments

- Probability of Node=false not given, can subtract from 1:

B	E	P(A=T)
T	T	.95

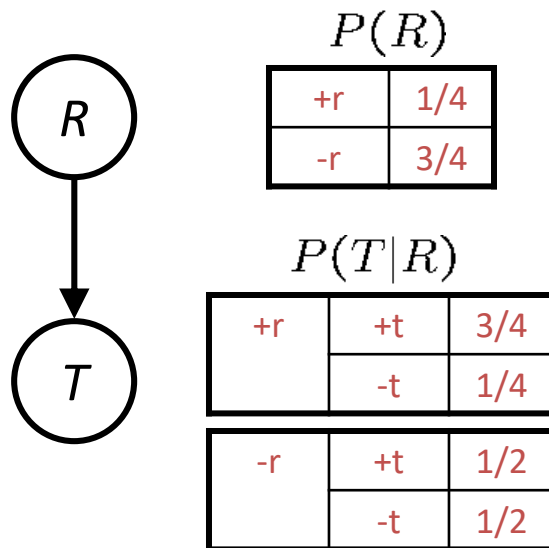


B	E	P(A=F)
T	T	.05

- CPT rows do not need to add up to one – they are NOT NORMALIZED. (convenient for inference)
- Example requires 10 parameters rather than $2^5 - 1 = 31$ for specifying the full joint distribution.
- Number of parameters in the CPT for a node is exponential in the number of parents (fan-in).

BNs = Causality?

- Causal direction: rain causes traffic

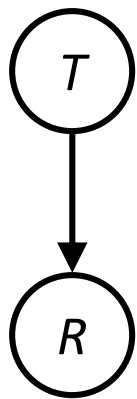


Joint probability distribution $P(T, R)$:

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

Example: Reverse Traffic

- Reverse causality?
“traffic causes rain”


$$P(T)$$

+t	9/16
-t	7/16

$$P(R|T)$$

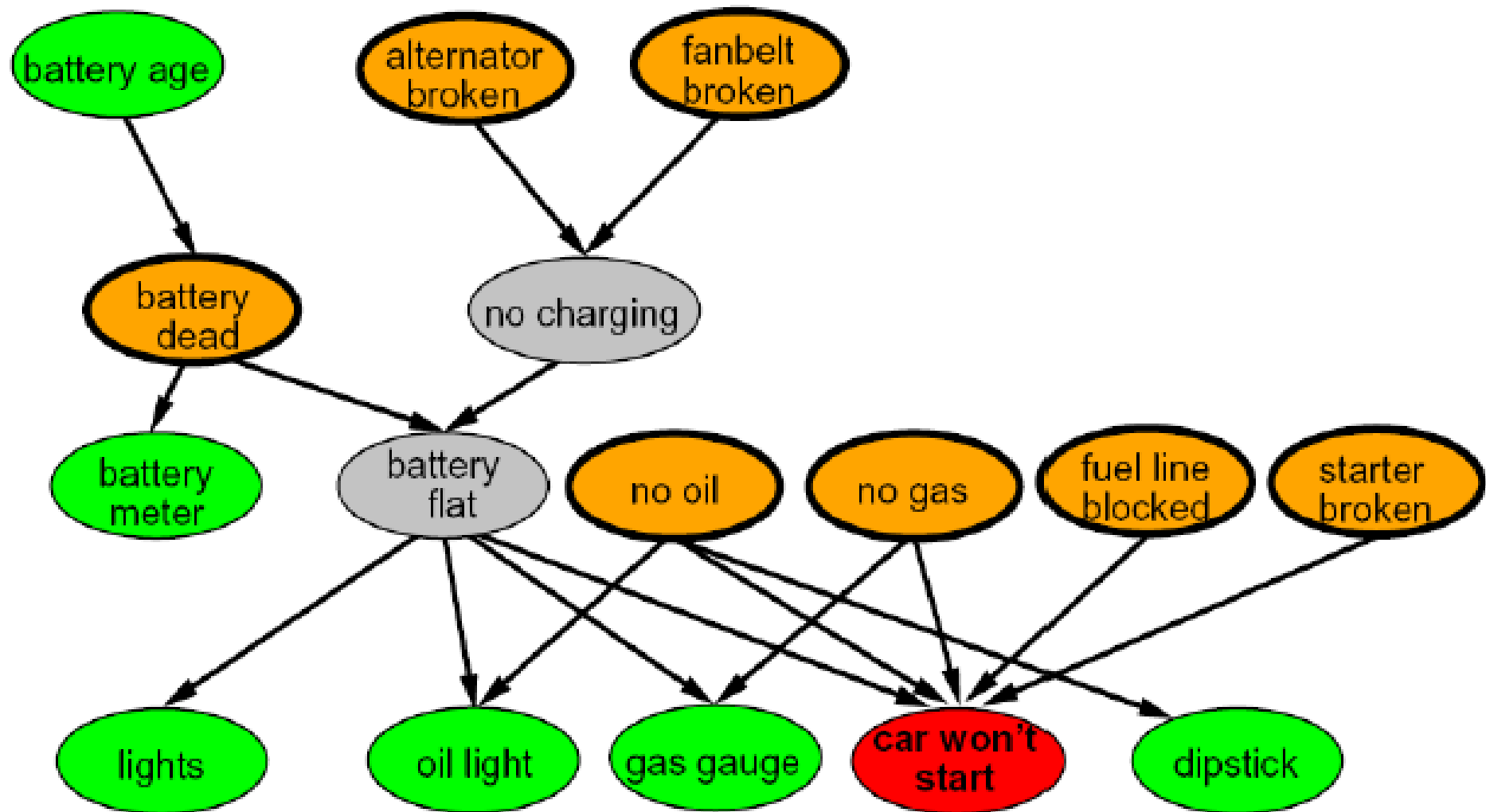
+t	+r	1/3
	-r	2/3
-t	+r	1/7
	-r	6/7

A cartoon illustration showing a red convertible car with a yellow bird-like character driving in the rain. A teal car is following behind it. Rain clouds and raindrops are depicted above the red car.

$$P(T, R)$$

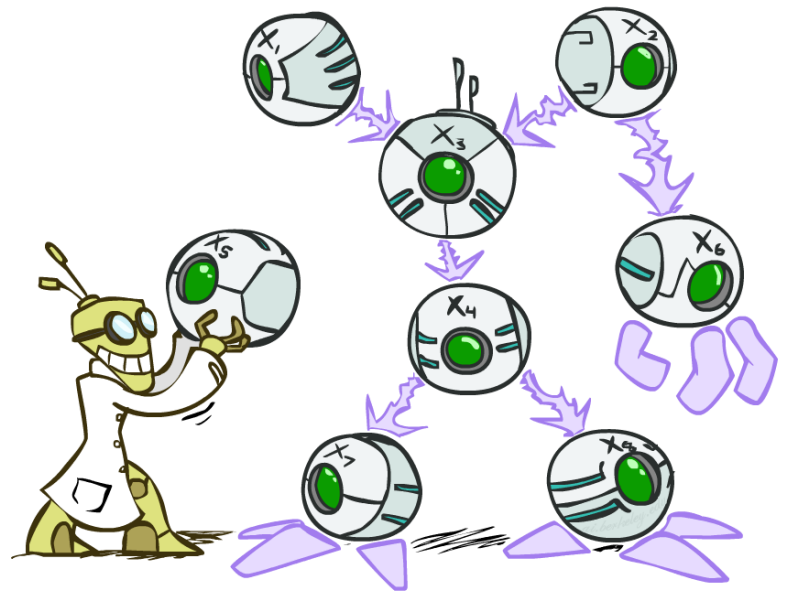
+r	+t	3/16
+r	-t	1/16
-r	+t	6/16
-r	-t	6/16

Example Bayes' Net: Car



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)

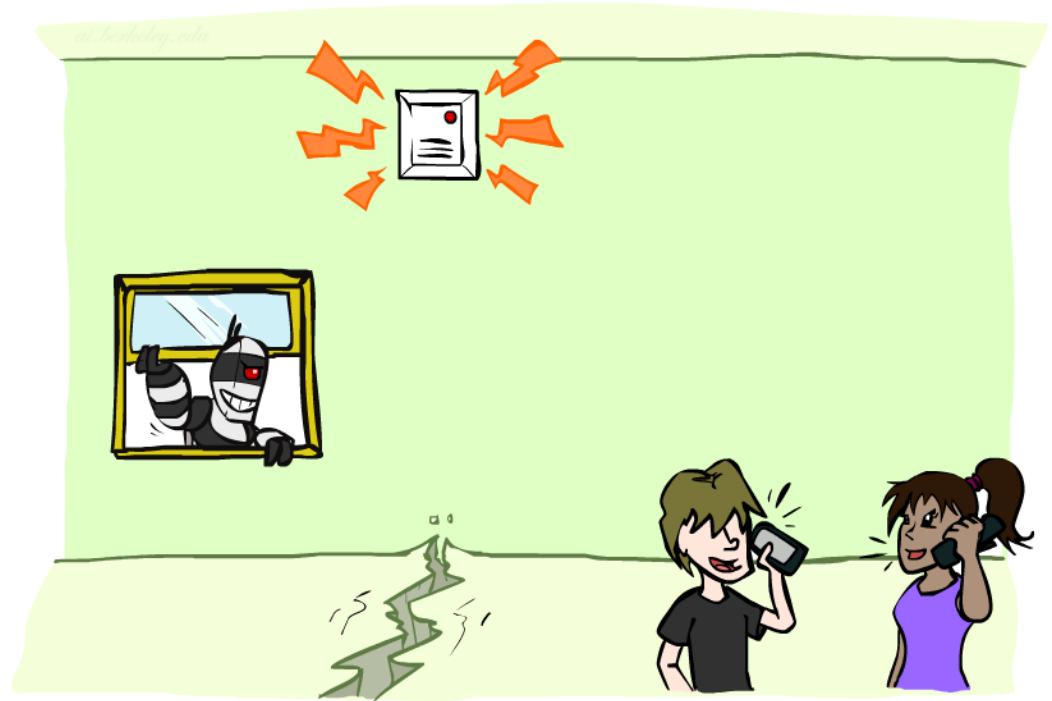


<https://www.bayesserver.com/Live.aspx>

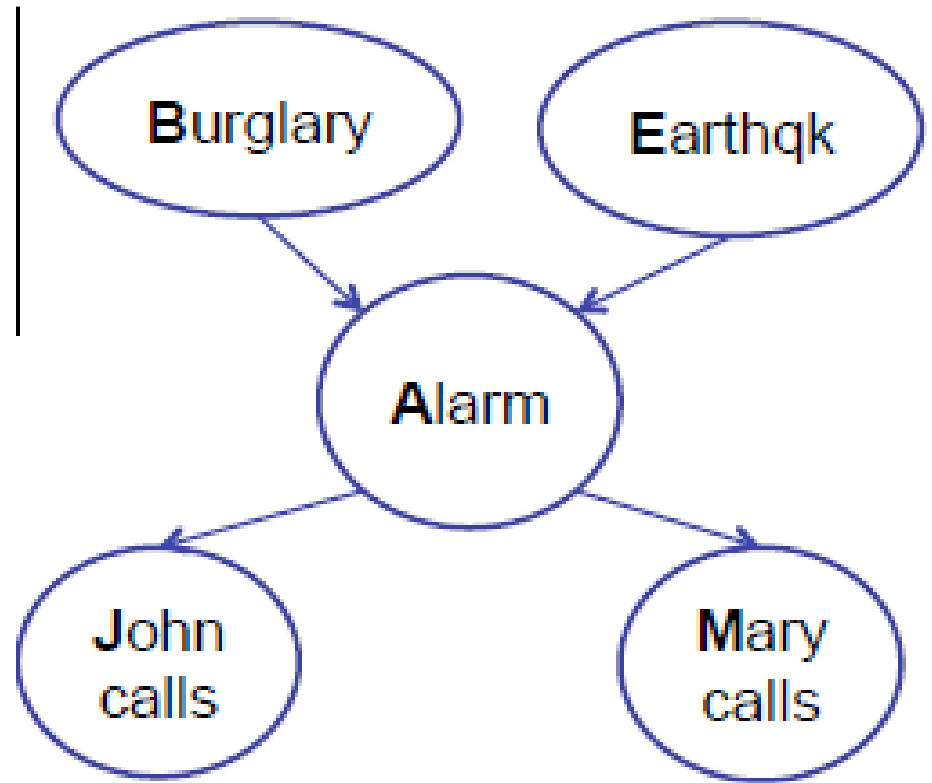
BN Inference: Alarm Network

- Variables

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

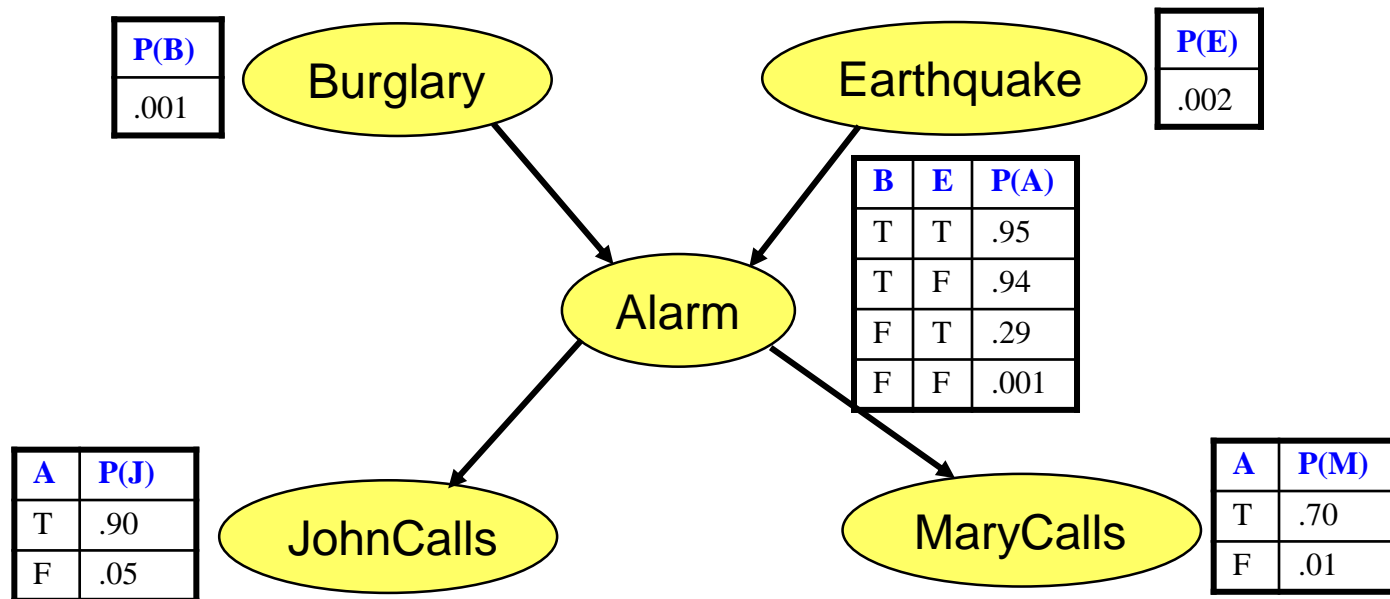


Board



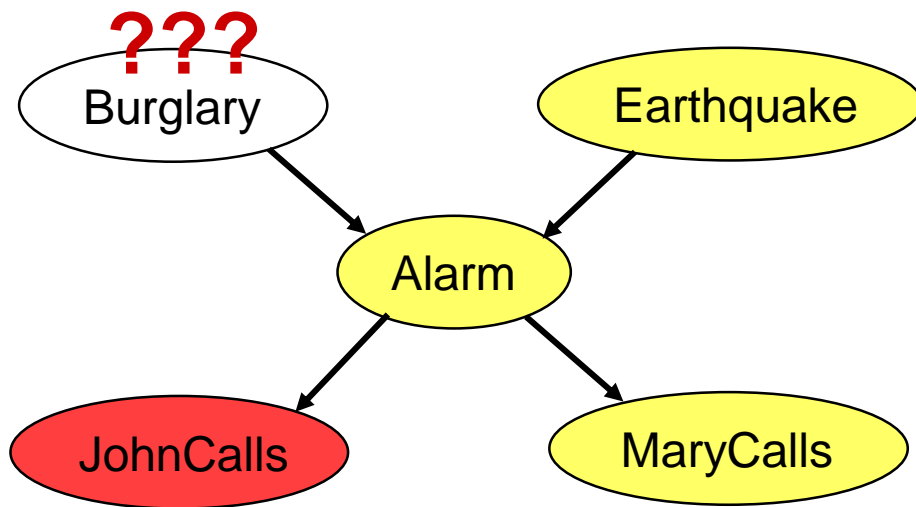
Example: Alarm Network with CPTs

- Each node has a **conditional probability table (CPT)** that gives the probability of each of its values given every possible combination of values for its parents (conditioning case).
 - Roots (sources) of the DAG that have no parents are given prior probabilities.



Bayes Net Inference

- Given known values for some **evidence variables**, determine the posterior probability of some **query variables**.
- Example: Given that John and Mary call, what is the probability that there is a Burglary?



Joint Distributions for Bayes Nets

- A Bayesian Network implicitly defines a joint distribution.

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

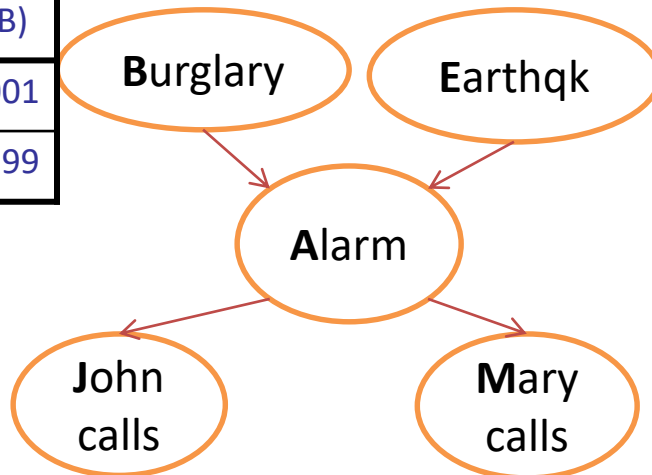
- Example

$$\begin{aligned} &P(J \wedge M \wedge A \wedge \neg B \wedge \neg E) \\ &= P(J \mid A)P(M \mid A)P(A \mid \neg B \wedge \neg E)P(\neg B)P(\neg E) \\ &= 0.9 \times 0.7 \times 0.001 \times 0.999 \times 0.998 = 0.00062 \end{aligned}$$

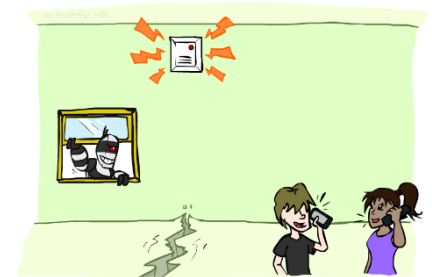
- An inefficient approach to inference is:
 - 1) Compute the joint distribution using this equation.
 - 2) Compute any desired conditional probability using the joint distribution.

Inference in Alarm Network

B	P(B)
+b	0.001
-b	0.999



E	P(E)
+e	0.002
-e	0.998



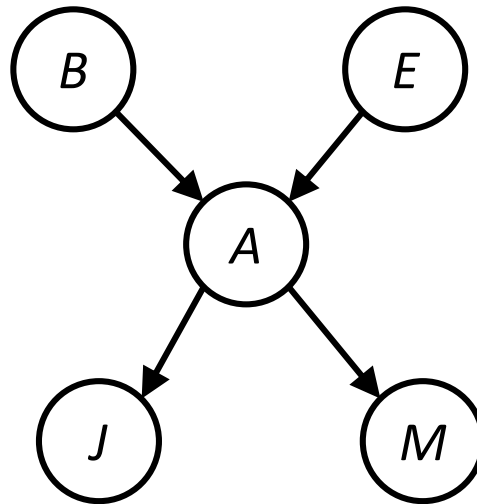
A	J	P(J A)
+a	+j	0.9
+a	-j	0.1
-a	+j	0.05
-a	-j	0.95

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

Inference Example

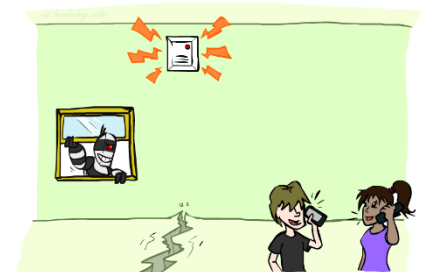
B	P(B)
+b	0.001
-b	0.999



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

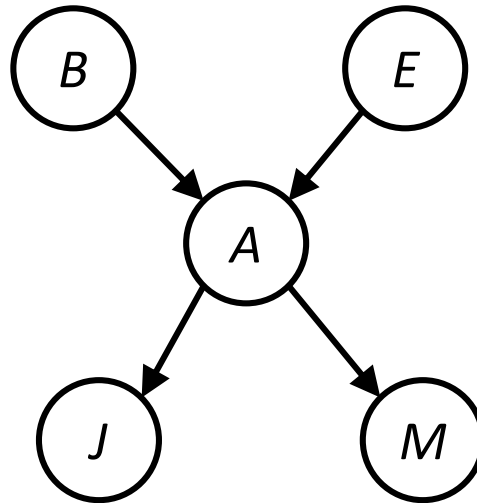


$$P(+b, -e, +a, -j, +m) =$$

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

Inference Example (2)

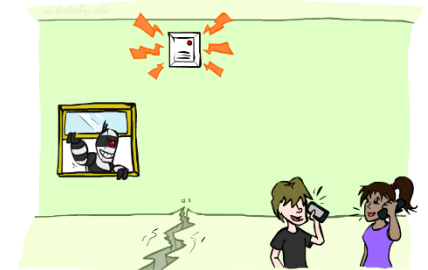
B	P(B)
+b	0.001
-b	0.999



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

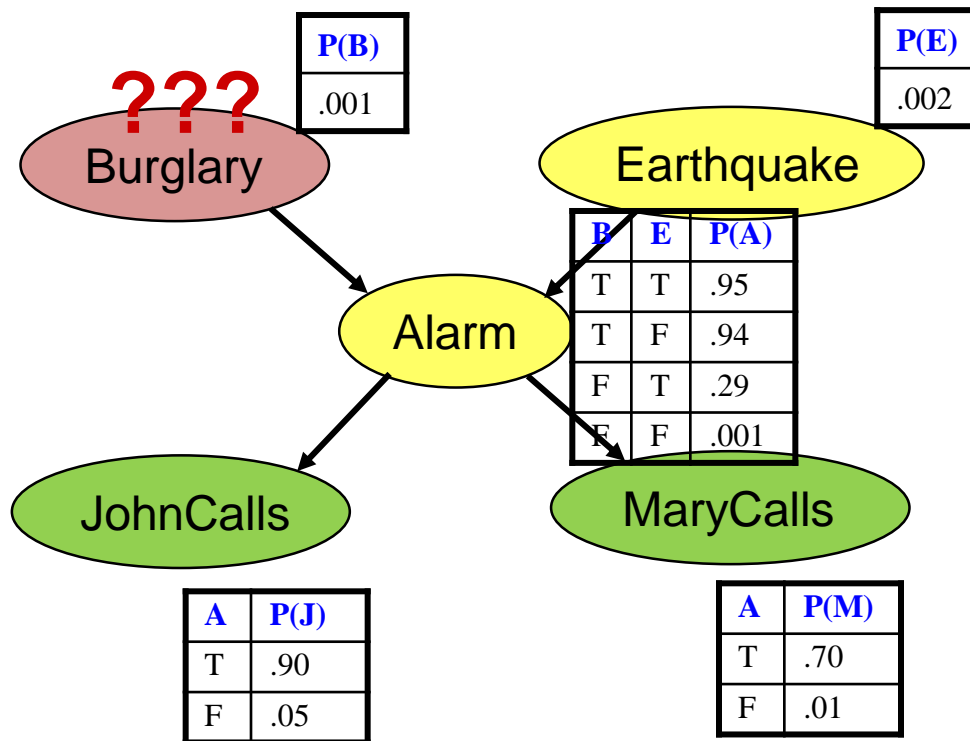


B	E	A	P(A B,E)
+b	+e	+a	0.95
+b	+e	-a	0.05
+b	-e	+a	0.94
+b	-e	-a	0.06
-b	+e	+a	0.29
-b	+e	-a	0.71
-b	-e	+a	0.001
-b	-e	-a	0.999

$$\begin{aligned}
 P(+b, -e, +a, -j, +m) &= \\
 P(+b)P(-e)P(+a|+b, -e)P(-j|+a)P(+m|+a) &= \\
 0.001 \times 0.998 \times 0.94 \times 0.1 \times 0.7 &
 \end{aligned}$$

Bayes Net Inference: Example

- Example: Given that John and Mary call (+j, +m), what is the probability that there is a Burglary (+b)?



$$P(+b \mid +j, +m) = ?$$

Inference by Enumeration

- General case:

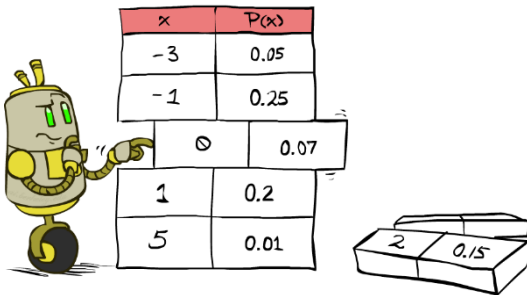
- Evidence variables: $E_1 \dots E_k = e_1 \dots e_k$
 - Query* variable: Q
 - Hidden variables: $H_1 \dots H_r$
- $$\left. \begin{array}{l} E_1 \dots E_k = e_1 \dots e_k \\ Q \\ H_1 \dots H_r \end{array} \right\} \begin{array}{l} X_1, X_2, \dots X_n \\ \text{All} \\ \text{variables} \end{array}$$

- We want:

** Works fine with multiple query variables, too*

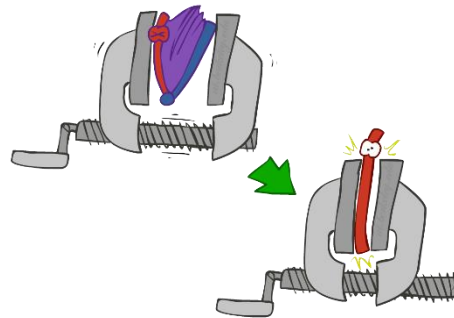
$$P(Q|e_1 \dots e_k)$$

- Step 1: Select the entries consistent with the evidence



x	P(x)
-3	0.05
-1	0.25
0	0.07
1	0.2
5	0.01

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



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

- Step 3: Normalize

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

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

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

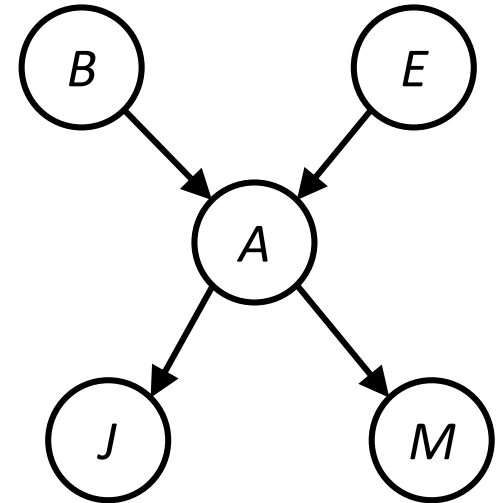
Inference by Enumeration in Bayes' Net

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

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

$$= \sum_{e,a} P(B, e, a, +j, +m)$$

$$= \sum_{e,a} P(B)P(e)P(a|B, e)P(+j|a)P(+m|a)$$



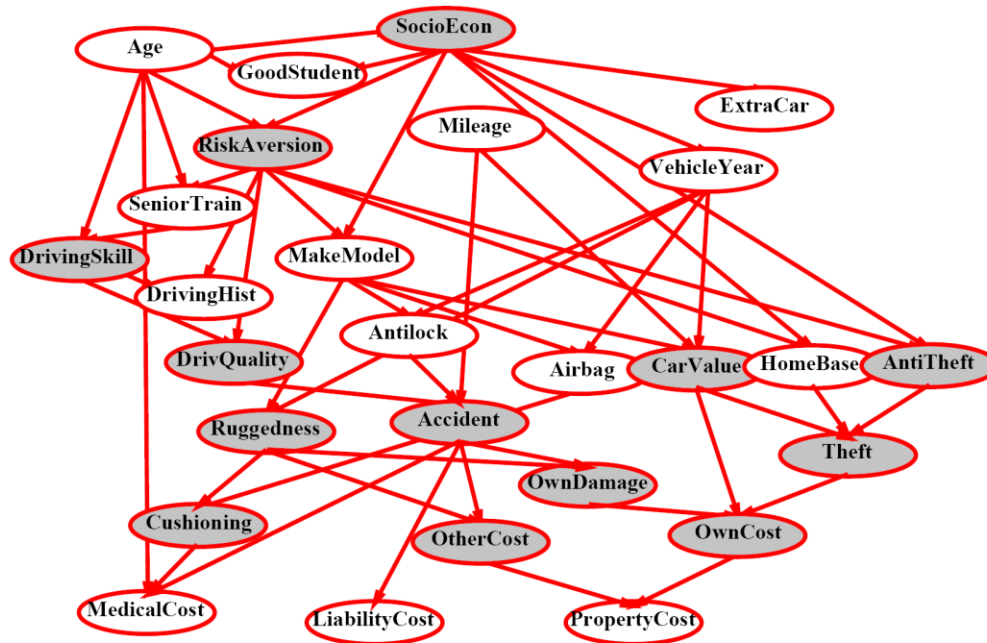
$$\begin{aligned} = & P(B)P(+e)P(+a|B, +e)P(+j|+a)P(+m|+a) + P(B)P(+e)P(-a|B, +e)P(+j|-a)P(+m|-a) \\ & P(B)P(-e)P(+a|B, -e)P(+j|+a)P(+m|+a) + P(B)P(-e)P(-a|B, -e)P(+j|-a)P(+m|-a) \end{aligned}$$

Exact Inference: By Enumeration

- Only need the CPTs to construct all possible ways query could be true (and sum them):

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

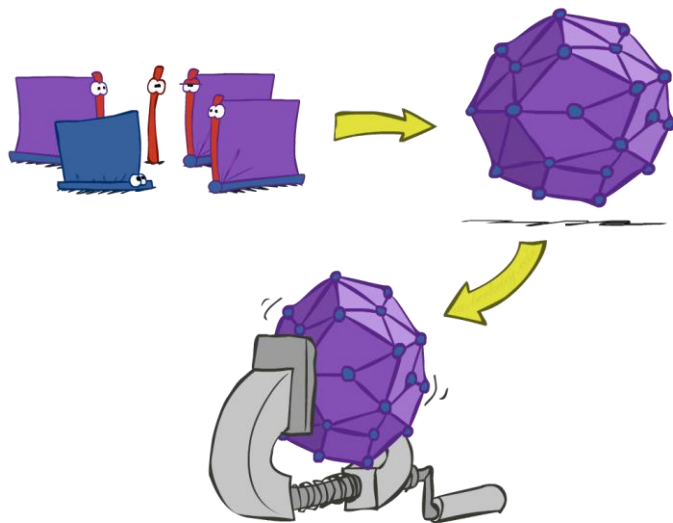
Inference by Enumeration?



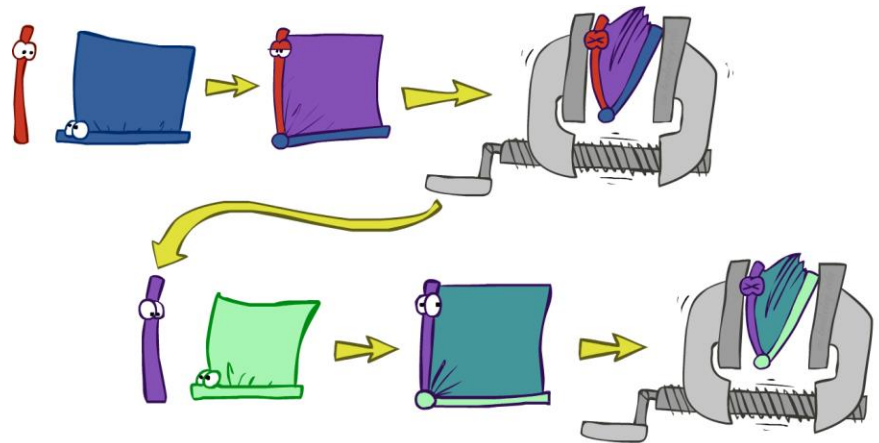
$$P(\text{Antilock} | \text{observed variables}) = ?$$

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

Inference by Enumeration: Procedural Outline

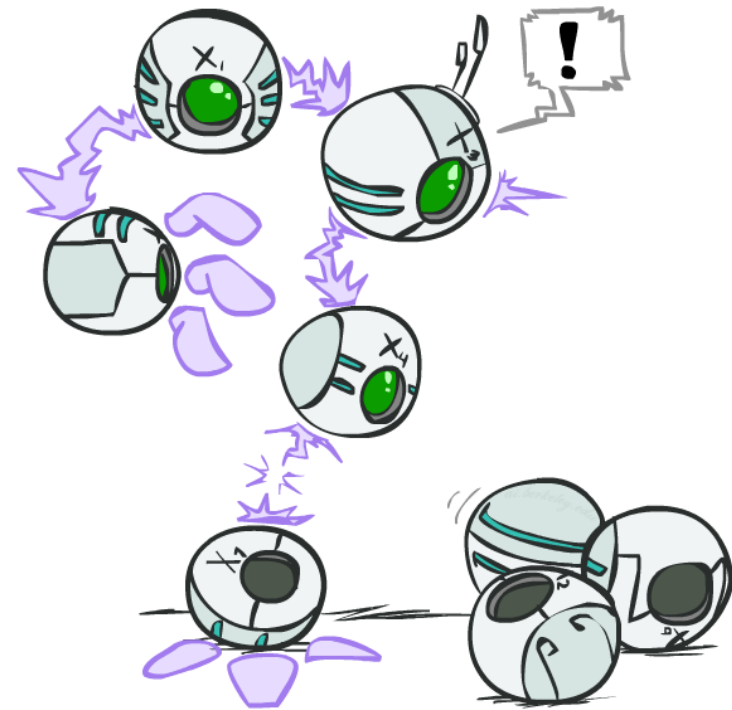
- Track objects called **factors**
- Initial factors are local CPTs (one per node)

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

- Any known values are selected
 - E.g. if we know $L = +l$, the initial factors are

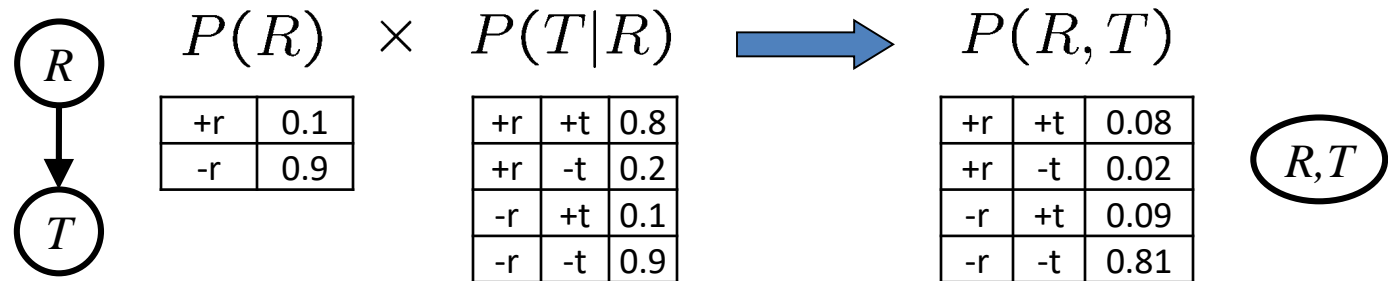
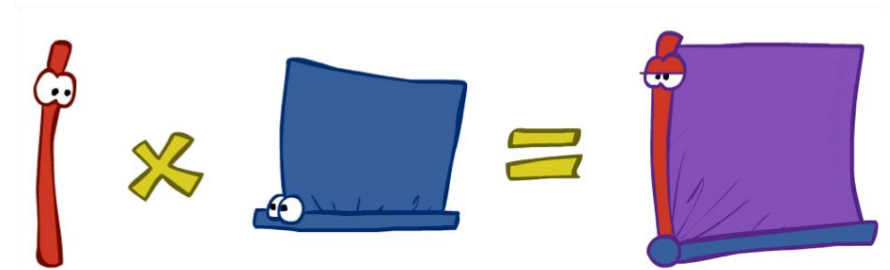
$P(R)$		$P(T R)$			$P(+l T)$		
+r	0.1	+r	+t	0.8	+t	+l	0.3
-r	0.9	+r	-t	0.2	-t	+l	0.1
		-r	+t	0.1			
		-r	-t	0.9			

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



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** $\forall r, t : P(r, t) = P(r) \cdot P(t|r)$

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:

$$P(R, T)$$

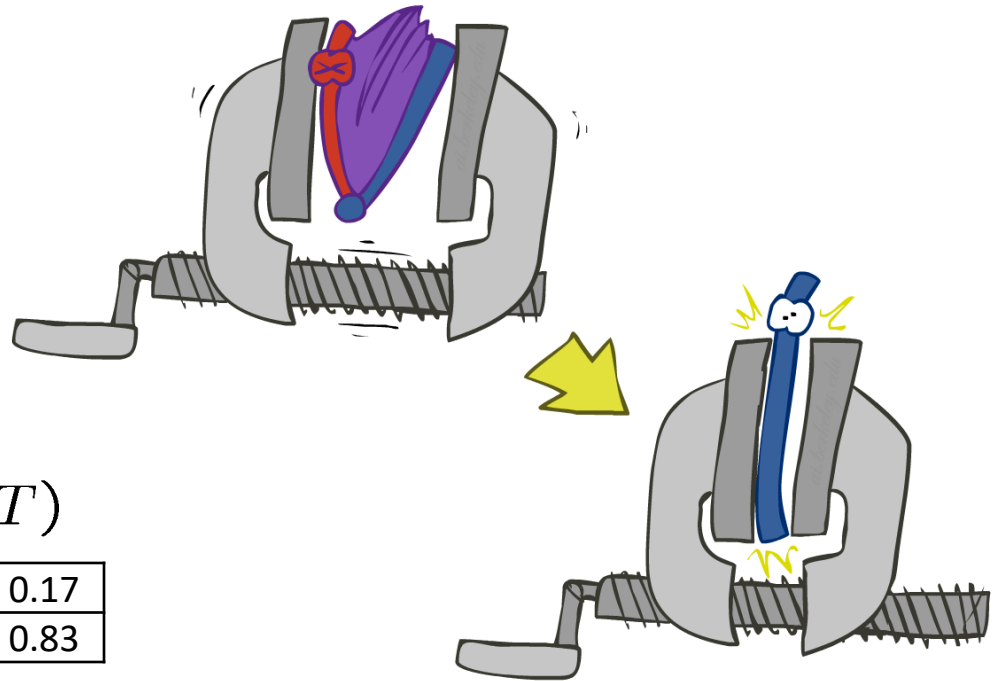
+r	+t	0.08
+r	-t	0.02
-r	+t	0.09
-r	-t	0.81

sum R

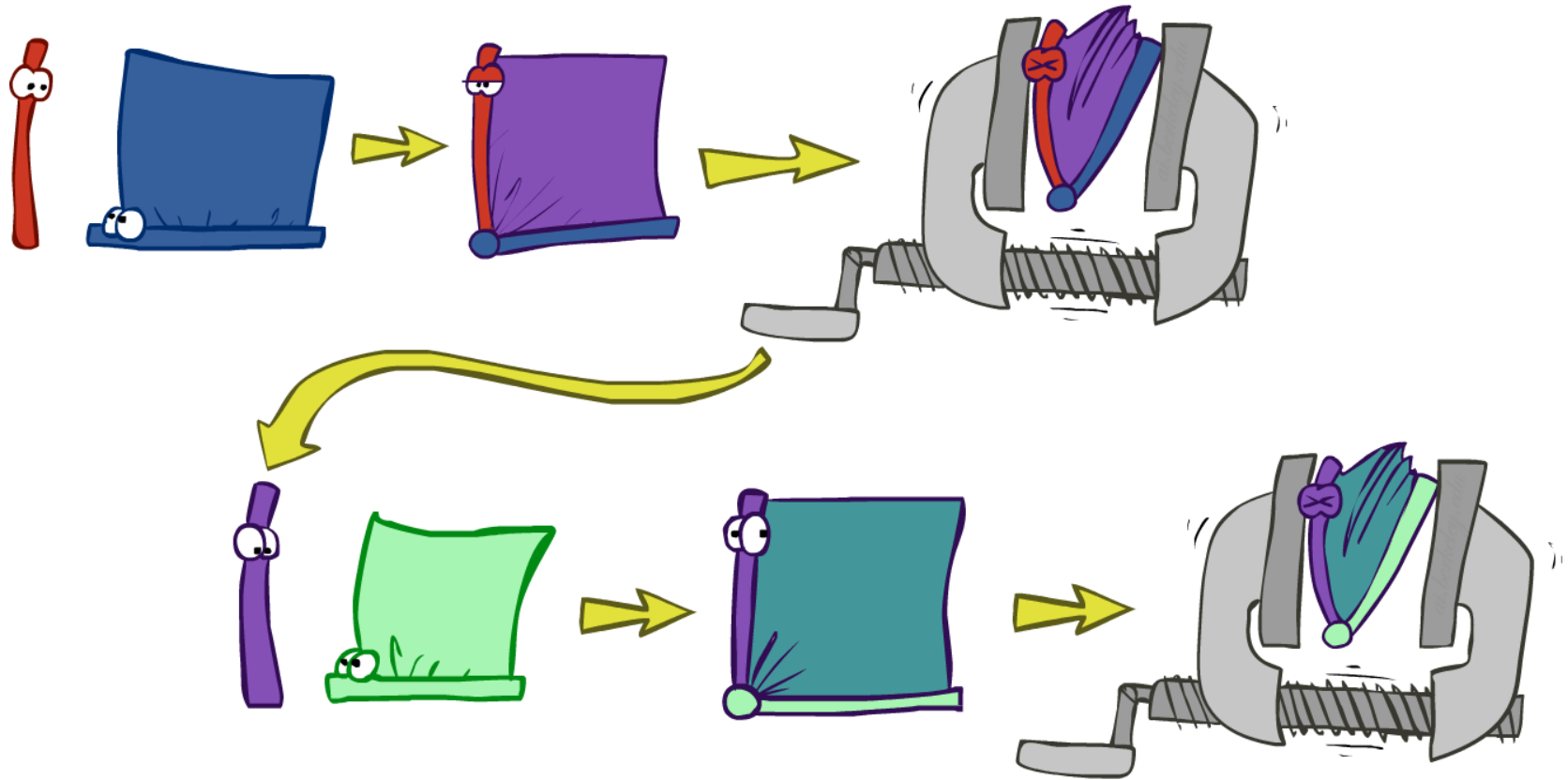


$$P(T)$$

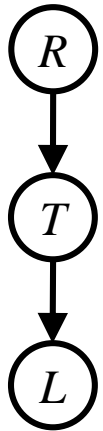
+t	0.17
-t	0.83



Marginalizing Early (= Variable Elimination)



Traffic Domain



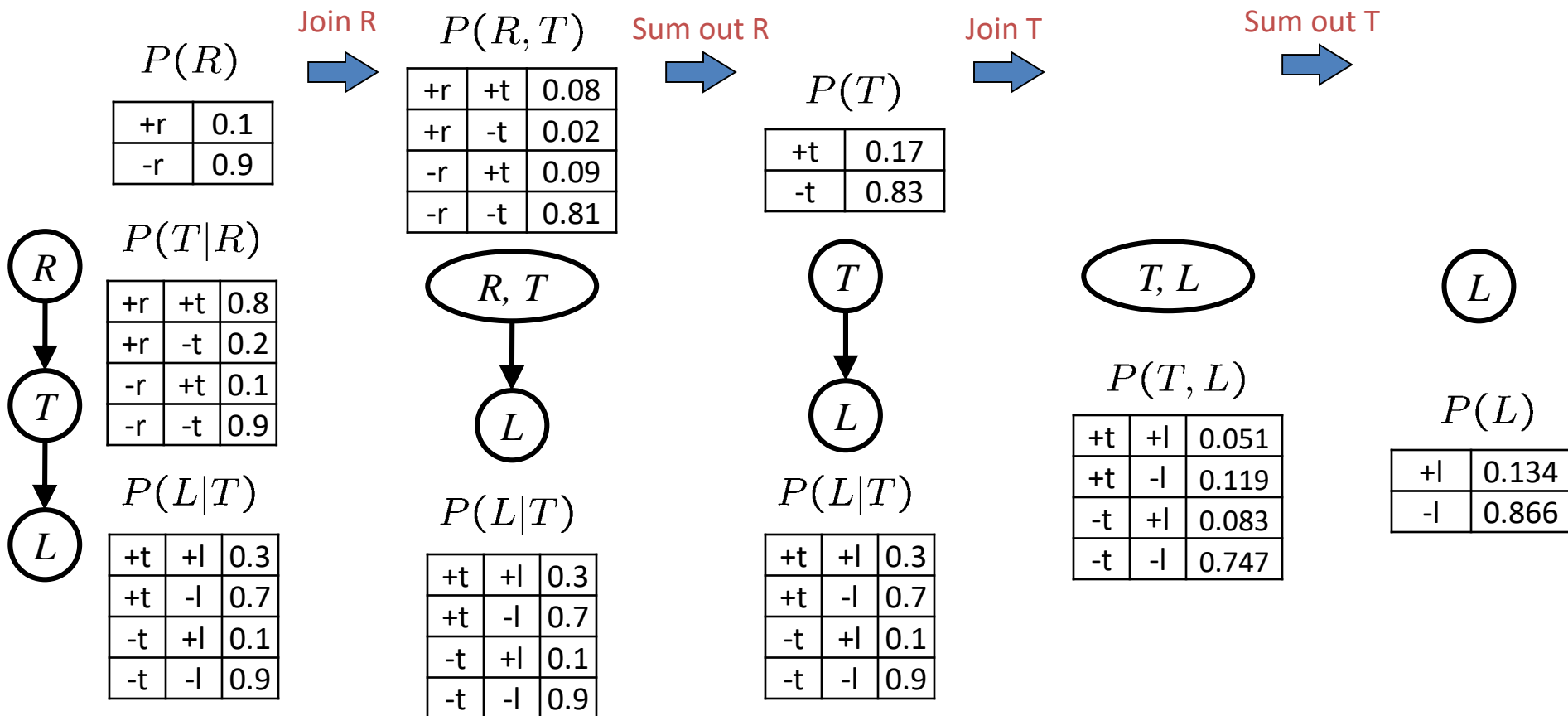
- $P(L) = ?$
• Inference by Enumeration

$$= \sum_t \sum_r \underbrace{P(L|t)P(r)P(t|r)}_{\text{Join on } r} \underbrace{}_{\text{Join on } t} \underbrace{}_{\text{Eliminate } r} \underbrace{}_{\text{Eliminate } t}$$

- Variable Elimination

$$= \sum_t P(L|t) \underbrace{\sum_r P(r)P(t|r)}_{\text{Join on } r} \underbrace{}_{\text{Eliminate } r} \underbrace{}_{\text{Join on } t} \underbrace{}_{\text{Eliminate } t}$$

Marginalizing Early! (aka VE)



Evidence

- If evidence, start with factors that select that evidence
 - No evidence uses these initial factors:

$$P(R)$$

+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

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

$$P(+r)$$

+r	0.1
----	-----

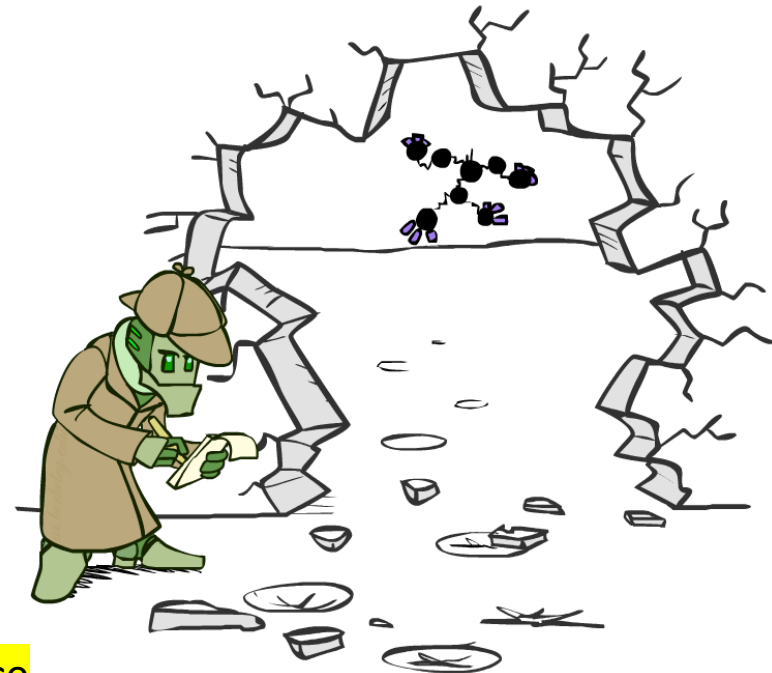
$$P(T|+r)$$

+r	+t	0.8
+r	-t	0.2

$$P(L|T)$$

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

- Eliminate all variables other than query + evidence



Evidence II

- Result will be a selected joint of query and evidence
 - E.g. for $P(L \mid +r)$, we would end up with:

$$P(+r, L)$$

+r	+l	0.026
+r	-l	0.074

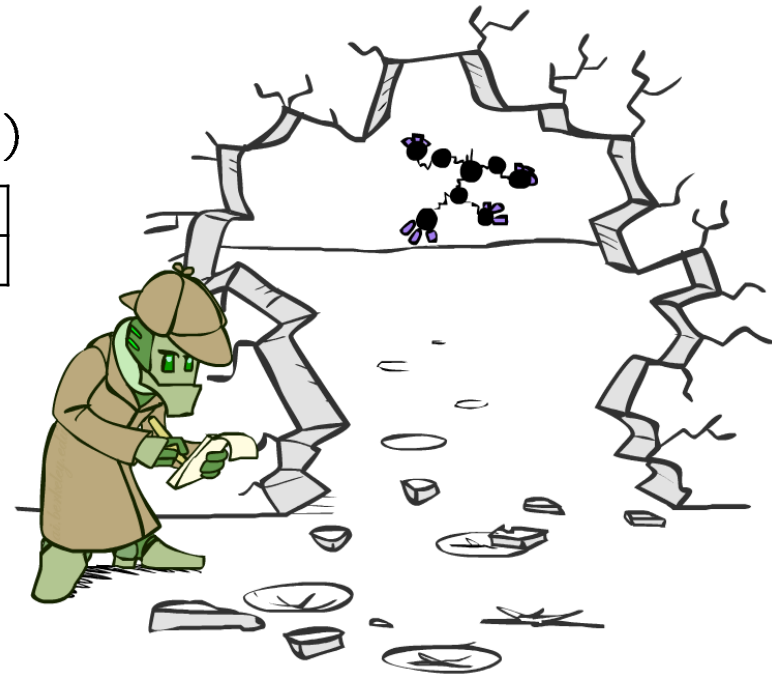
Normalize



$$P(L \mid +r)$$

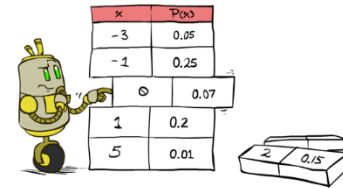
+l	0.26
-l	0.74

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



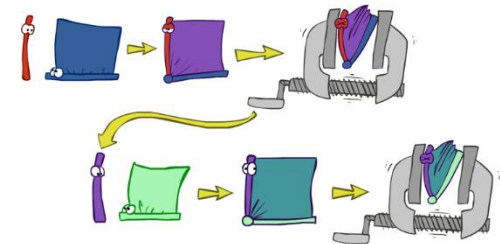
General Variable Elimination

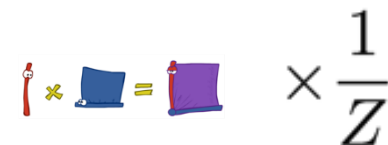
- Query: $P(Q|E_1 = e_1, \dots, E_k = e_k)$
- Start with initial factors:
 - Local CPTs (but instantiated by evidence)
- While there are still hidden variables (not Q or evidence):
 - Pick a hidden variable H
 - Join all factors mentioning H
 - Eliminate (sum out) H
- Join all remaining factors and normalize



x	P(x)
-3	0.05
-1	0.25
0	0.07
1	0.2
5	0.01

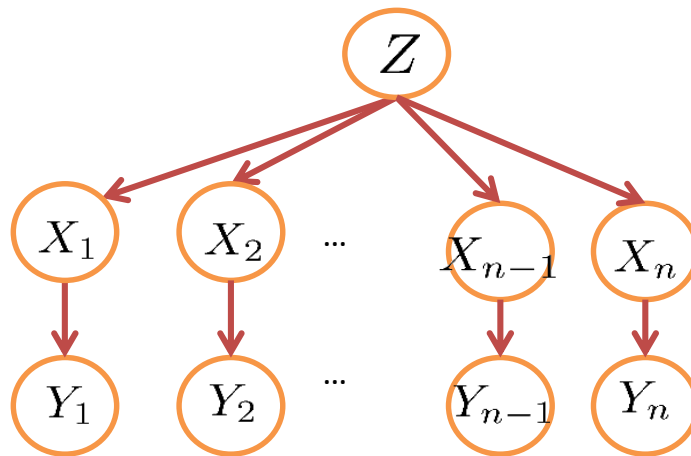
2	0.15
---	------




$$\text{red factor} \times \text{blue factor} = \text{purple factor} \times \frac{1}{Z}$$

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?



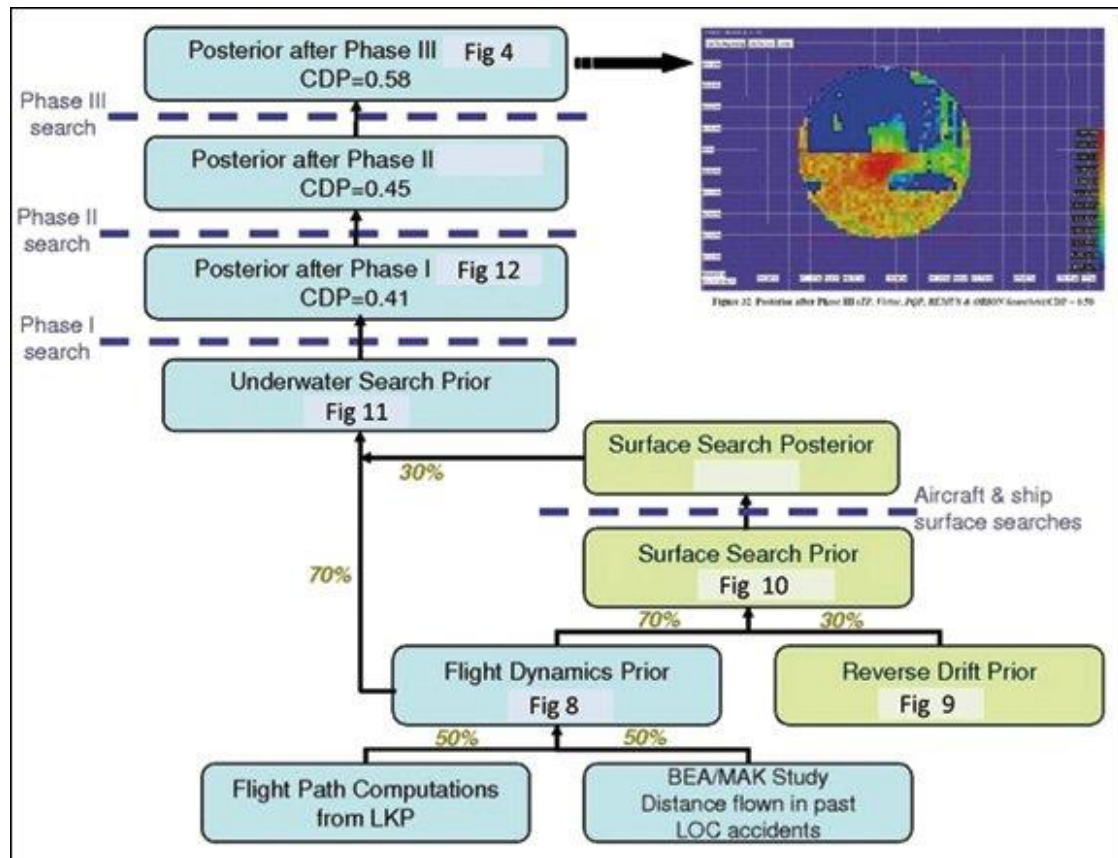
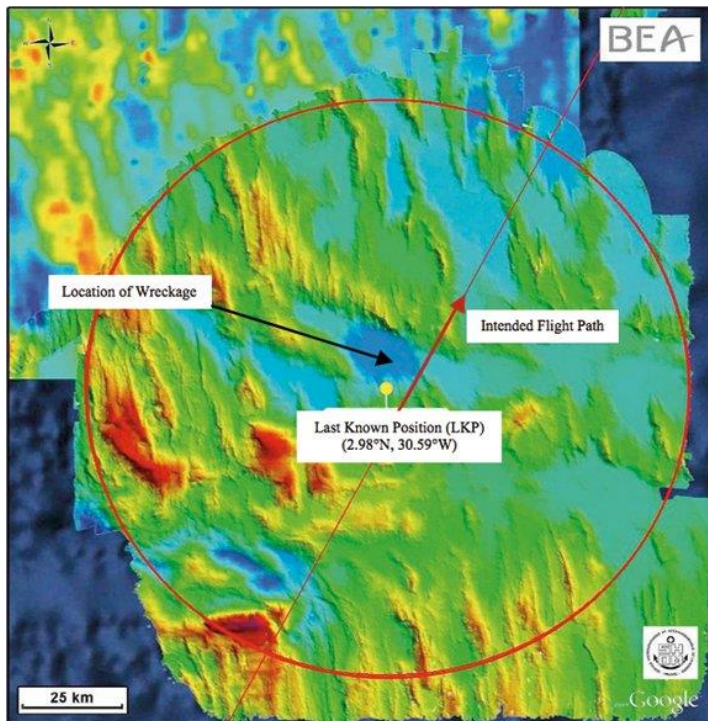
- Answer: 2^{n+1} versus 2^2 (assuming binary)
- In general: the ordering can greatly affect efficiency.

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.
 - Try to marginalize (eliminate) small factors first
- Does there always exist an ordering that only results in small factors?
 - No!

Finding Air France Flight 447

- <http://fivethirtyeight.com/features/how-statisticians-could-help-find-flight-370/>
- <https://www.informs.org/ORMS-Today/Public-Articles/August-Volume-38-Number-4/In-Search-of-Air-France-Flight-447>



Complexity of Bayes Net Inference

- In general, the problem of Bayes Net inference is NP-hard (exponential in the size of the graph).
- For **singly-connected networks** or **polytrees** in which there are no undirected loops, there are linear-time algorithms based on **belief propagation**.
 - Each node sends local evidence messages to their children and parents.
 - Each node updates belief in each of its possible values based on incoming messages from its neighbors and propagates evidence on to its neighbors.
- There are approximations to inference for general networks based on **loopy belief propagation** that iteratively refines probabilities that converge to accurate values in the limit.