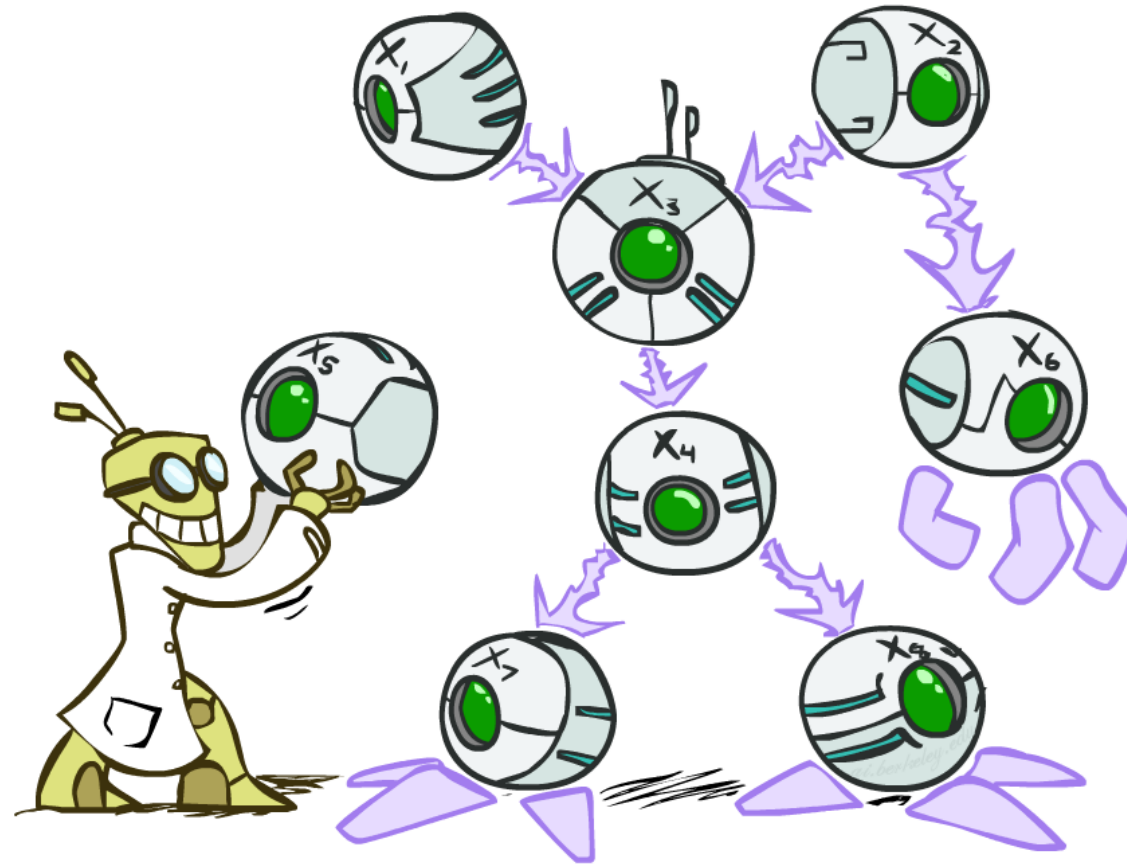


# CS 188: Artificial Intelligence

## Bayesian Networks



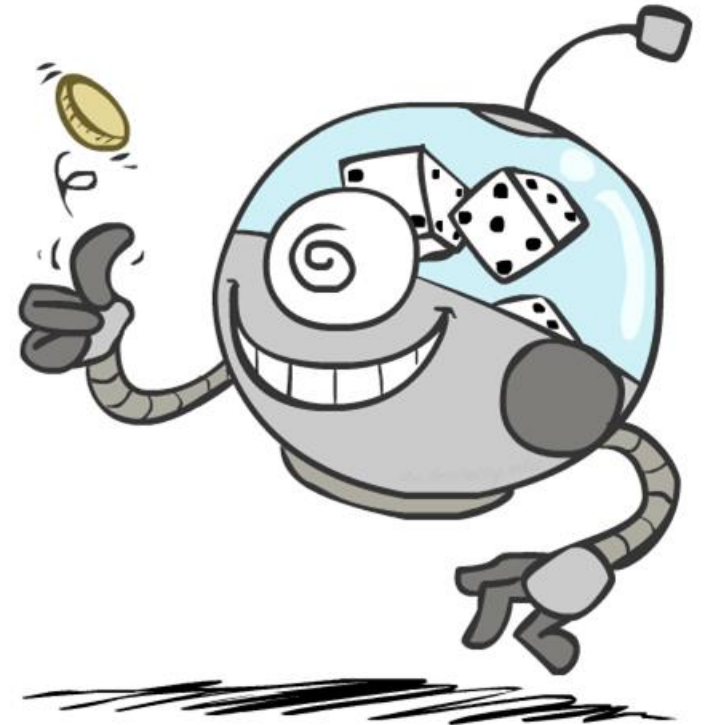
Instructor: Evgeny Pobachienko — UC Berkeley

[Slides credit: Dan Klein, Pieter Abbeel, Anca Dragan, Stuart Russell, Satish Rao, and many others]

# Recall: Random Variables

---

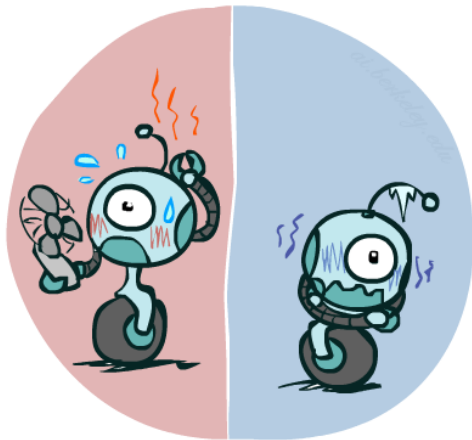
- Recall: random variable is some aspect of the world about which we (may) have uncertainty
  - R = Is it raining?
  - T = Is it hot?
  - D = How long will it take to drive to work?
- Capital letters: Random variables
- Lowercase letters: values that the R.V. can take
  - $r \in \{+r, -r\}$
  - $t \in \{+t, -t\}$
  - $d \in [0, \infty)$



# Probability Distributions

- Associate a probability with each value

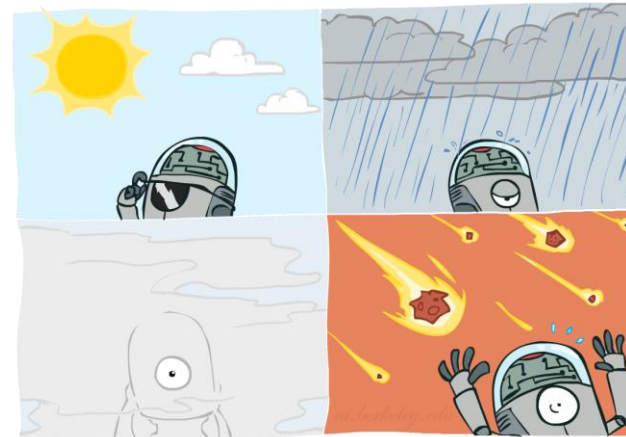
- Temperature:



$P(T)$

T	P
hot	0.5
cold	0.5

- Weather:



$P(W)$

W	P
sun	0.6
rain	0.1
fog	0.3
meteor	0.0

# Joint Distributions

- A *joint distribution* over a set of random variables:  $X_1, X_2, \dots, X_n$  specifies a real number for each assignment (or *outcome*):

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

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

- Must obey:  $P(x_1, x_2, \dots, x_n) \geq 0$  (non-negativity)

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

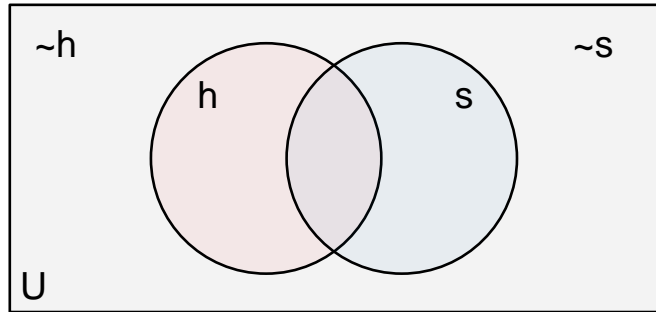
- Size of distribution if  $n$  variables with domain sizes  $d$ ?
  - For all but the smallest distributions, impractical to write out!

$P(T, W)$

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

# Probability

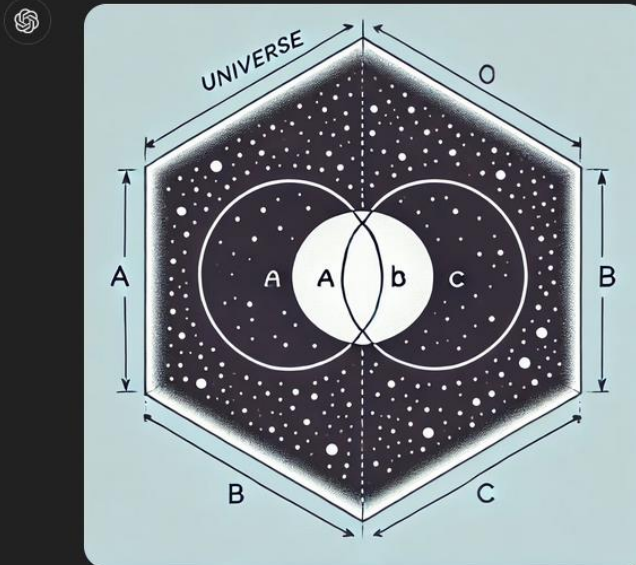
---



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

# AI to teach AI

I'm teaching probability. Make a diagram with a rectangle representing the universe, and intersecting circles A, B, C as events



Here is the diagram representing the universe with intersecting circles for events A, B, and C.

If you need any adjustments or further details, let me know!

Venn diagram style, with universe meaning all possible probability events.



**You've reached your image creation limit.**

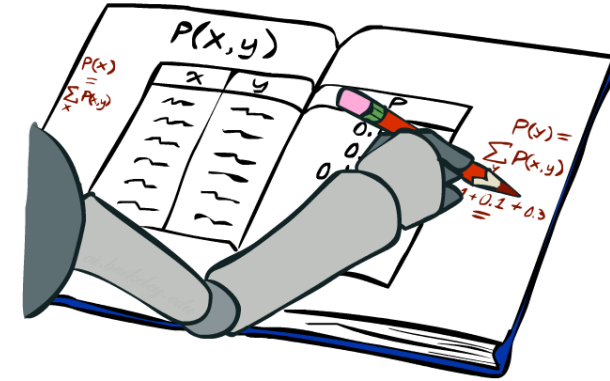
Upgrade to ChatGPT Plus or try again tomorrow after 7:11 PM.

Get Plus



# Marginal Distributions

- Marginal distributions are sub-tables which eliminate variables
- Marginalization (summing out): Combine collapsed rows by adding



$P(T, W)$

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



$$P(t) = \sum_s P(t, s)$$



$$P(s) = \sum_t P(t, s)$$

$P(T)$

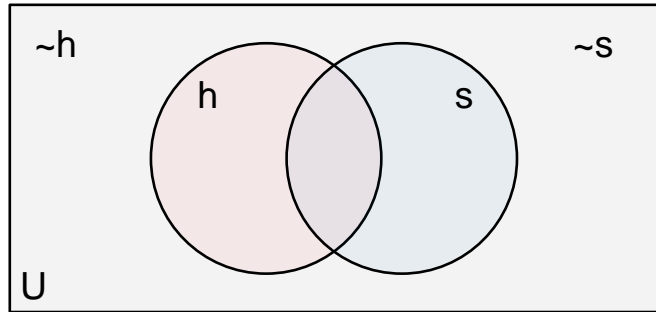
T	P
hot	0.5
cold	0.5

$P(W)$

W	P
sun	0.6
rain	0.4

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

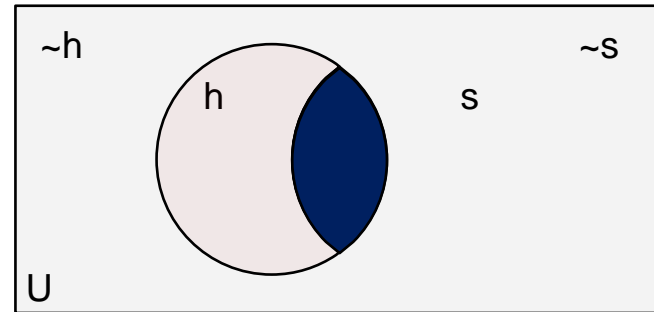
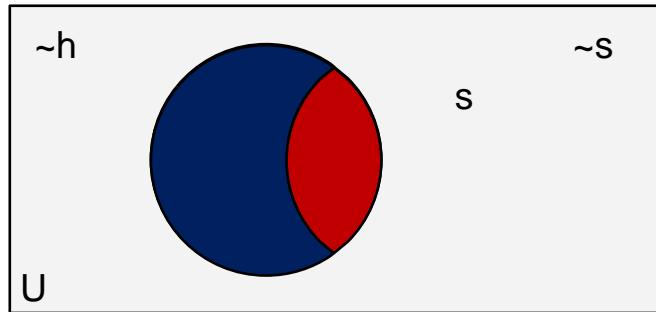
# Probability



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

$$P(h) = P(h, s) + P(h, \sim s)$$

$$P(s|h) = \frac{P(s, h)}{P(h)}$$





# Conditional Distributions

- Conditional distributions are probability distributions over some variables given fixed values of others

Conditional Distributions

$P(W|T)$

$P(W T = hot)$	
W	P
sun	0.8
rain	0.2

$P(W T = cold)$	
W	P
sun	0.4
rain	0.6

Joint Distribution

$P(T, W)$

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

# Normalization Trick

$P(T, W)$

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

$$\begin{aligned}P(W = s|T = c) &= \frac{P(W = s, T = c)}{P(T = c)} \\ &= \frac{P(W = s, T = c)}{P(W = s, T = c) + P(W = r, T = c)} \\ &= \frac{0.2}{0.2 + 0.3} = 0.4\end{aligned}$$



$$\begin{aligned}P(W = r|T = c) &= \frac{P(W = r, T = c)}{P(T = c)} \\ &= \frac{P(W = r, T = c)}{P(W = s, T = c) + P(W = r, T = c)} \\ &= \frac{0.3}{0.2 + 0.3} = 0.6\end{aligned}$$

$P(W|T = c)$

W	P
sun	0.4
rain	0.6

# Normalization Trick

$P(T, W)$

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

**SELECT** the joint probabilities matching the evidence



$P(c, W)$

T	W	P
cold	sun	0.2
cold	rain	0.3

**NORMALIZE** the selection (make it sum to one)



$P(W|T = c)$

W	P
sun	0.4
rain	0.6

# To Normalize

○ (Dictionary) To bring or restore to a **normal** condition

○ Procedure:

- Step 1: Compute  $Z = \text{sum over all entries}$
- Step 2: Divide every entry by  $Z$

All entries sum to ONE

○ Example

W	P
sun	0.2
rain	0.3

Normalize  
Z = 0.5

W	P
sun	0.4
rain	0.6

# Probabilistic Inference

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



# Inference by Enumeration

---

- $P(W)$ ?

S	T	W	P
summer	hot	sun	0.30
summer	hot	rain	0.05
summer	cold	sun	0.10
summer	cold	rain	0.05
winter	hot	sun	0.10
winter	hot	rain	0.05
winter	cold	sun	0.15
winter	cold	rain	0.20

# Inference by Enumeration

---

- $P(W)$ ?

S	T	W	P
summer	hot	sun	0.30
summer	hot	rain	0.05
summer	cold	sun	0.10
summer	cold	rain	0.05
winter	hot	sun	0.10
winter	hot	rain	0.05
winter	cold	sun	0.15
winter	cold	rain	0.20

# Inference by Enumeration

---

- $P(W)$ ?

S	T	W	P
summer	hot	sun	0.30
summer	hot	rain	0.05
summer	cold	sun	0.10
summer	cold	rain	0.05
winter	hot	sun	0.10
winter	hot	rain	0.05
winter	cold	sun	0.15
winter	cold	rain	0.20



# Inference by Enumeration

---

○ P(W)?

$$P(\text{sun}) = .3 + .1 + .1 + .15 = .65$$

S	T	W	P
summer	hot	sun	0.30
summer	hot	rain	0.05
summer	cold	sun	0.10
summer	cold	rain	0.05
winter	hot	sun	0.10
winter	hot	rain	0.05
winter	cold	sun	0.15
winter	cold	rain	0.20

# Inference by Enumeration

---

○  $P(W)$ ?

$$P(\text{sun}) = .3 + .1 + .1 + .15 = .65$$

$$P(\text{rain}) = 1 - .65 = .35$$

S	T	W	P
summer	hot	sun	0.30
summer	hot	rain	0.05
summer	cold	sun	0.10
summer	cold	rain	0.05
winter	hot	sun	0.10
winter	hot	rain	0.05
winter	cold	sun	0.15
winter	cold	rain	0.20

# Inference by Enumeration

---

- $P(W \mid \text{winter, hot})?$

S	T	W	P
summer	hot	sun	0.30
summer	hot	rain	0.05
summer	cold	sun	0.10
summer	cold	rain	0.05
winter	hot	sun	0.10
winter	hot	rain	0.05
winter	cold	sun	0.15
winter	cold	rain	0.20

# Inference by Enumeration

---

- $P(W \mid \text{winter, hot})?$

$P(\text{sun} \mid \text{winter, hot}) \sim .1$

$P(\text{rain} \mid \text{winter, hot}) \sim .05$

S	T	W	P
summer	hot	sun	0.30
summer	hot	rain	0.05
summer	cold	sun	0.10
summer	cold	rain	0.05
winter	hot	sun	0.10
winter	hot	rain	0.05
winter	cold	sun	0.15
winter	cold	rain	0.20

# Inference by Enumeration

---

- $P(W \mid \text{winter, hot})?$

$P(\text{sun} \mid \text{winter, hot}) \sim .1$   
 $P(\text{rain} \mid \text{winter, hot}) \sim .05$   
 $P(\text{sun} \mid \text{winter, hot}) = 2/3$   
 $P(\text{rain} \mid \text{winter, hot}) = 1/3$

S	T	W	P
summer	hot	sun	0.30
summer	hot	rain	0.05
summer	cold	sun	0.10
summer	cold	rain	0.05
winter	hot	sun	0.10
winter	hot	rain	0.05
winter	cold	sun	0.15
winter	cold	rain	0.20

# Inference by Enumeration

- General case:

- Evidence variables:  $E_1 \dots E_k = e_1 \dots e_k$
  - Query\* variable:  $Q$
  - Hidden variables:  $H_1 \dots H_r$
- }  $X_1, X_2, \dots, X_n$   
All variables

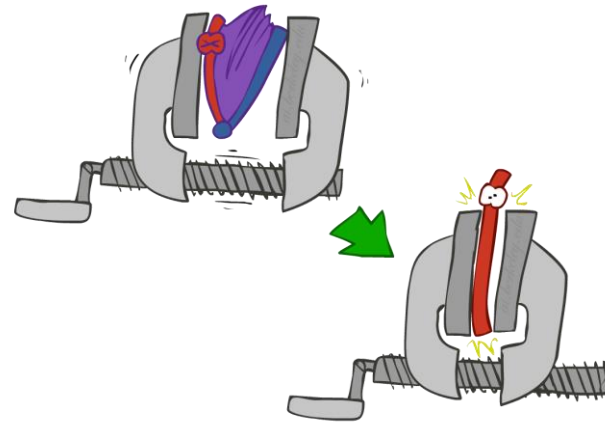
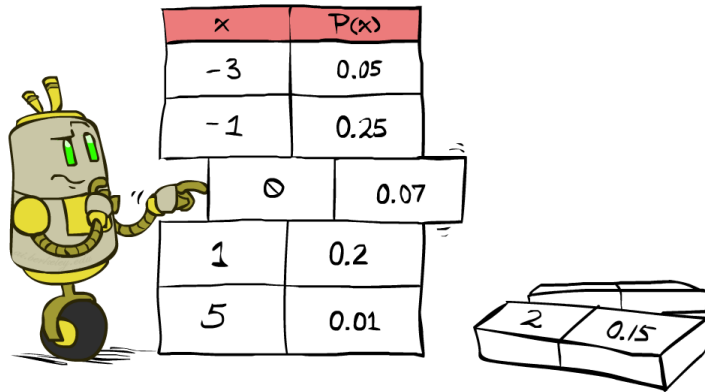
- We want:

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

- Step 1: Select the entries consistent with the evidence

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

- Step 3: Normalize



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

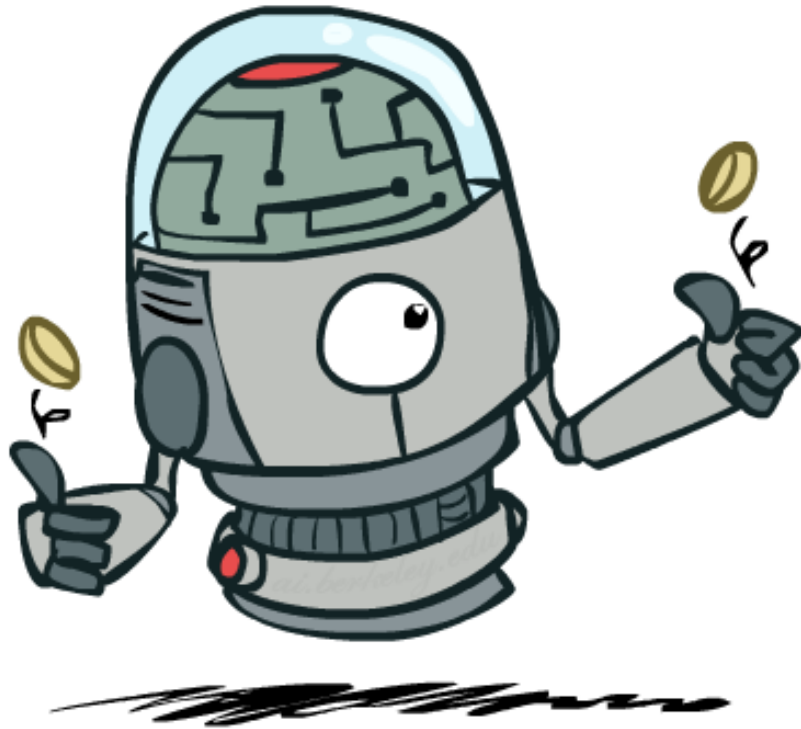
$$P(Q, e_1 \dots e_k) = \sum_{h_1 \dots h_r} \underbrace{P(Q, h_1 \dots h_r, e_1 \dots e_k)}_{X_1, X_2, \dots, X_n}$$

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

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

# Independence

---



# Independence

- Two variables are *independent* if:

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

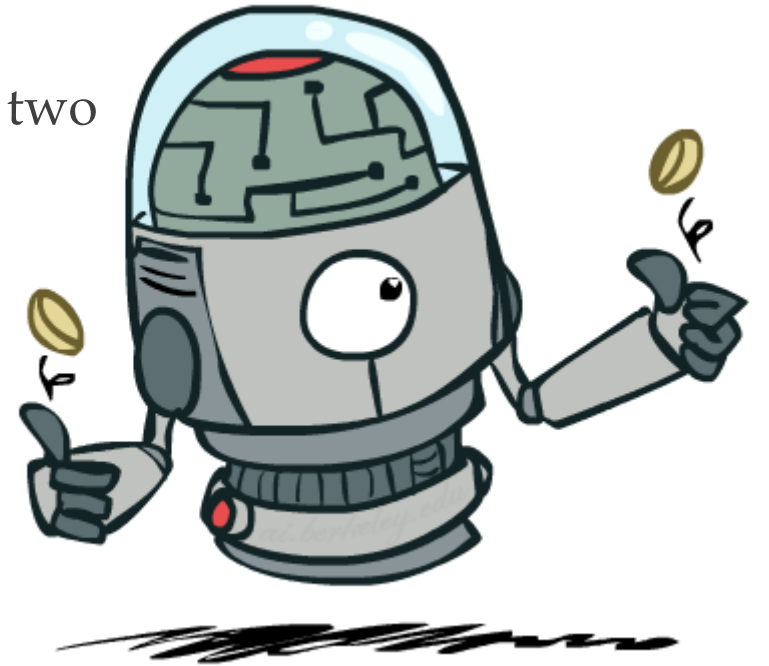
- This says that their joint distribution *factors* into a product of two simpler distributions
- Another form:

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

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

- Independence is a simplifying *modeling assumption*

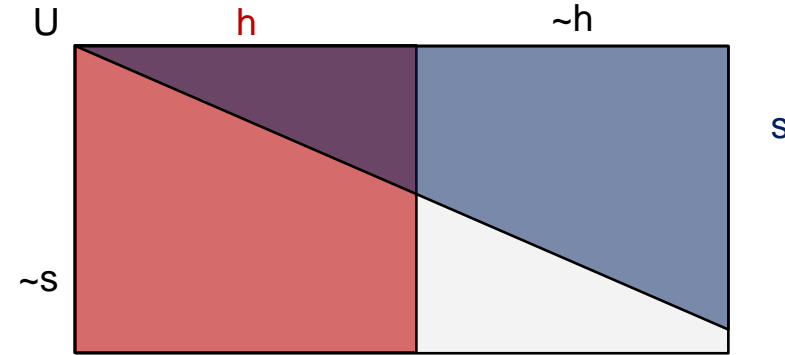
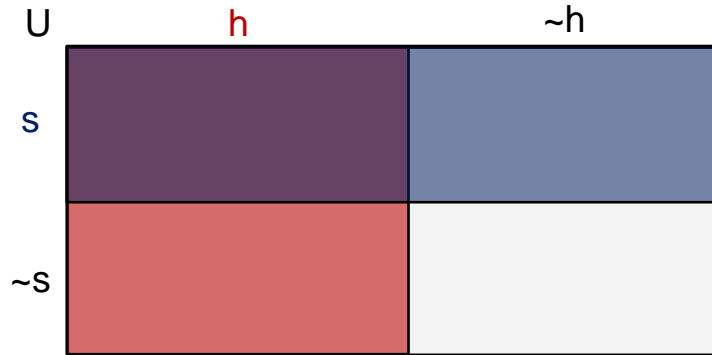
- *Empirical* joint distributions: at best “close” to independent
- What could we assume for {Weather, Traffic, Cavity, Toothache}?





# Independence

---



$$P(s|h) = \frac{P(s, h)}{P(h)}$$

$$P(s, h) = P(s|h) * P(h)$$

# Example: Independence?

---

$P_1(T, W)$

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

$P(T)$

T	P
hot	0.5
cold	0.5

$P_2(T, W)$

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

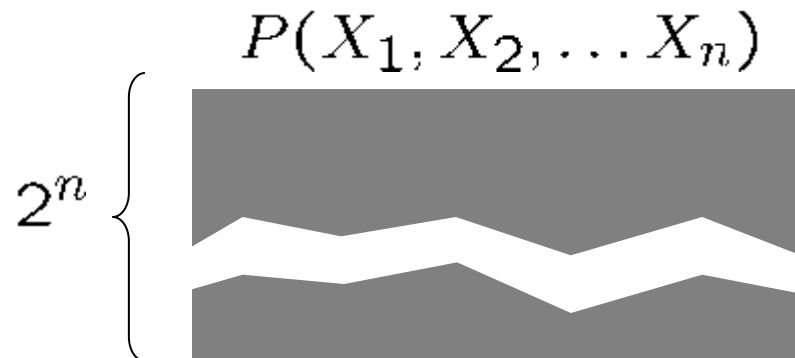
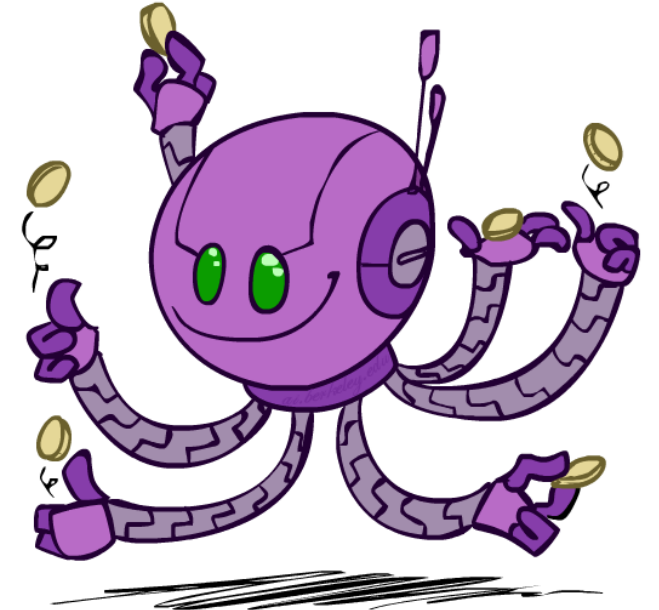
$P(W)$

W	P
sun	0.6
rain	0.4

# Example: Independence

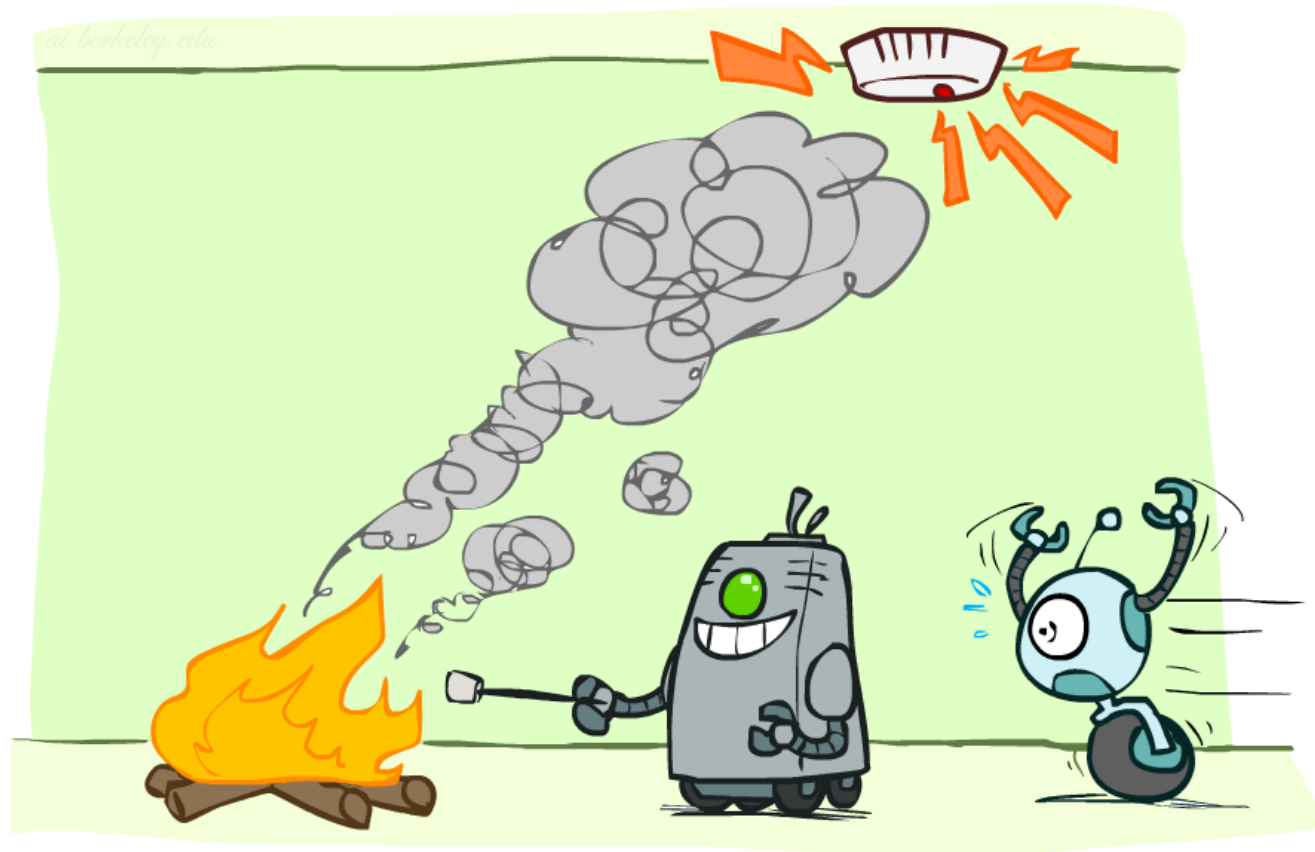
- N fair, independent coin flips:

$P(X_1)$		$P(X_2)$		...	$P(X_n)$	
H	0.5	H	0.5		H	0.5
T	0.5	T	0.5		T	0.5



# Conditional Independence

---



# 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

---

- 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

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

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

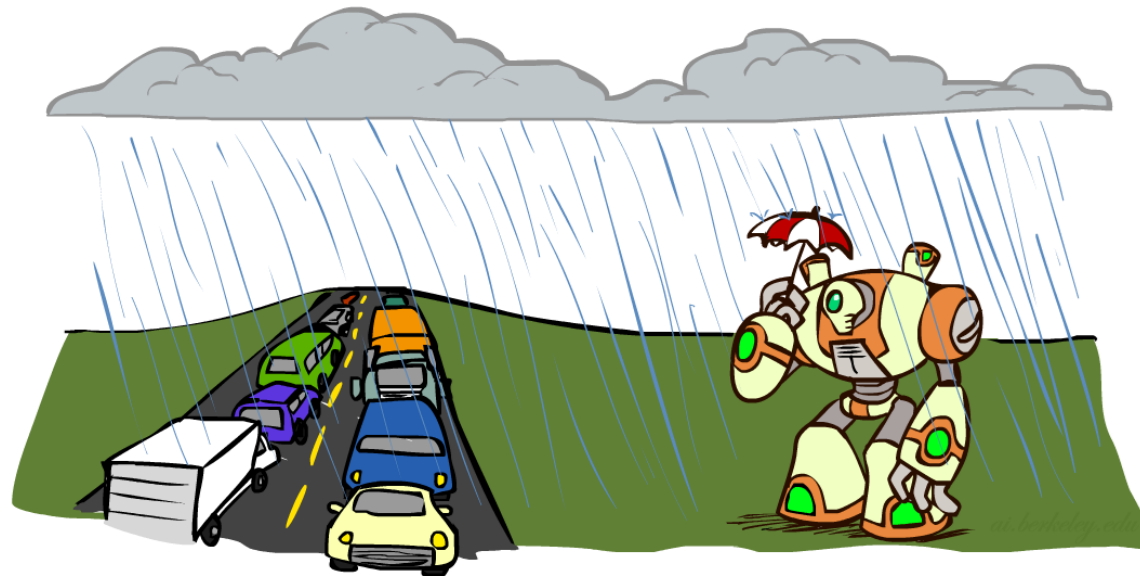
$$= \frac{P(x, y|z)P(z)}{P(y|z)P(z)}$$

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

# Conditional Independence

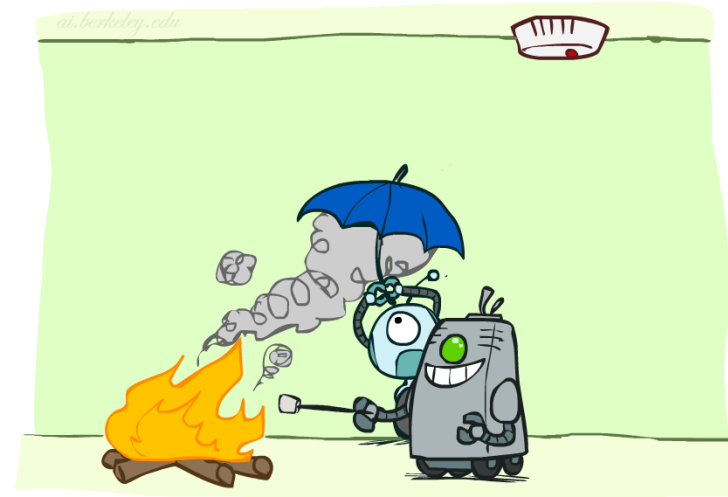
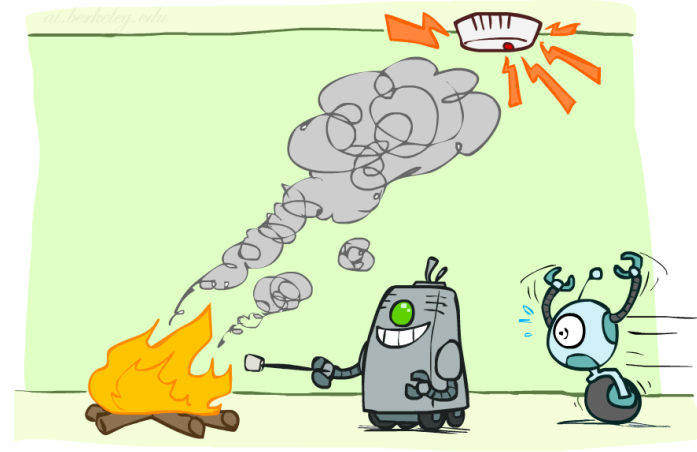
---

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



# Conditional Independence

- What about this domain:
  - Fire
  - Smoke
  - Alarm





# Conditional Independence and the Chain Rule

---

○ Chain rule:  $P(X_1, X_2, \dots, X_n) = P(X_1)P(X_2|X_1)P(X_3|X_1, X_2) \dots$

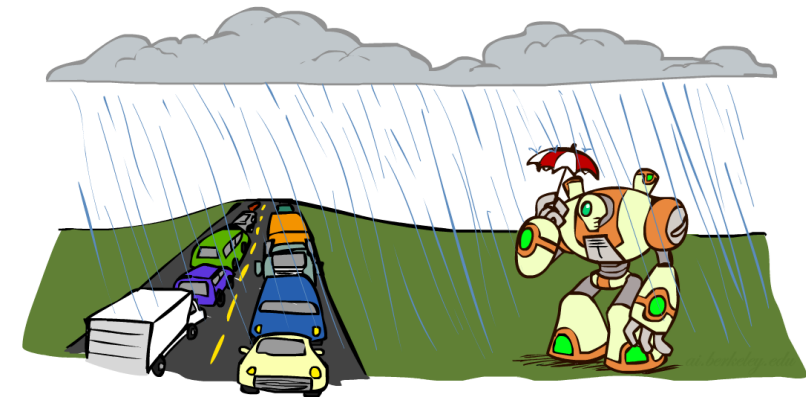
○ Trivial decomposition:

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

○ With assumption of conditional independence:

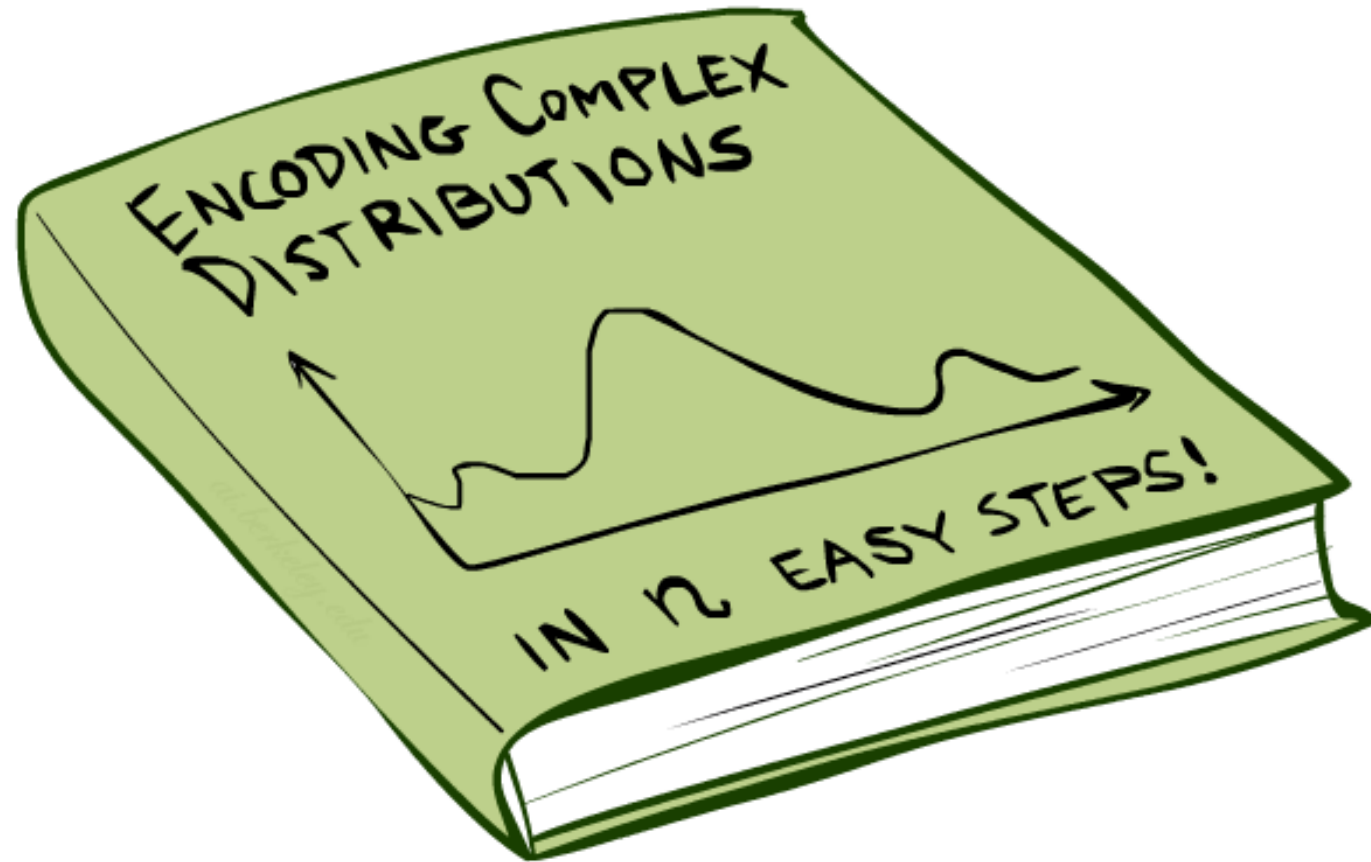
$$P(\text{Traffic, Rain, Umbrella}) = P(\text{Rain})P(\text{Traffic}|\text{Rain})P(\text{Umbrella}|\text{Rain})$$

○ Bayesian Networks/graphical models help us express conditional independence assumptions



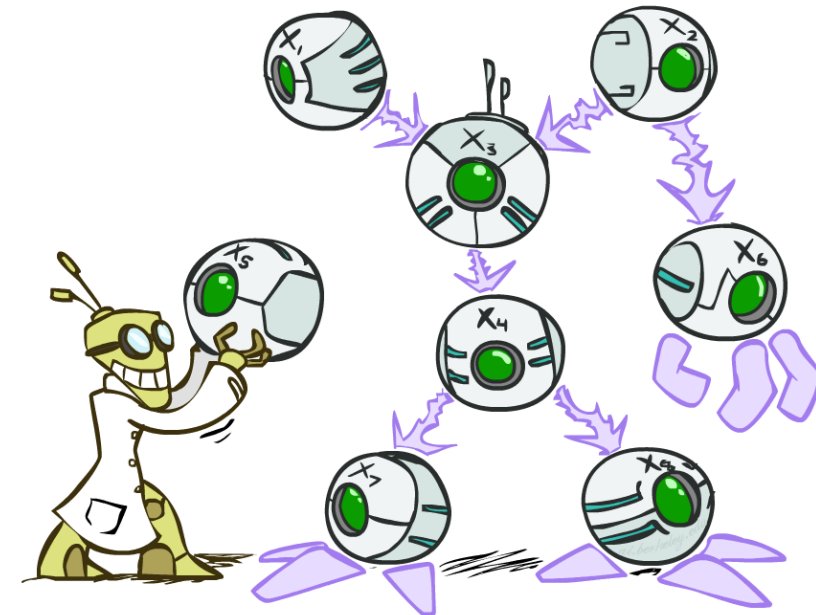
# Bayesian Networks: The Big Picture

---

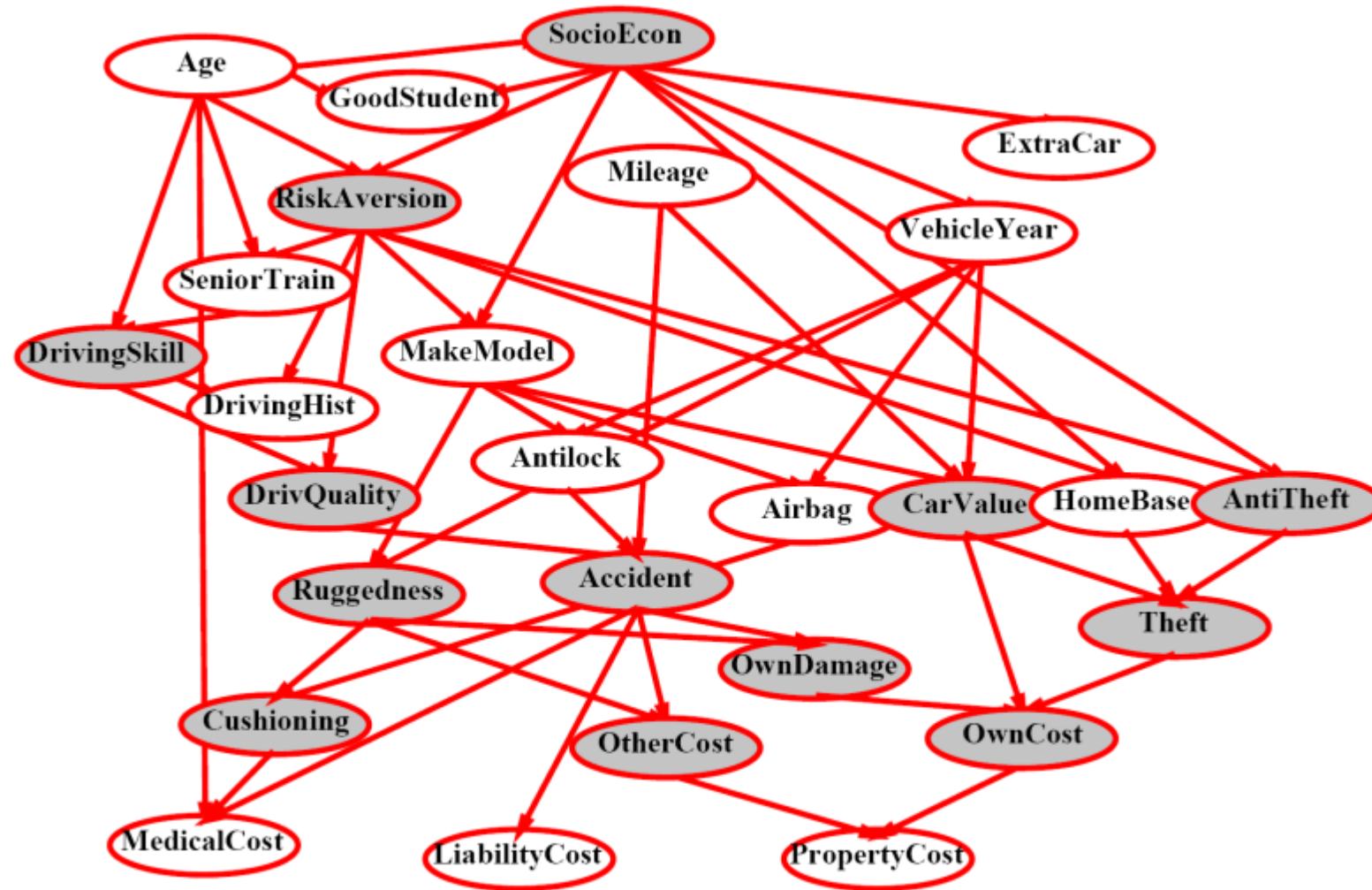


# Bayesian Networks: The Big Picture

- 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
- **Bayesian Networks:** a technique for describing complex joint distributions (models) using simple, local distributions (conditional probability tables, or CPTs)
  - More properly called **graphical models**
  - We describe how variables locally interact
  - Local interactions chain together to give global, indirect interactions

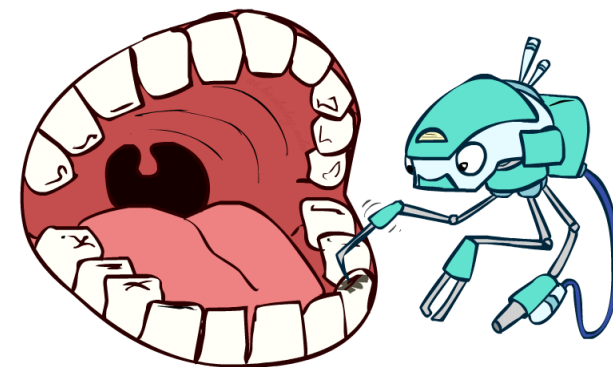
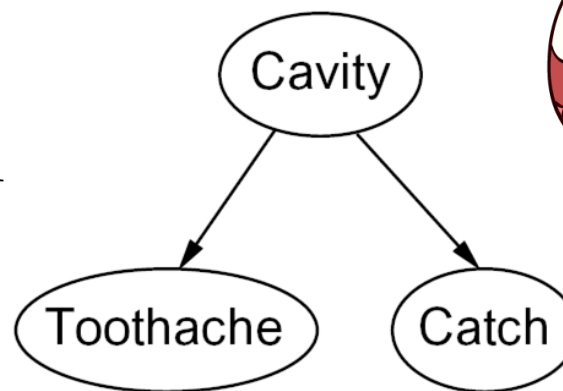
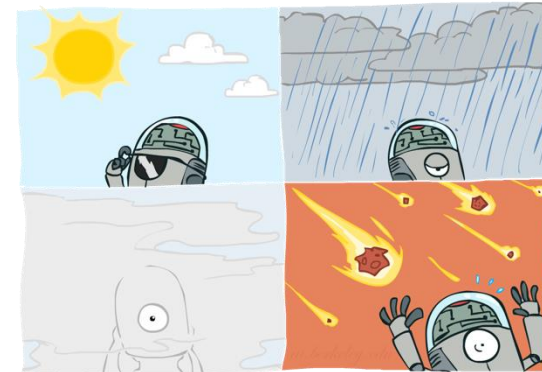


# Example Bayes Net: Insurance



# Graphical Model Notation

- Nodes: variables (with domains)
  - Can be assigned (observed) or unassigned (unobserved)
- Arcs: interactions
  - MAY indicate influence between variables
  - Formally: encode conditional independence relationships (more later)
- For now: arrows mean that there **may be** a causal relationship between the two variables



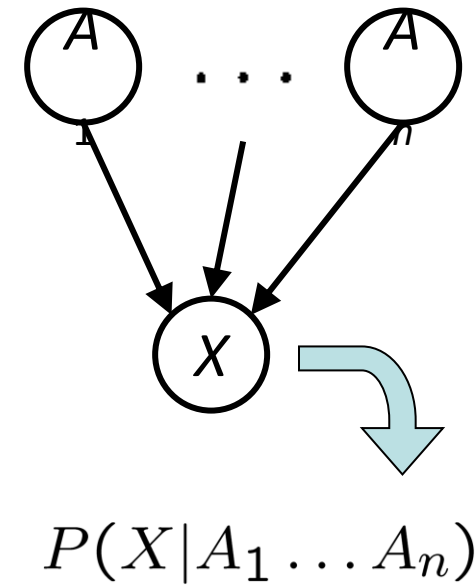
# Bayes Net Semantics



- 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

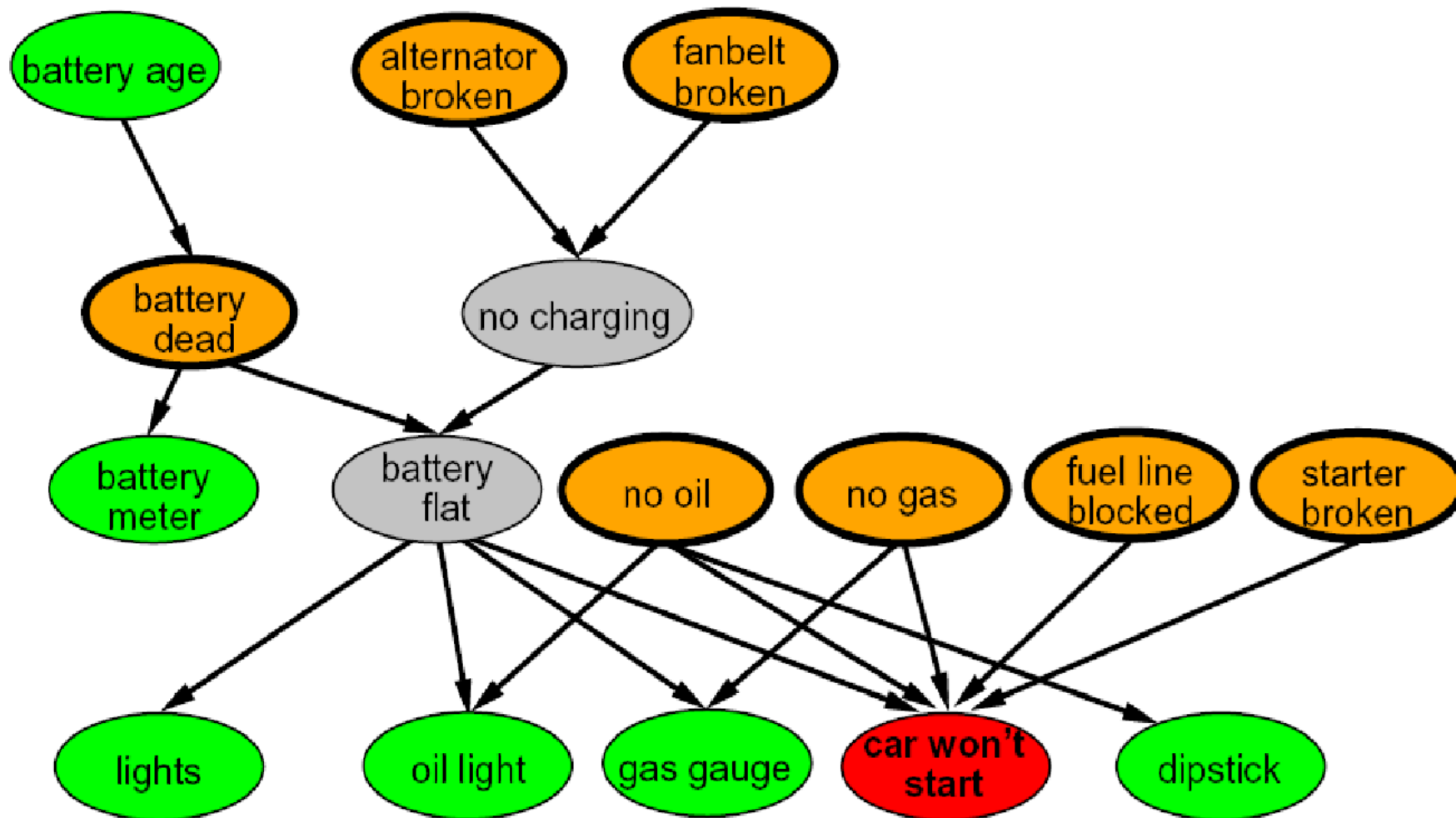
$$P(X|a_1 \dots a_n)$$

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



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

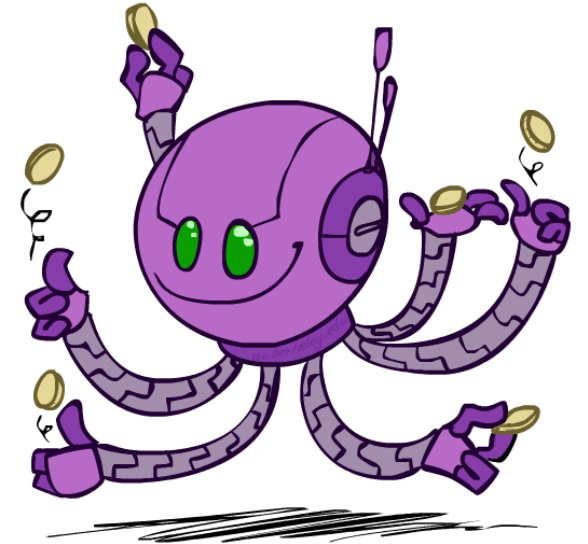
# Example Bayes Net: Car



# Example: Coin Flips

---

- N independent coin flips

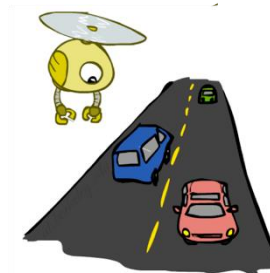


- No interactions between variables: **absolute independence**

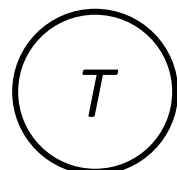
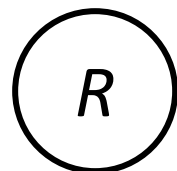


# Example: Traffic

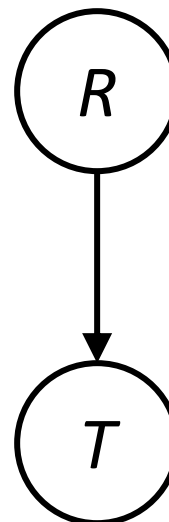
- Variables:
  - R: It rains
  - T: There is traffic



- Model 1: independence



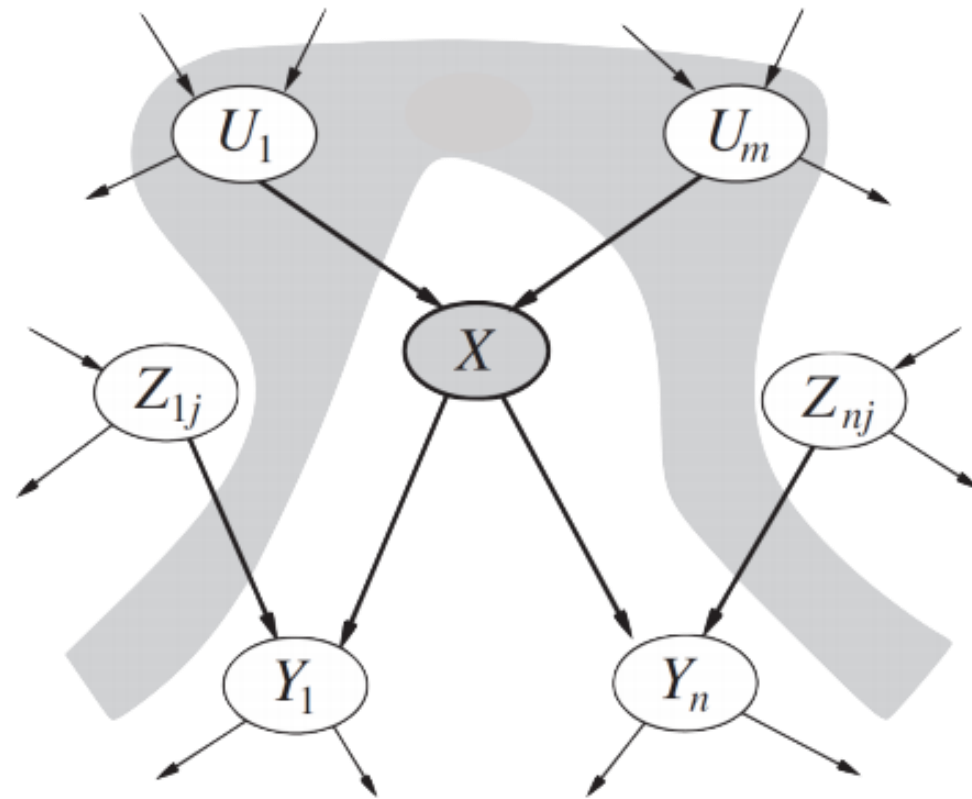
- Model 2: rain may cause traffic



- Why is an agent using model 2 better?

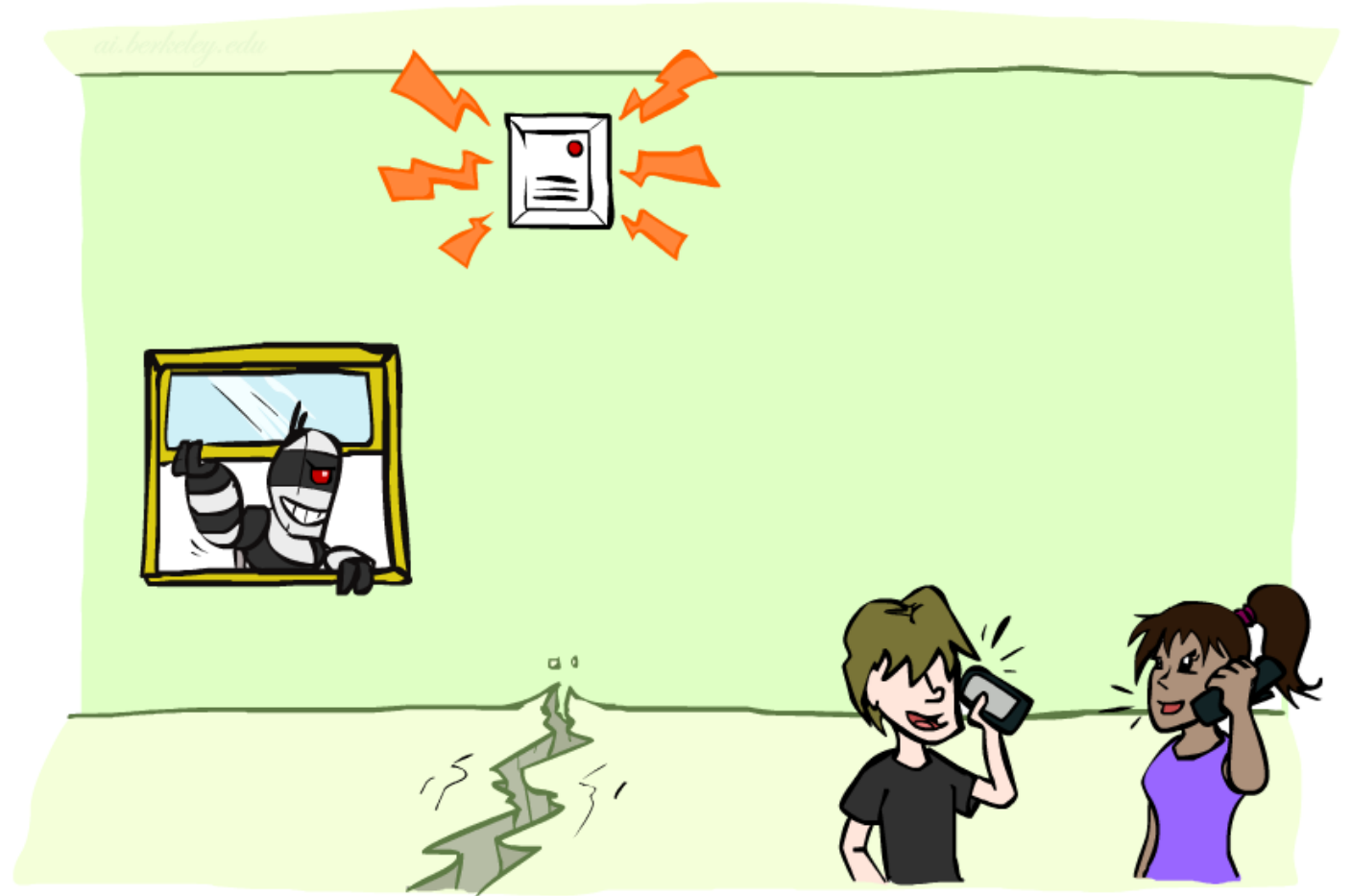
# Bayes Net: DAG + CPTs

---



# Example: Alarm Network

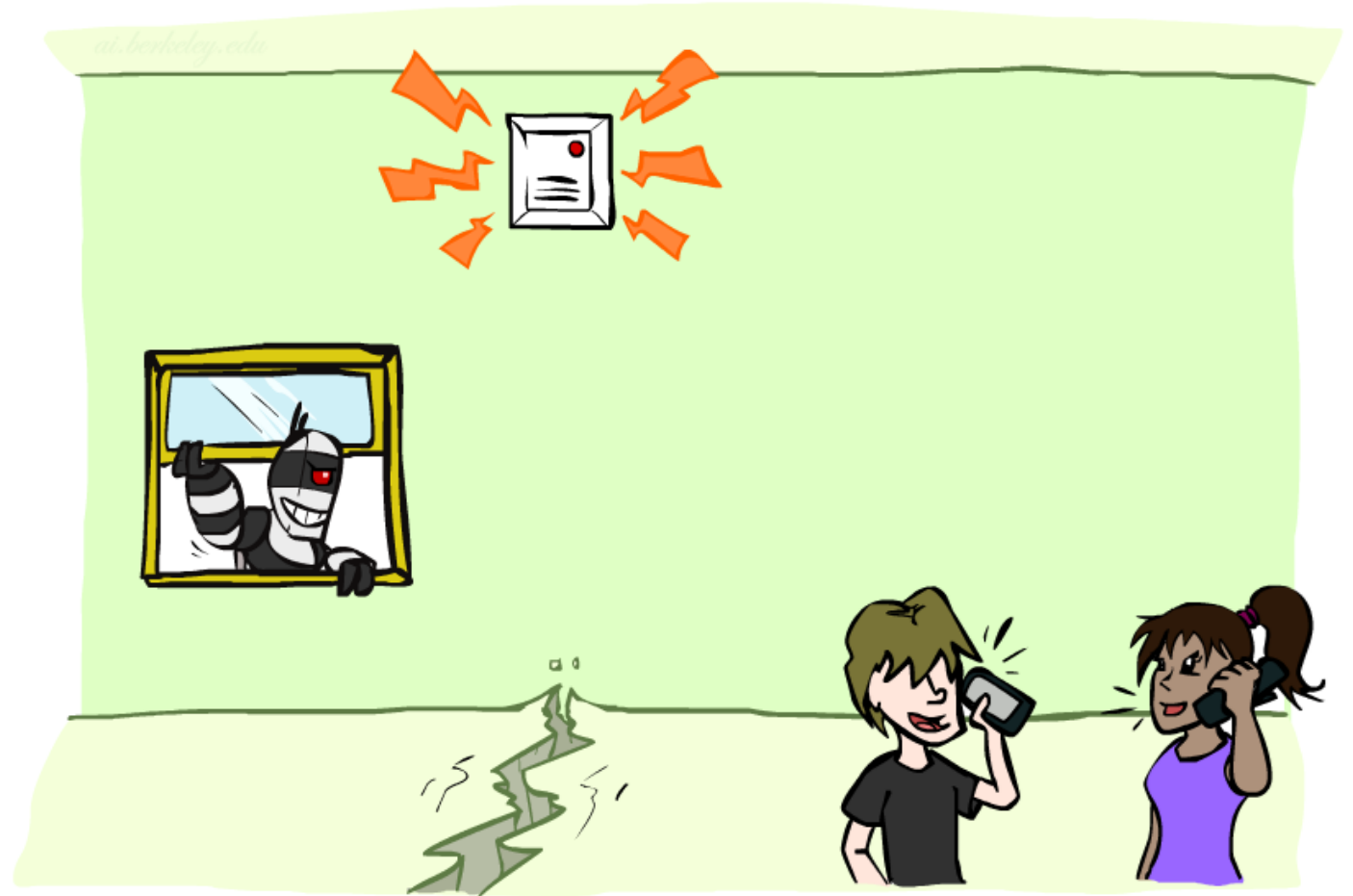
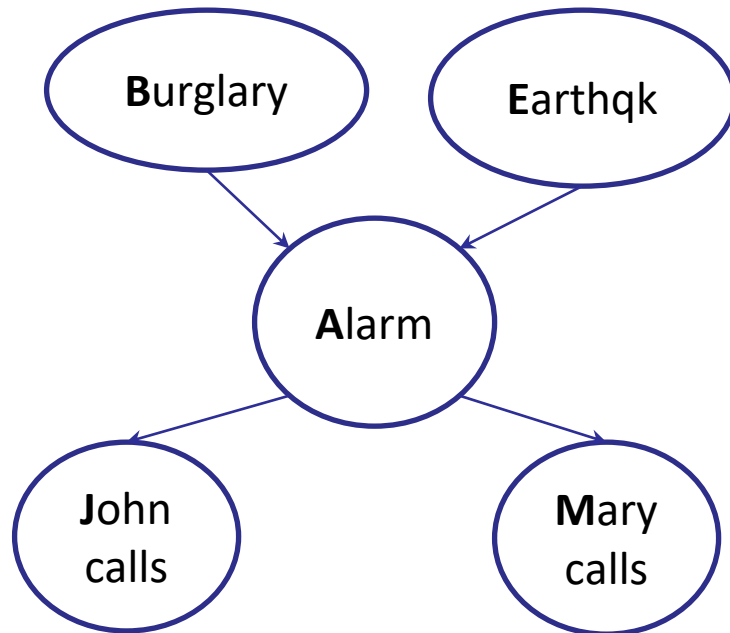
- Variables
  - B: Burglary
  - A: Alarm goes off
  - M: Mary calls
  - J: John calls
  - E: Earthquake!



# Example: Alarm Network

- Variables

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



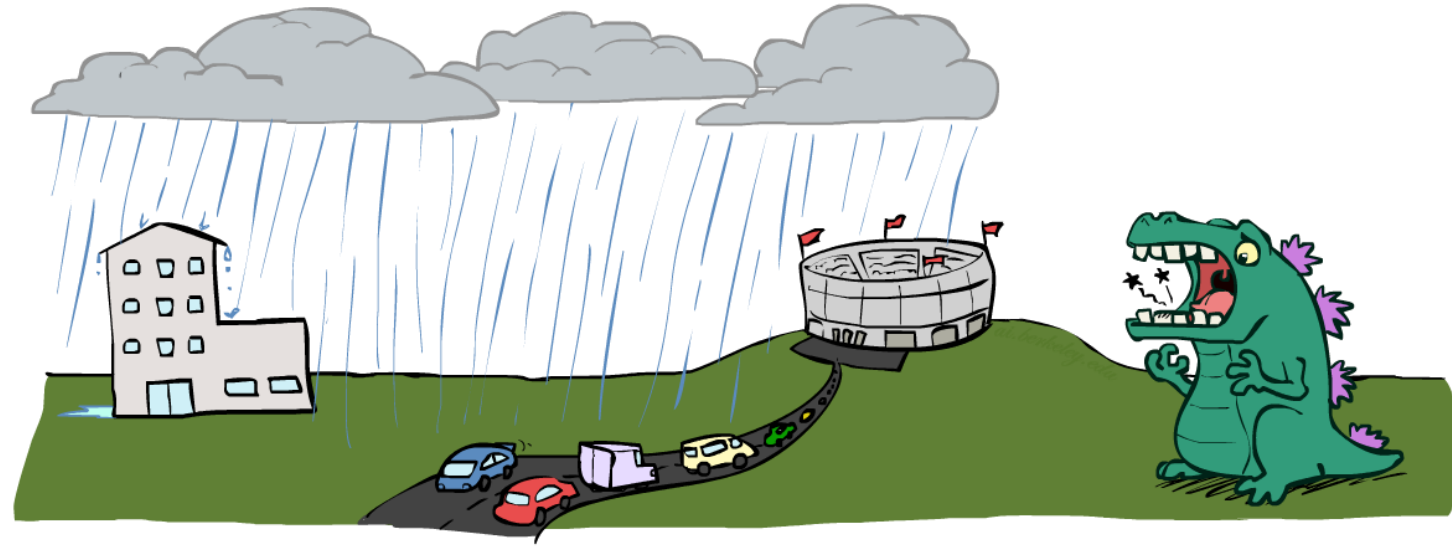
# Example: Humans

---

- G: human's goal / human's reward parameters
- S: state of the physical world
- A: human's action

# Example: Traffic II

- Variables
  - T: Traffic
  - R: It rains
  - L: Low pressure
  - D: Roof drips
  - B: Ballgame
  - C: Cavity



# Bayesian Network Semantics

---



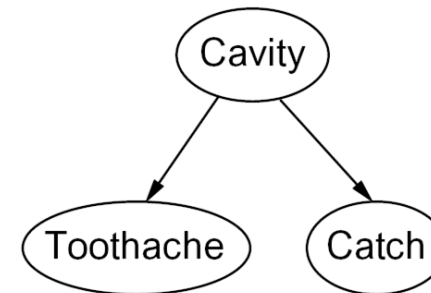
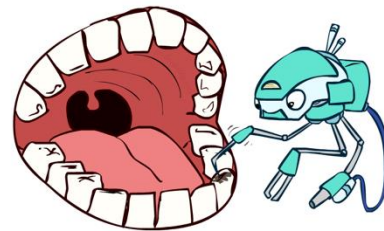
# Probabilities in BNs



- 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:



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



# 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):  $P(x_1, x_2, \dots, x_n) = \prod_{i=1}^n P(x_i | x_1 \dots x_{i-1})$

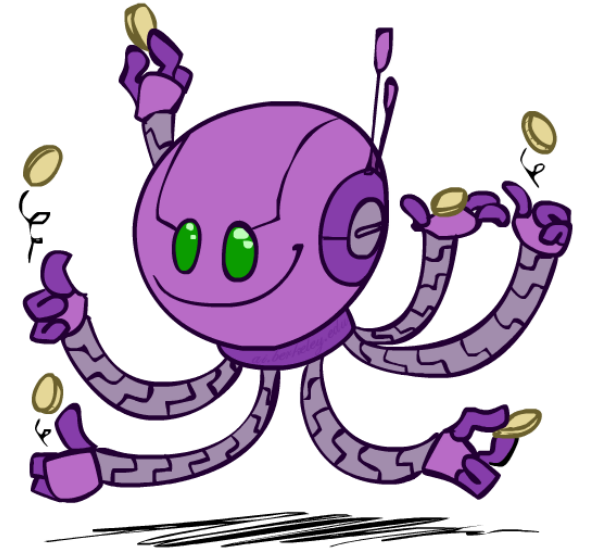
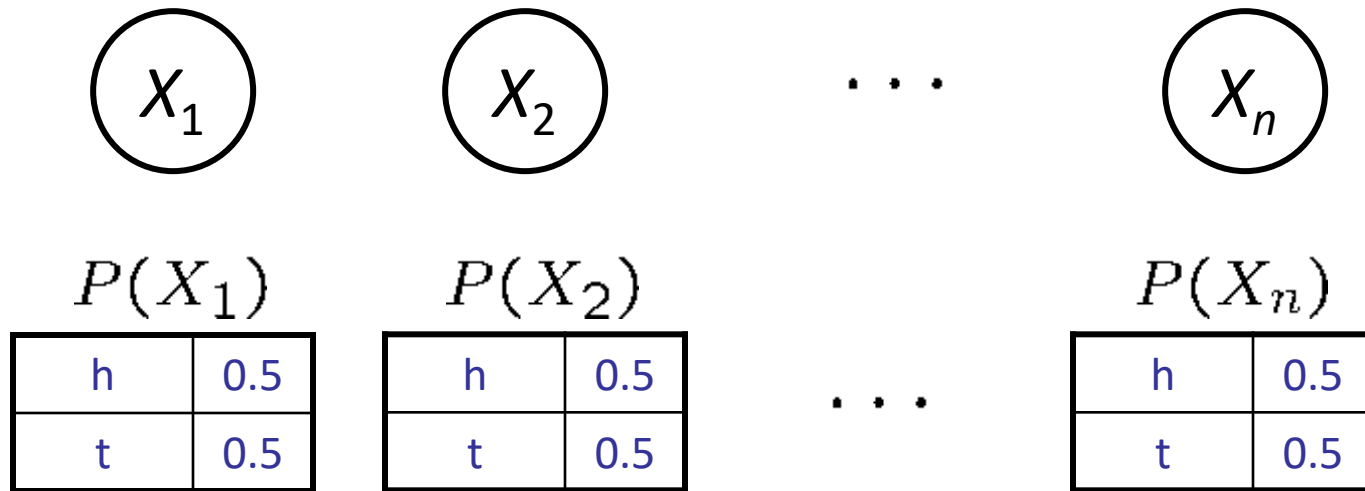
- Assume conditional independences:  $P(x_i | x_1, \dots, x_{i-1}) = P(x_i | \text{parents}(X_i))$

- ⊙ 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

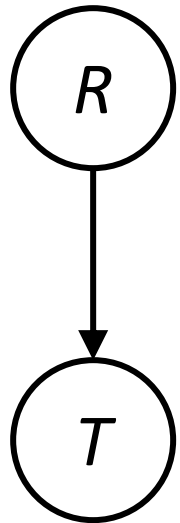
# Example: Coin Flips



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

*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

$$P(T|R)$$

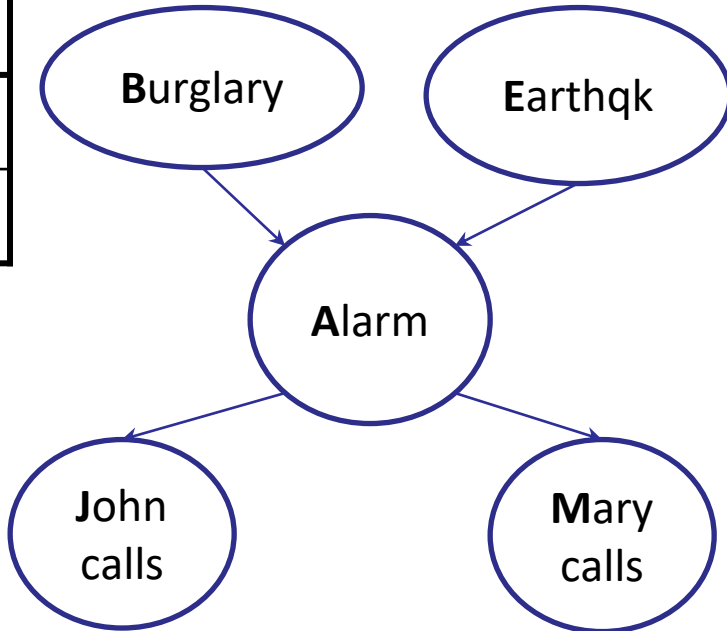
+r	+t	3/4
+r	-t	1/4
-r	+t	1/2
-r	-t	1/2

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

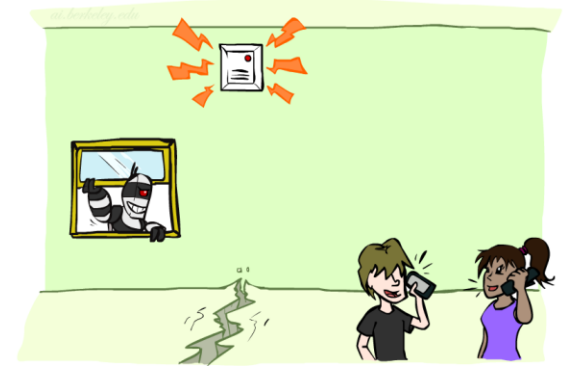


# Example: Alarm Network

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

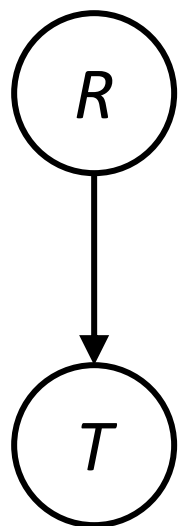
A	M	P(M A)
+a	+m	0.7
+a	-m	0.3
-a	+m	0.01
-a	-m	0.99

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

$$\begin{aligned}
 &P(M|A)P(J|A) \\
 &P(A|B,E)P(E) \\
 &P(B)
 \end{aligned}$$

# Example: Traffic

- Causal direction



$P(R)$

+r	1/4
-r	3/4

$P(T|R)$

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

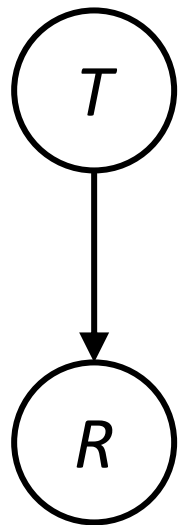
-r	+t	1/2
	-t	1/2

$P(T, R)$

+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
	-r	2/3

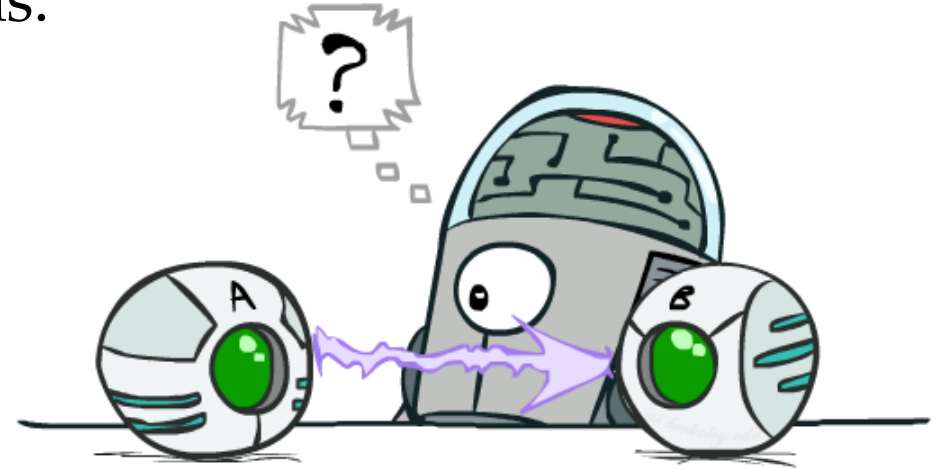
-t	+r	1/7
	-r	6/7

$P(T, R)$

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

# Causality?

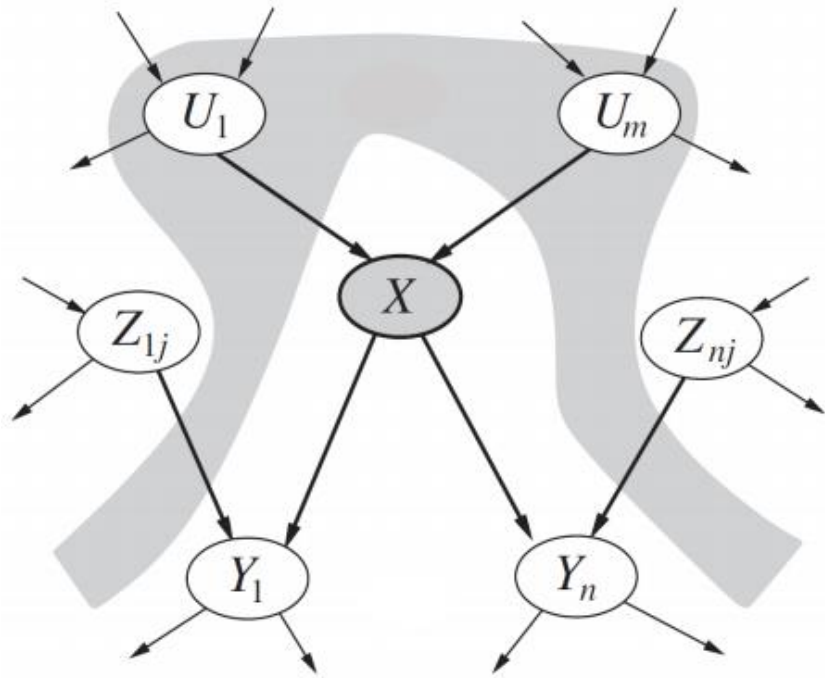
- When Bayes' nets reflect the true causal patterns:
  - Often simpler (nodes have fewer parents)
  - Often easier to think about
  - Often easier to elicit from experts
- BNs need not actually be causal
  - Sometimes no causal net exists over the domain (especially if variables are missing)
  - E.g. 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**



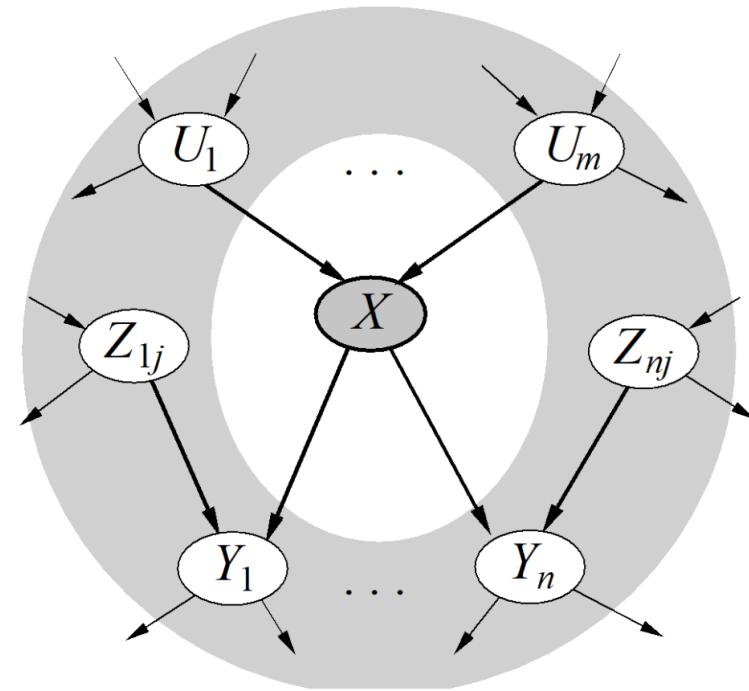
$$P(x_i|x_1, \dots, x_{i-1}) = P(x_i|\text{parents}(X_i))$$

# Conditional Independence Assumptions

- Each node, given its parents, is conditionally independent of all its non-descendants in the graph



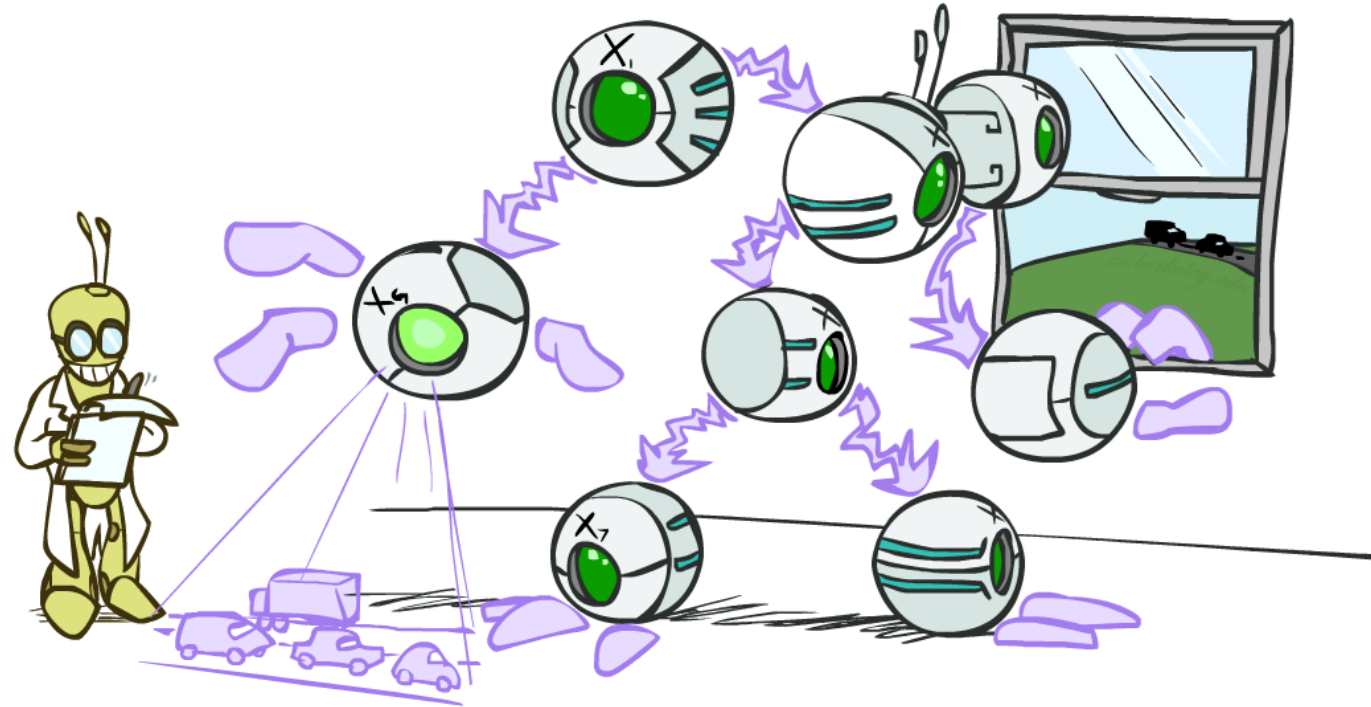
- Each node, given its MarkovBlanket, is conditionally independent of all other nodes in the graph



MarkovBlanket refers to the parents, children, and children's other parents.



# Inference with Bayesian Networks



# Inference

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

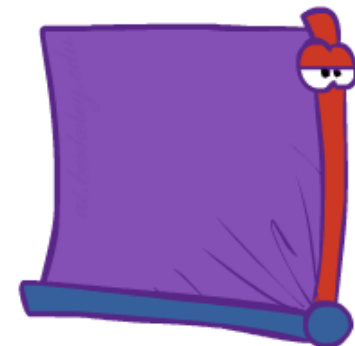
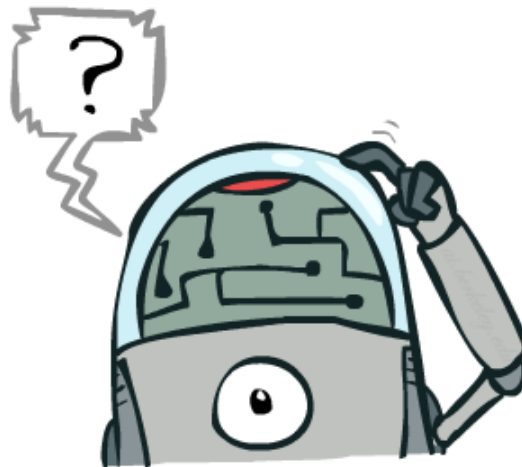
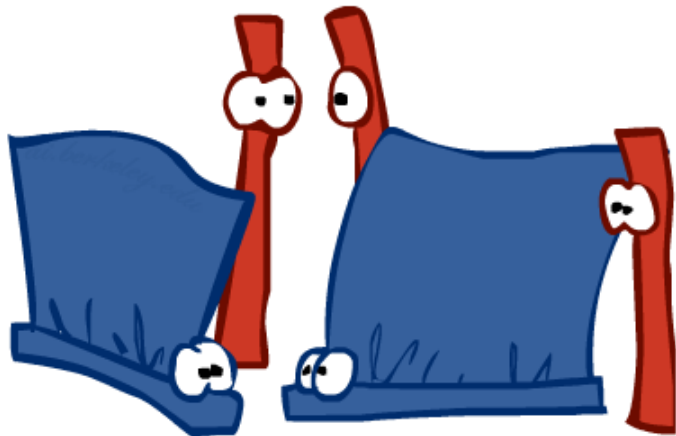
- Examples:

- 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:  $H_1 \dots H_r$
- }  $X_1, X_2, \dots, X_n$   
All variables

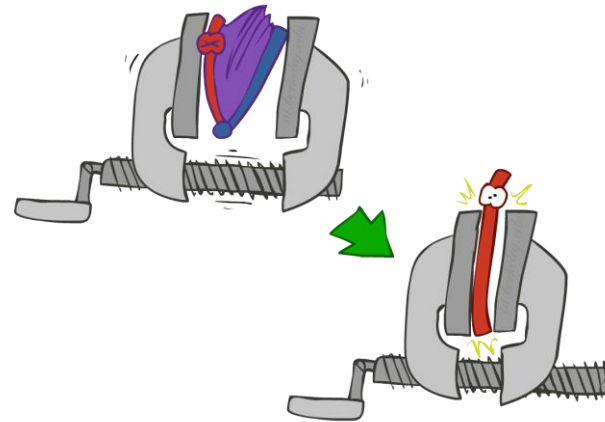
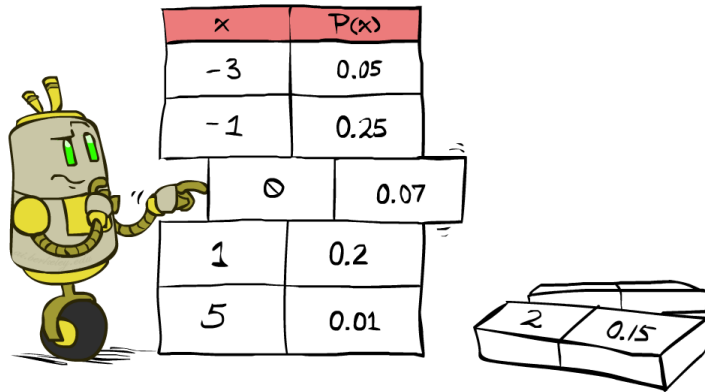
- We want:

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

- Step 1: Select the entries consistent with the evidence

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

- Step 3: Normalize



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

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

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

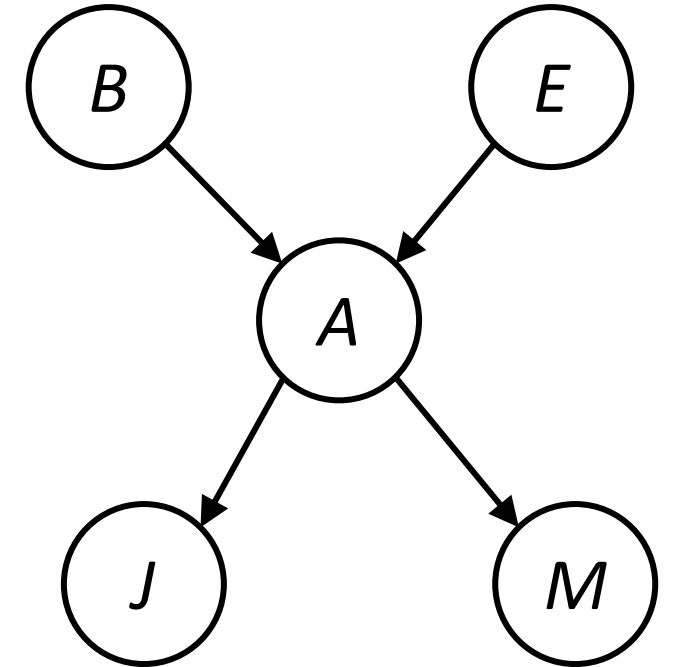
$$P(Q|e_1 \dots e_k) = \frac{1}{Z} P(Q, e_1 \dots e_k)$$

# Inference by Enumeration in Bayes' Net

---

- Given unlimited time, inference in BNs is easy

$$P(A|B, E)P(B)P(E) = P(A|B, E)P(B, E) = P(A, B, E)$$

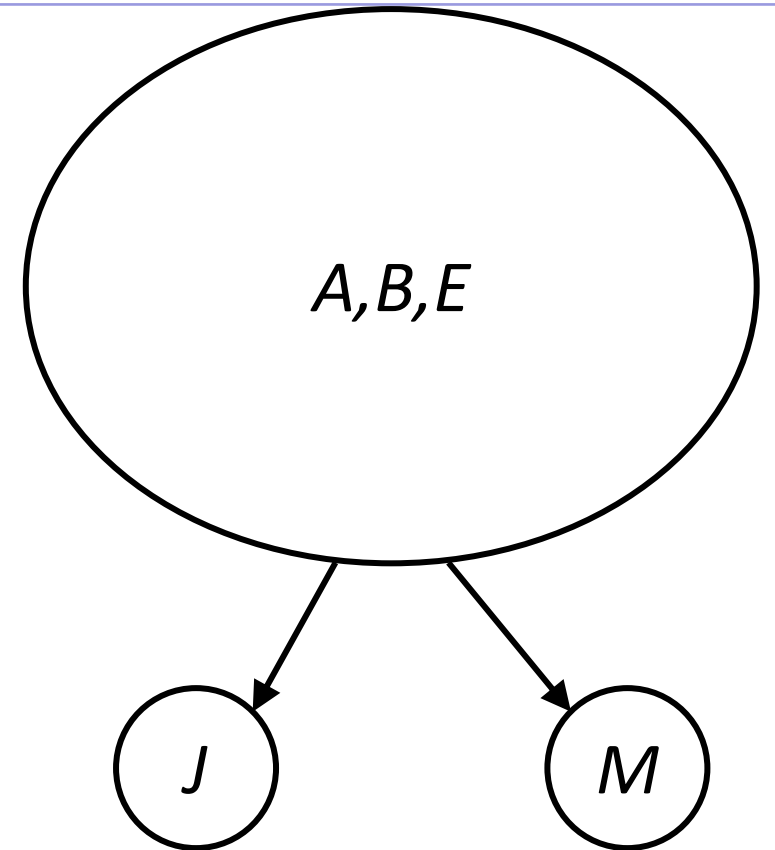


# Inference by Enumeration in Bayes' Net

---

- Given unlimited time, inference in BNs is easy

$$P(A|B, E)P(B)P(E) = P(A|B, E)P(B, E) = P(A, B, E)$$

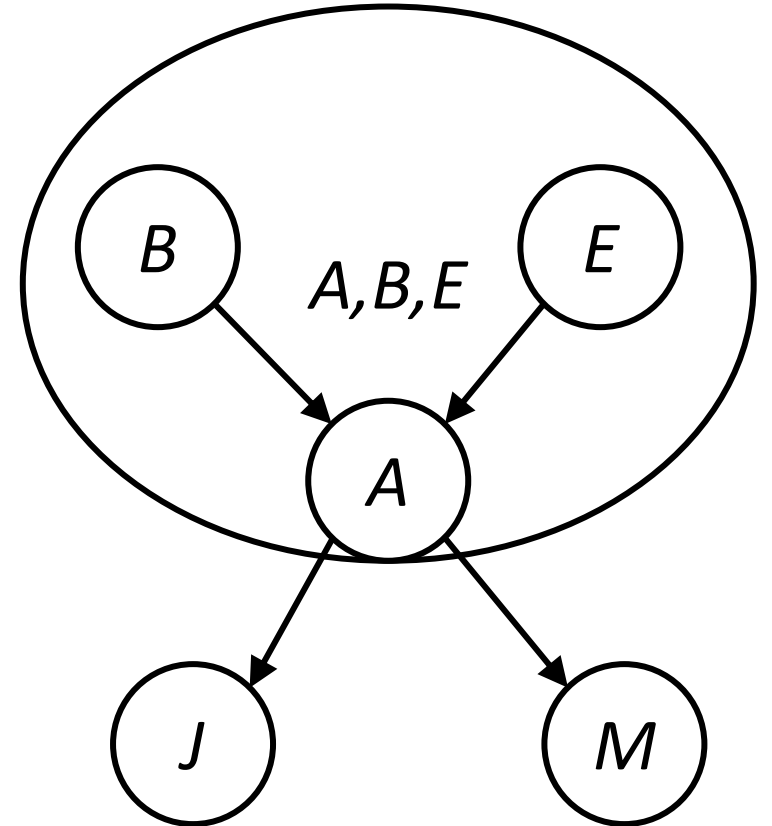


# Inference by Enumeration in Bayes' Net

- Given unlimited time, inference in BNs is easy

$$P(A|B, E)P(B)P(E) = P(A|B, E)P(B, E) = P(A, B, E)$$

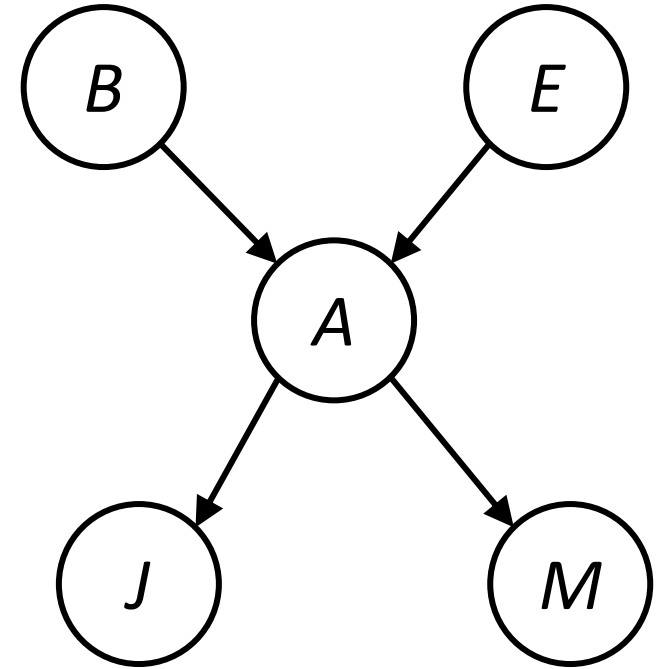
$$\begin{aligned} &P(J|A)P(M|A)P(A, B, E) \\ &= P(J, M|A)P(A, B, E) \\ &= P(J, M|A, B, E)P(A, B, E) \\ &= P(J, M, A, B, E) \end{aligned}$$



# Inference by Enumeration in Bayes' Net

- Given unlimited time, inference in BNs is easy

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



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

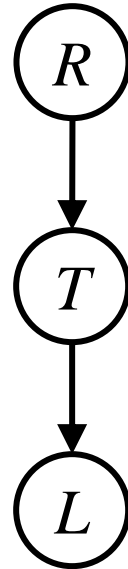
# Example: Traffic Domain

- 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

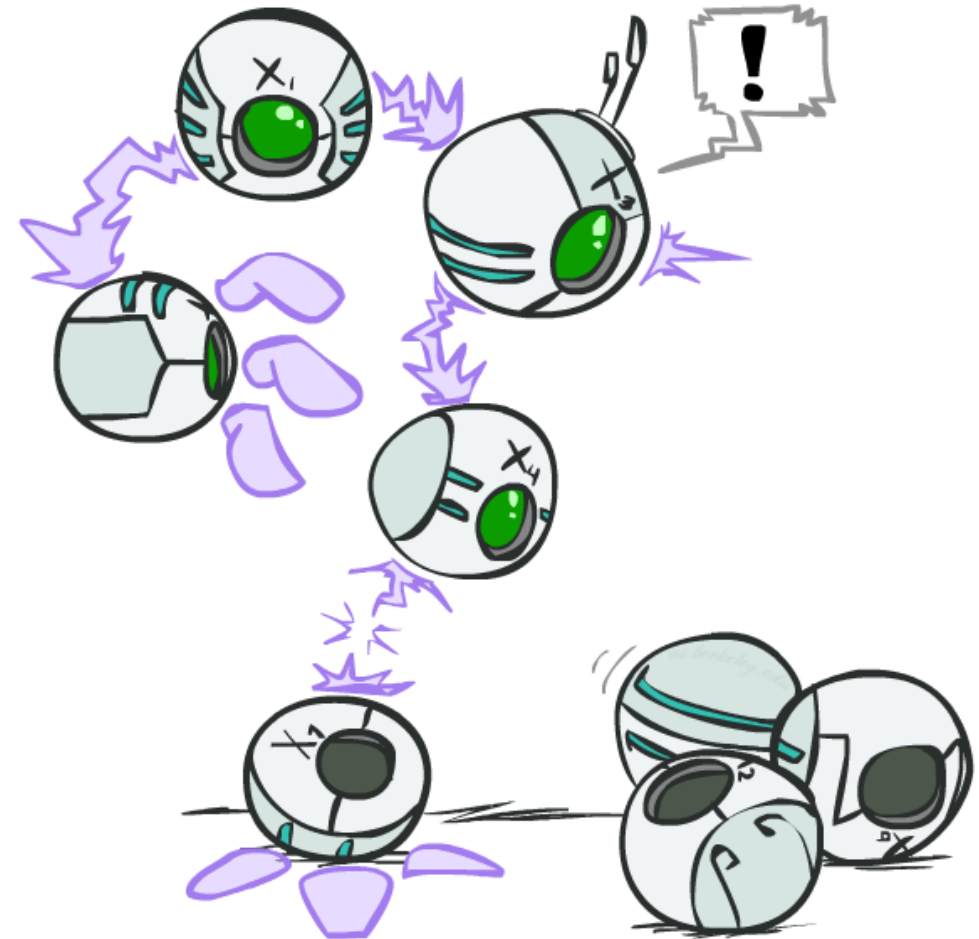
$$P(T|R)$$

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

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

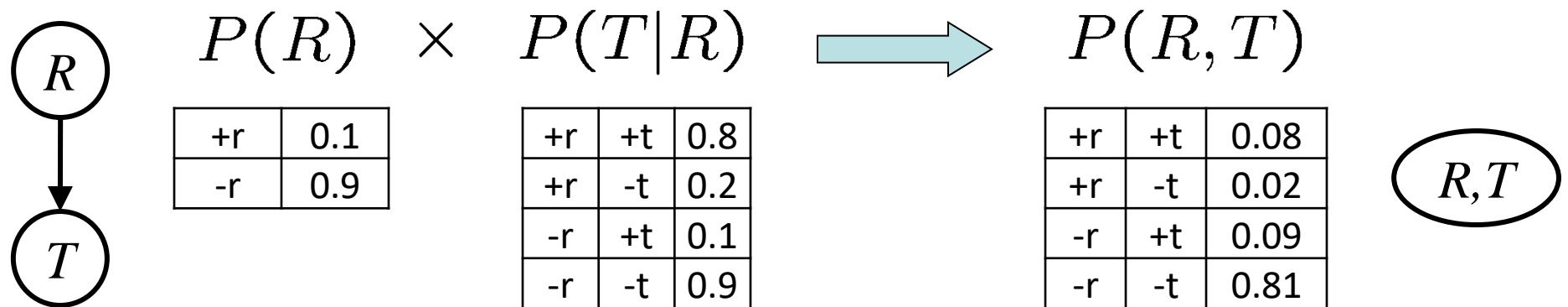
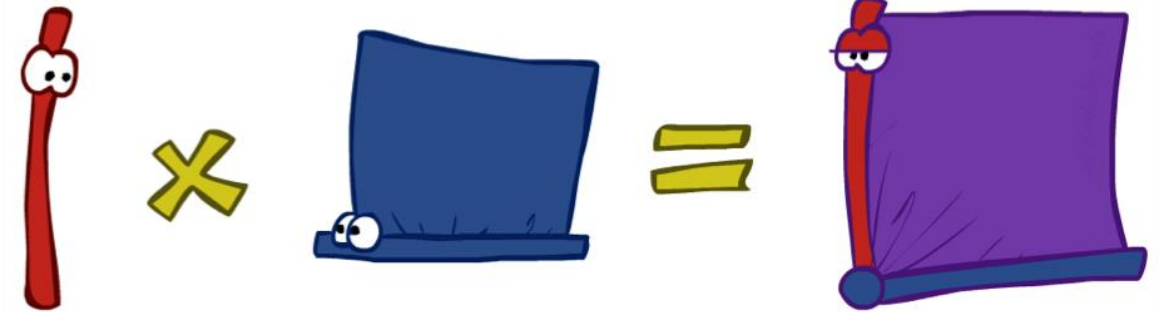
+t	+l	0.3
-t	+l	0.1

- Procedure: Join all factors, then sum out all hidden variables



# Operation 1: Join Factors

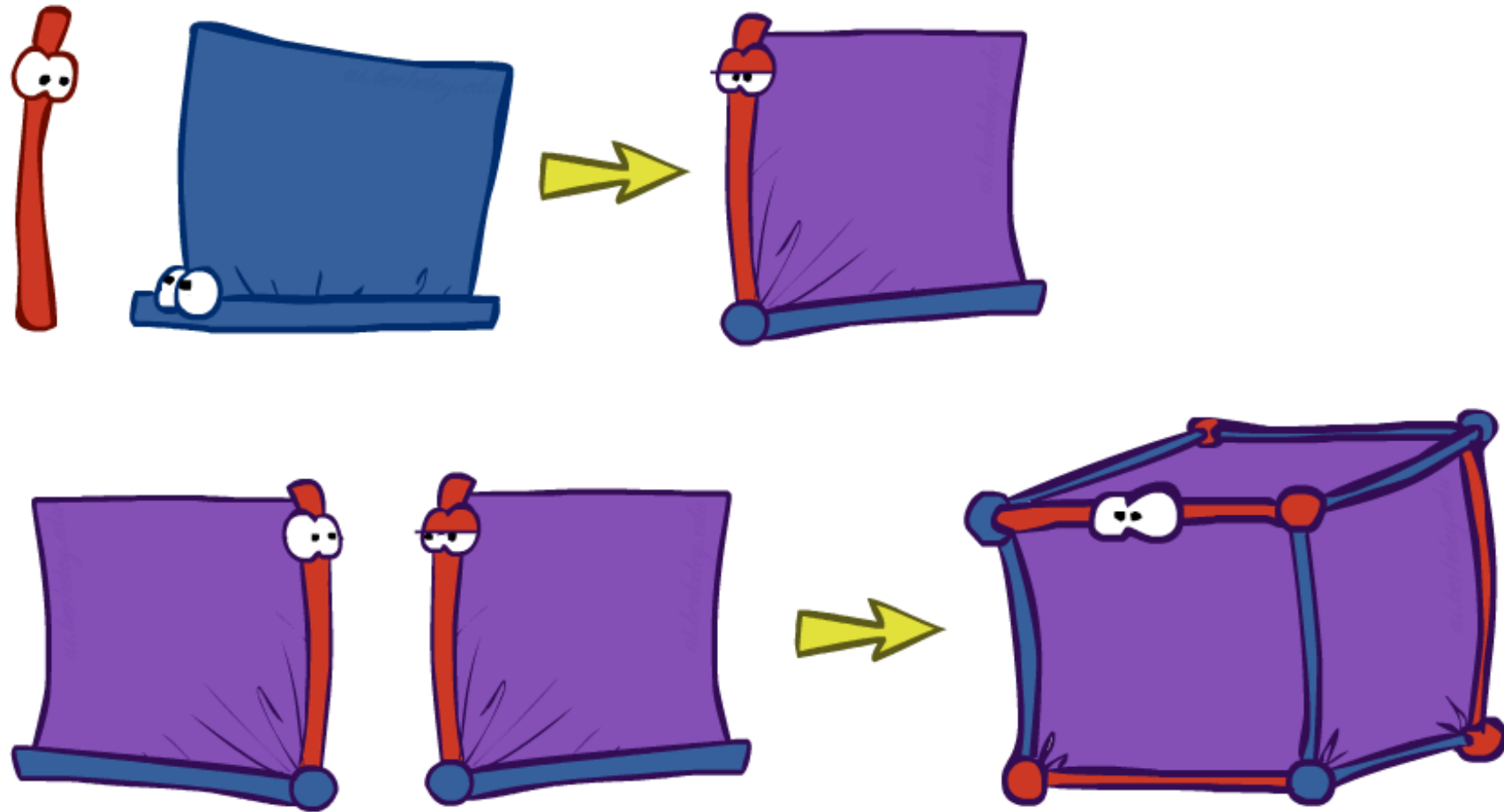
- First basic operation: **joining factors**
- Combining factors:
  - **Just like a database join**
  - Get all factors over the joining variable
  - Build a new factor over the union of the variables involved
- Example: Join on R



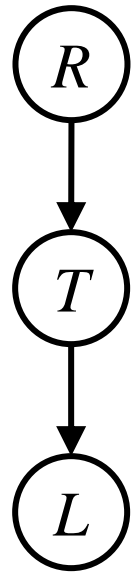
- Computation for each entry: pointwise products  $\forall r, t : P(r, t) = P(r) \cdot P(t|r)$

# Example: Multiple Joins

---



# Example: Multiple Joins



$P(R)$

+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

Join R

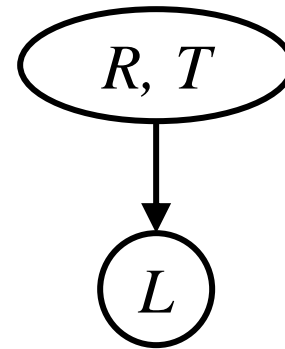


$P(R, T)$

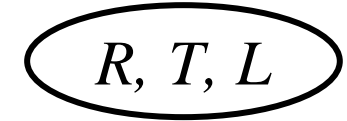
+r	+t	0.08
+r	-t	0.02
-r	+t	0.09
-r	-t	0.81

$P(L|T)$

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

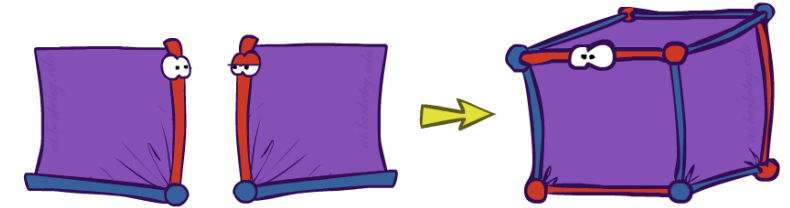
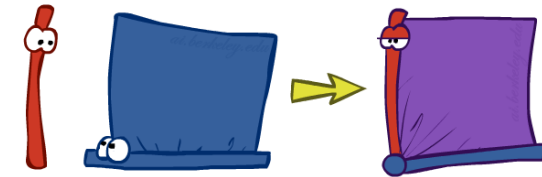


Join T



$P(R, T, L)$

+r	+t	+l	0.024
+r	+t	-l	0.056
+r	-t	+l	0.002
+r	-t	-l	0.018
-r	+t	+l	0.027
-r	+t	-l	0.063
-r	-t	+l	0.081
-r	-t	-l	0.729



# 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:

$P(R, T)$

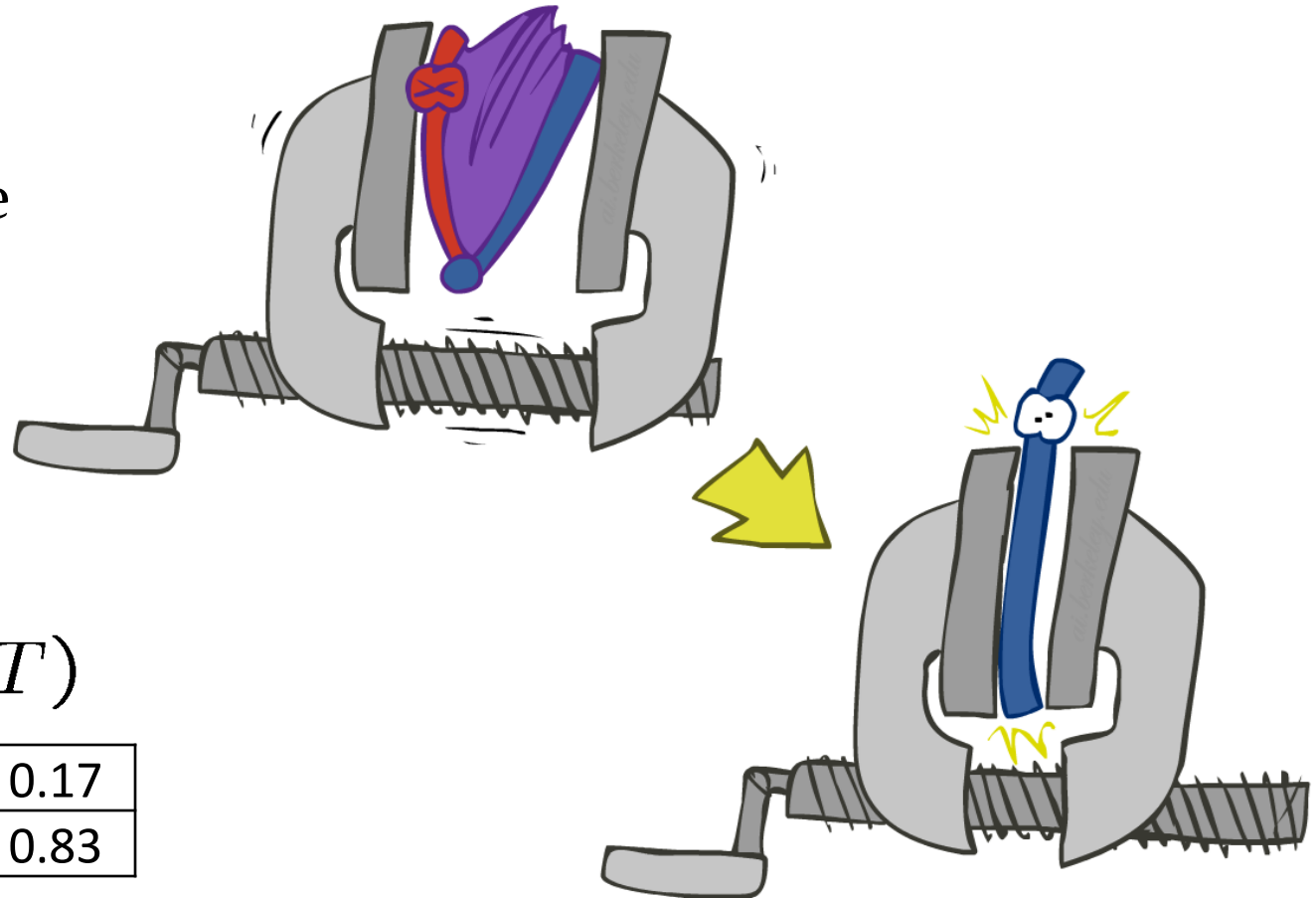
+r	+t	0.08
+r	-t	0.02
-r	+t	0.09
-r	-t	0.81

sum  $R$



$P(T)$

+t	0.17
-t	0.83



# Multiple Elimination

$R, T, L$

$T, L$

$L$

$P(R, T, L)$

+r	+t	+l	0.024
+r	+t	-l	0.056
+r	-t	+l	0.002
+r	-t	-l	0.018
-r	+t	+l	0.027
-r	+t	-l	0.063
-r	-t	+l	0.081
-r	-t	-l	0.729

Sum  
out R



$P(T, L)$

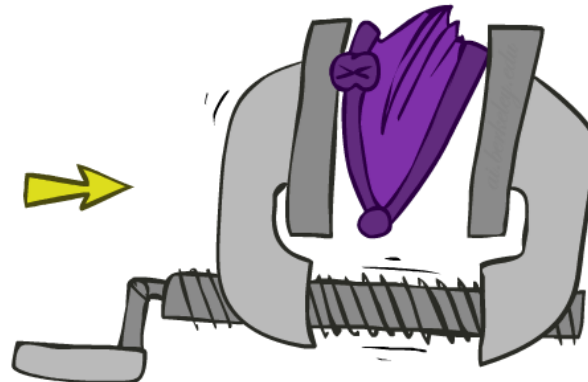
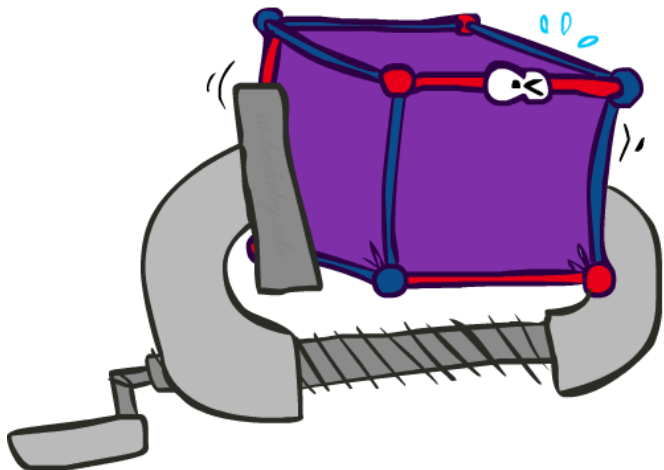
+t	+l	0.051
+t	-l	0.119
-t	+l	0.083
-t	-l	0.747

Sum  
out T



$P(L)$

+l	0.134
-l	0.866



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

---

$P(R)$

$P(T|R)$



$P(R, T, L)$



$P(L)$

$P(L|T)$

# Recap: Inference by Enumeration

- General case:

- Evidence variables:  $E_1 \dots E_k = e_1 \dots e_k$
  - Query\* variable:  $Q$
  - Hidden variables:  $H_1 \dots H_r$
- }  $X_1, X_2, \dots, X_n$   
All variables

- We want:

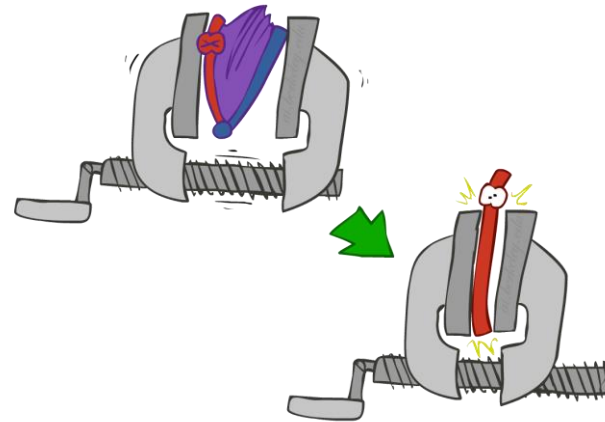
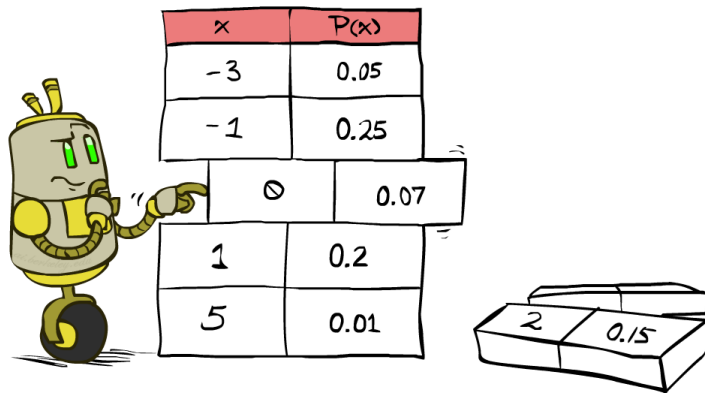
$$P(Q|e_1 \dots e_k)$$

*\* Works fine with multiple query variables, too*

- Step 1: Select the entries consistent with the evidence

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

- Step 3: Normalize



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

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

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

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



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

---

$$P(R)$$

$$P(T|R)$$

$$P(L|T)$$



- Compute joint

$$P(R, T, L)$$

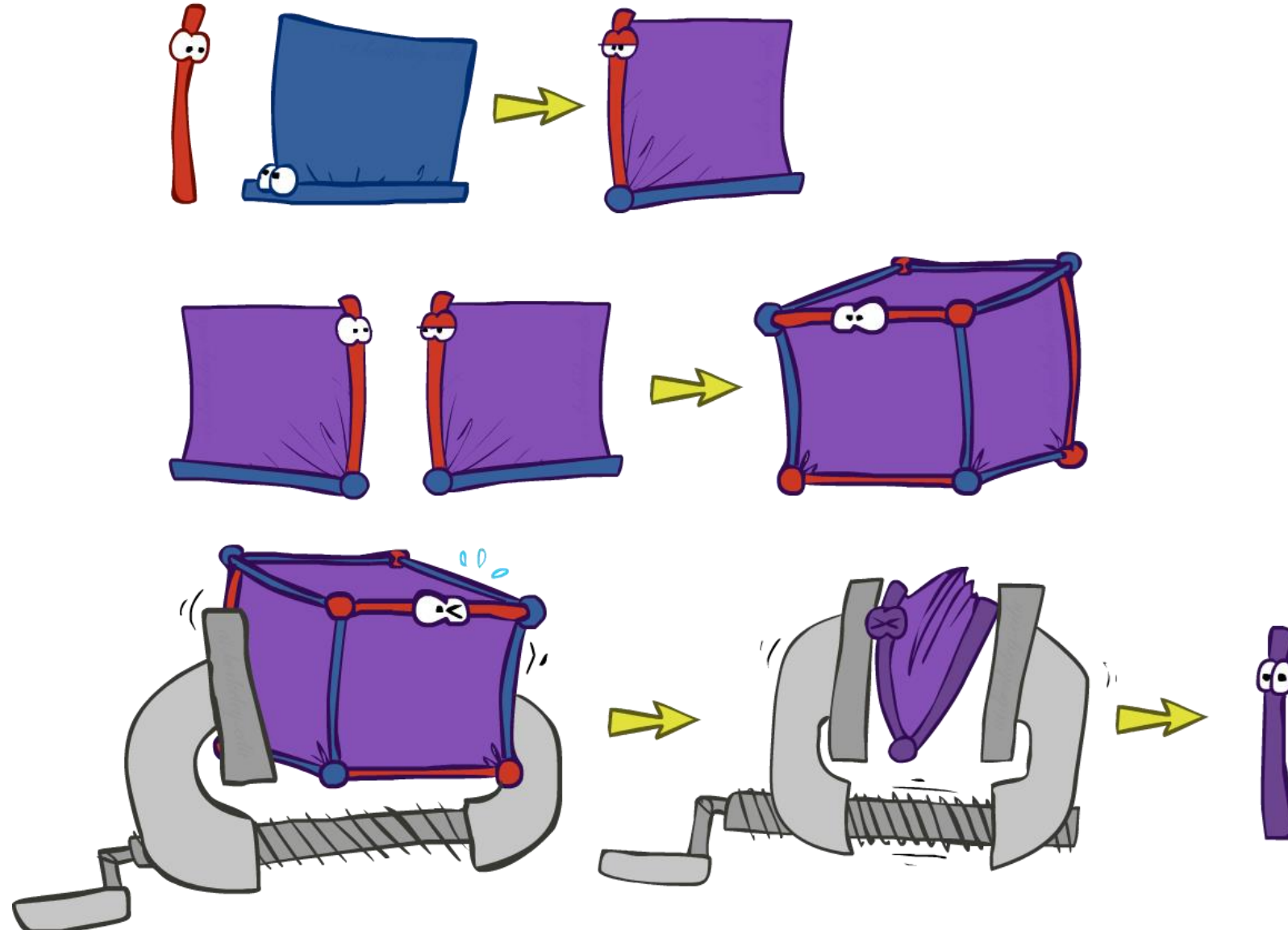
- Sum out hidden variables

$$P(L)$$

- [Step 3: Normalize]

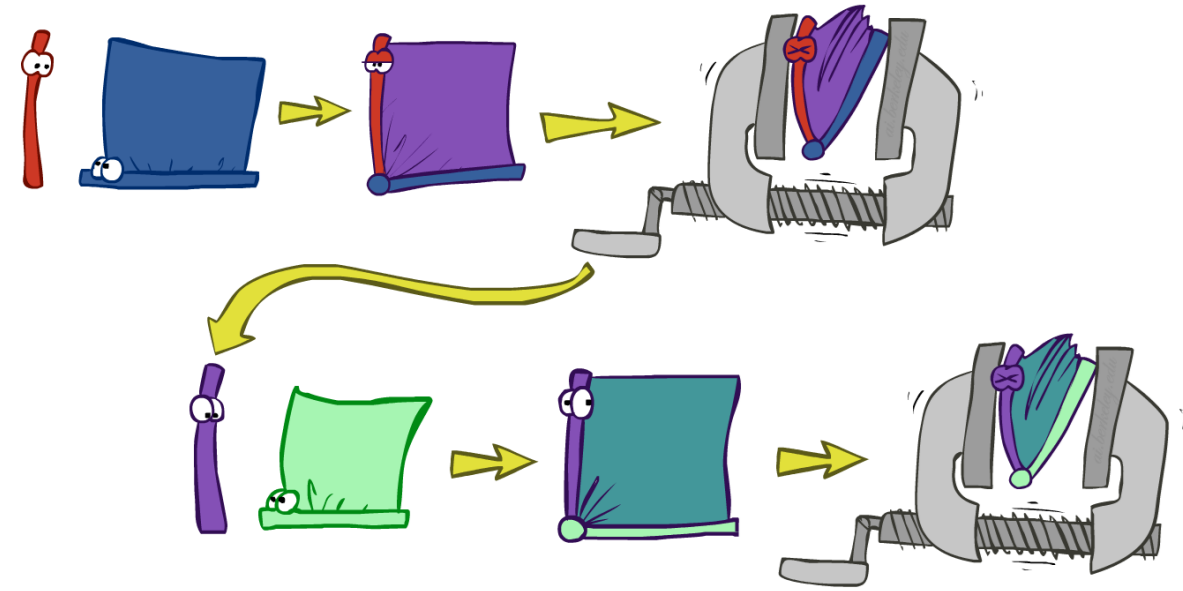
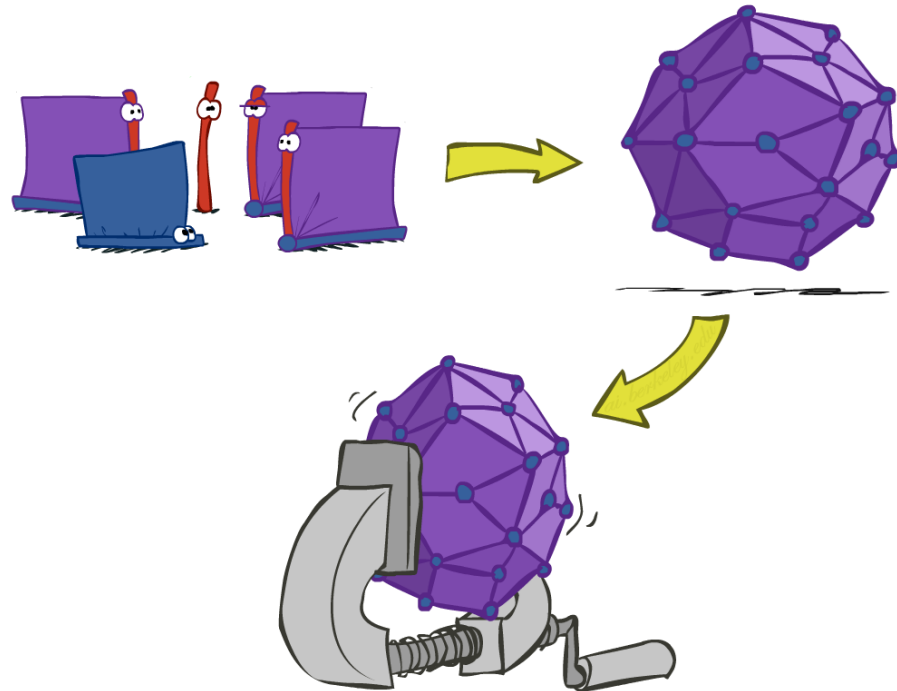
# Thus Far: Multiple Join, Multiple Eliminate (= Inference by Enumeration)

---

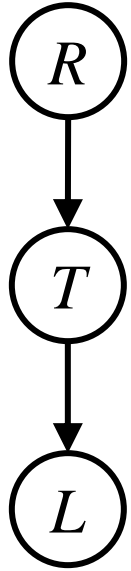


# Inference by Enumeration vs. Variable Elimination

- Why is inference by enumeration 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



# Traffic Domain



$$P(L) = ?$$

- Inference by Enumeration
- Variable Elimination

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

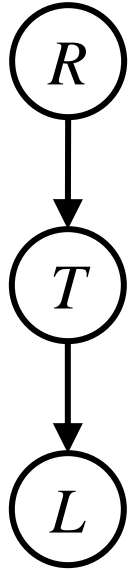
Diagram illustrating Inference by Enumeration. Red brackets and labels show the process of eliminating variables: 'Eliminate r' (bracket around  $P(r)P(t|r)$ ), 'Join on t' (bracket around  $P(L|t)P(t|r)$ ), and 'Eliminate t' (bracket around the final sum).

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

Diagram illustrating Variable Elimination. Red brackets and labels show the process of eliminating variables: 'Join on r' (bracket around  $P(r)P(t|r)$ ), 'Eliminate r' (bracket around the inner sum), 'Join on t' (bracket around  $P(L|t)$ ), and 'Eliminate t' (bracket around the final sum).

# Traffic Domain

---



$$P(L) = ?$$

- Inference by Enumeration
- Variable Elimination

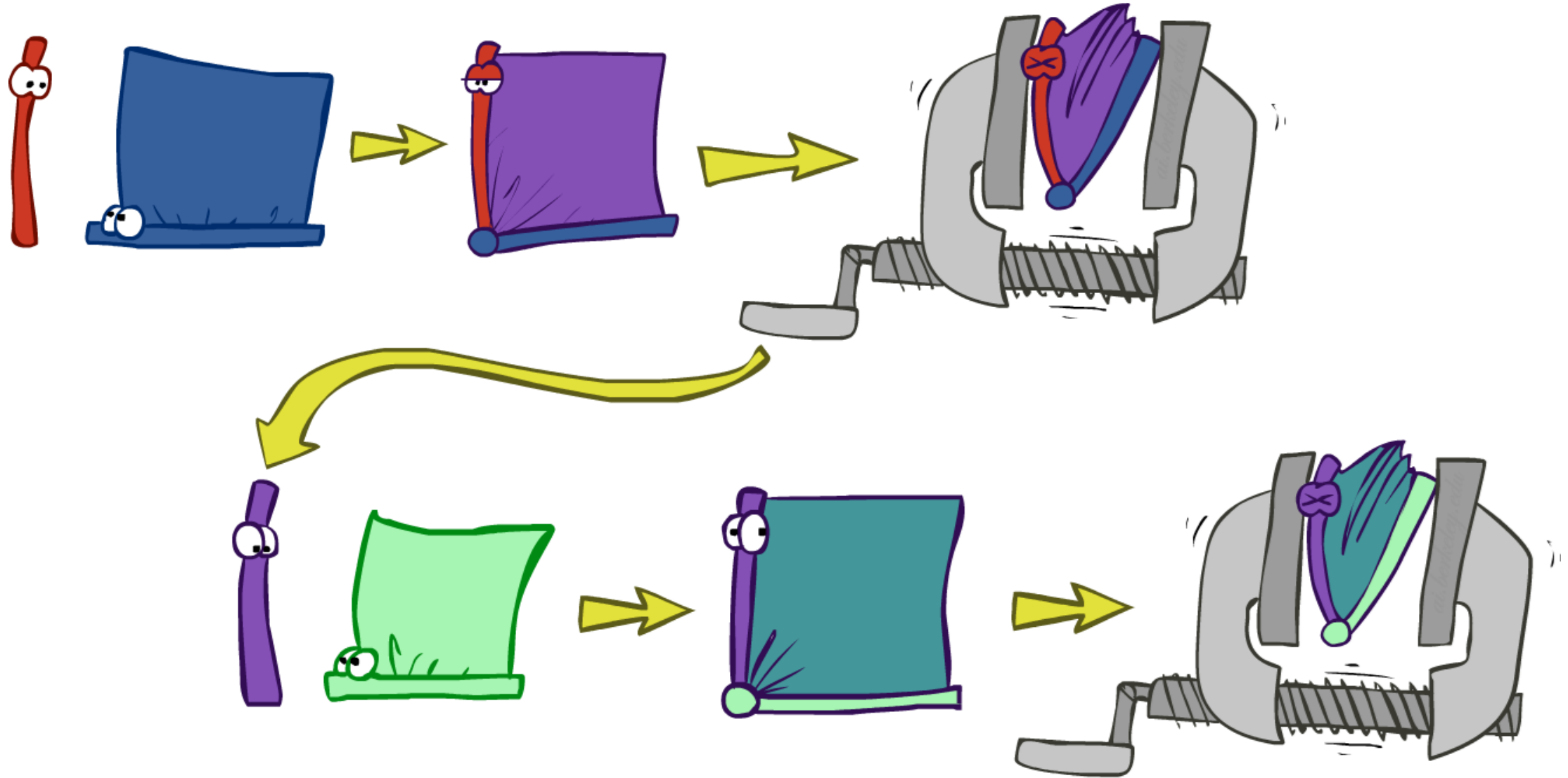
$$= \sum_t \sum_r P(L|t)P(r)P(t|r)$$

$$(5a) + (5b)$$

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

$$5(a + b)$$

# Marginalizing Early (Variable Elimination)



# Variable Elimination



$P(R)$

T	L
+r	0.1
-r	0.9

$P(T|R)$

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

$P(L|T)$

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

Join R

$P(R, T)$

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

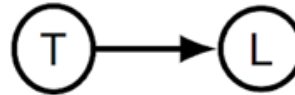


$P(L|T)$

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

Sum out R

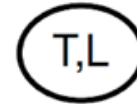
	+t	-t
	0.17	0.83



$P(L|T)$

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

Join T



$P(T, L)$

	+l	-l
+t	0.051	0.119
-t	0.083	0.747

Sum out T



	+l	-l
	0.134	0.866

# 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
----	-----

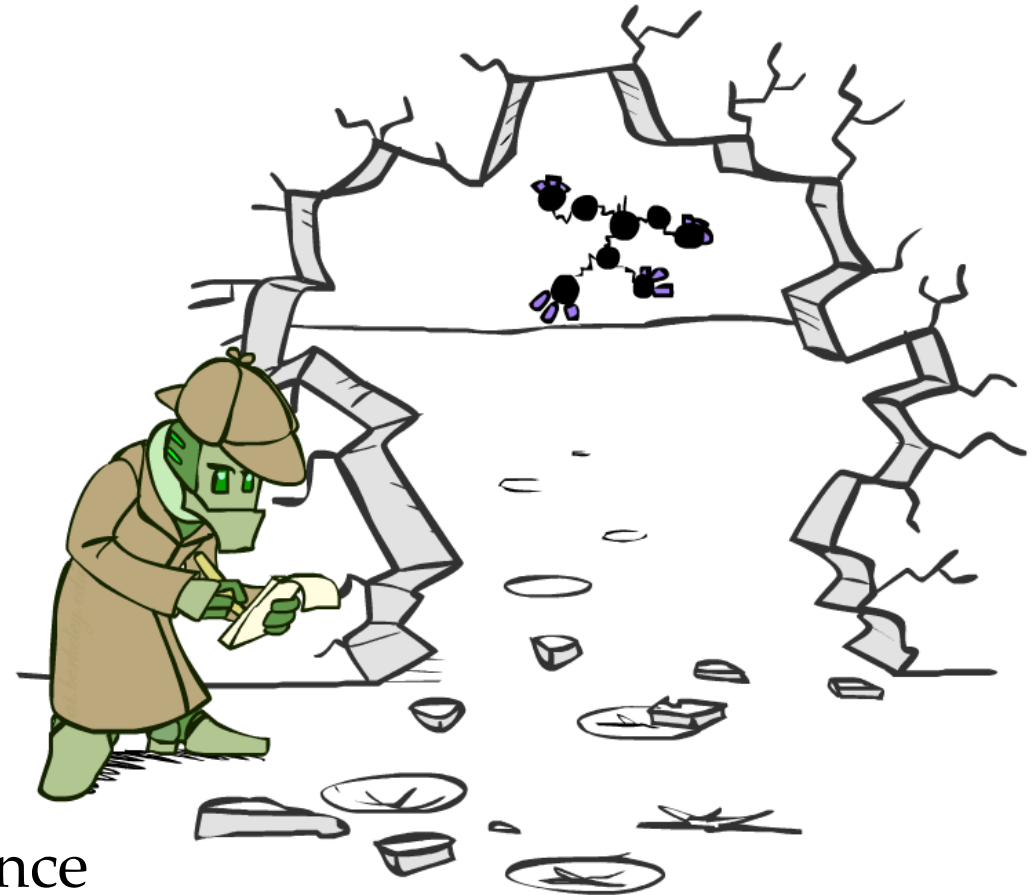
$$P(T|+r)$$

+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

- We eliminate all vars other than query + evidence





# Evidence

- 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

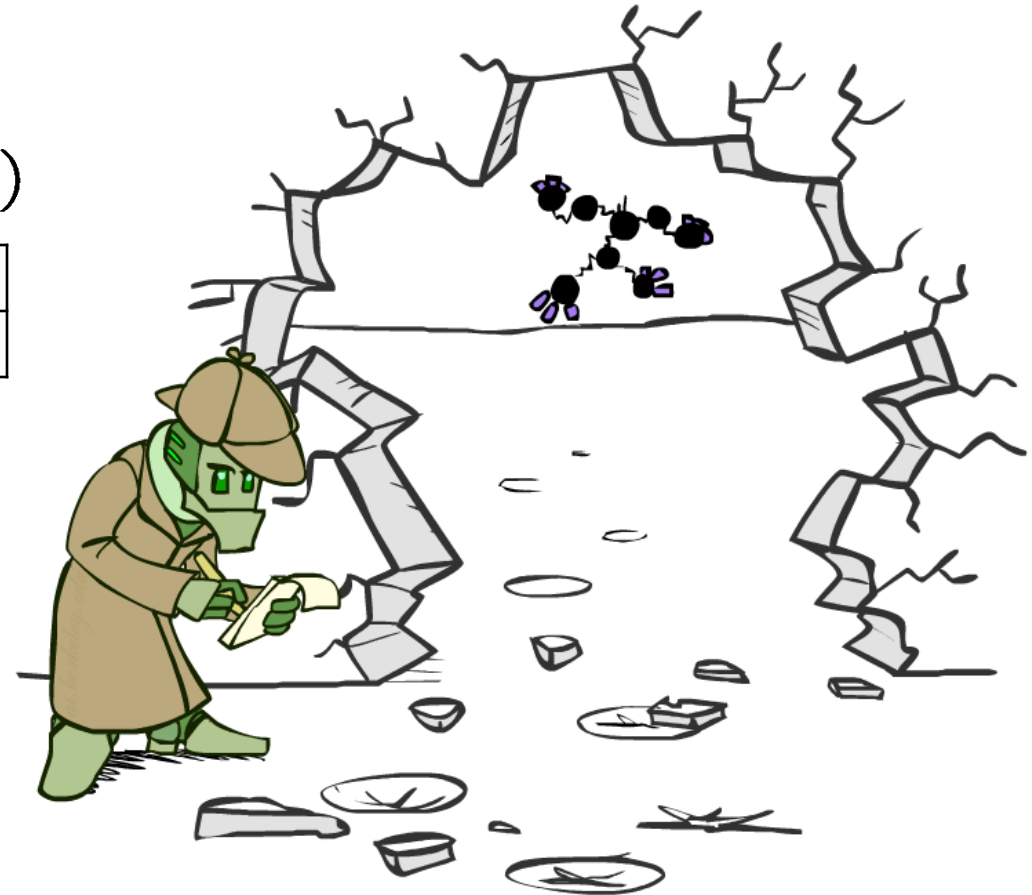
Normalize



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

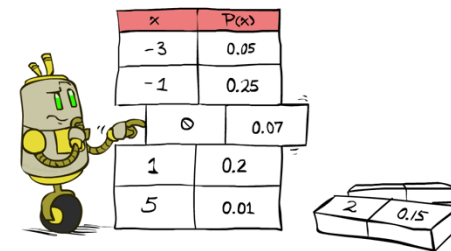
+l	0.26
-l	0.74

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



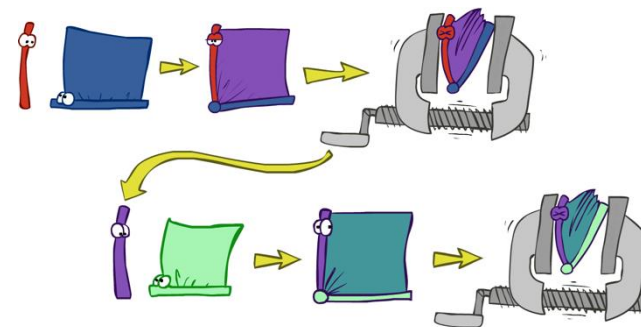
# General Variable Elimination

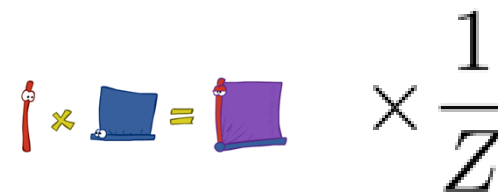
- Query:  $P(Q|E_1 = e_1, \dots, E_k = e_k)$
- Start with initial factors:
  - Local CPTs (but instantiated by evidence)
- While there are still hidden variables (not Q or evidence):
  - Pick a hidden variable H
  - Join all factors mentioning H
  - Eliminate (sum out) H
- Join all remaining factors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



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

2 0.15

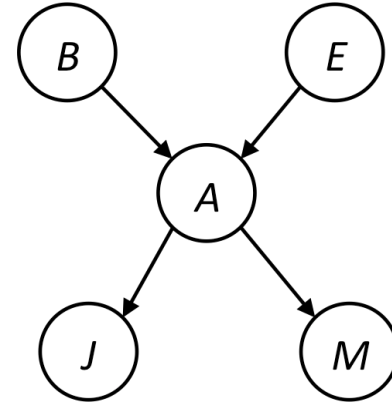



$$\text{stick} \times \text{blue square} = \text{purple square} \times \frac{1}{Z}$$

# Example

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

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



$$\begin{aligned}
 P(B|j, m) &\propto P(B, j, m) \\
 &= \sum_{e, a} P(B, j, m, e, a) \\
 &= \sum_{e, a} P(B)P(e)P(a|B, e)P(j|a)P(m|a) \\
 &= \sum_{e, a} P(B)P(e) \sum_a P(a|B, e)P(j|a)P(m|a) \\
 &= \sum_e P(B)P(e) f_1(j, m|B, e) \\
 &= P(B) \sum_e P(e) f_1(j, m|B, e) \\
 &= P(B) f_2^e(j, m|B)
 \end{aligned}$$

marginal can be obtained from joint by summing

use Bayes' net joint distribution expression

use  $x^*(y+z) = xy + xz$

joining on a, and then summing out gives  $f_1$

use  $x^*(y+z) = xy + xz$

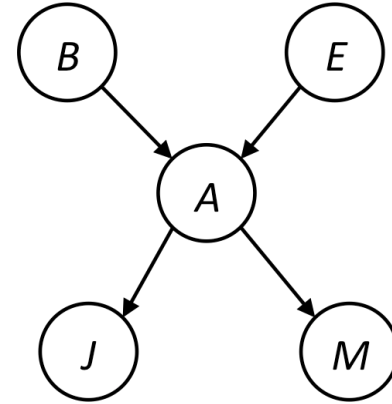
joining on e, and then summing out gives  $f_2$

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

# Example

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

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

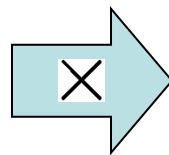


Choose A

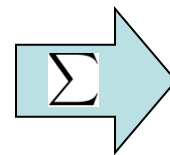
$$P(A|B, E)$$

$$P(j|A)$$

$$P(m|A)$$



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



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

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

# Example

$$P(B) \quad P(E) \quad P(j, m|B, E)$$

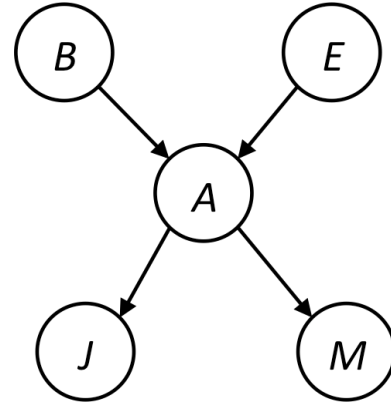
Choose E

$$\begin{array}{l} P(E) \\ P(j, m|B, E) \end{array} \xrightarrow{\times} P(j, m, E|B) \xrightarrow{\Sigma} P(j, m|B)$$

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

Finish with B

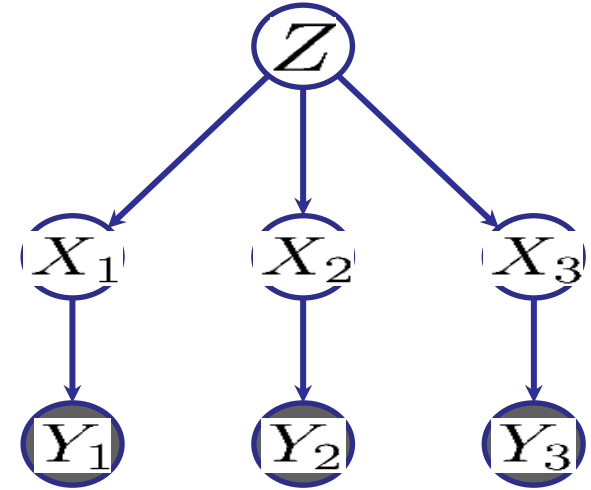
$$\begin{array}{l} P(B) \\ P(j, m|B) \end{array} \xrightarrow{\times} P(j, m, B) \xrightarrow{\text{Normalize}} P(B|j, m)$$



# Variable Elimination Example

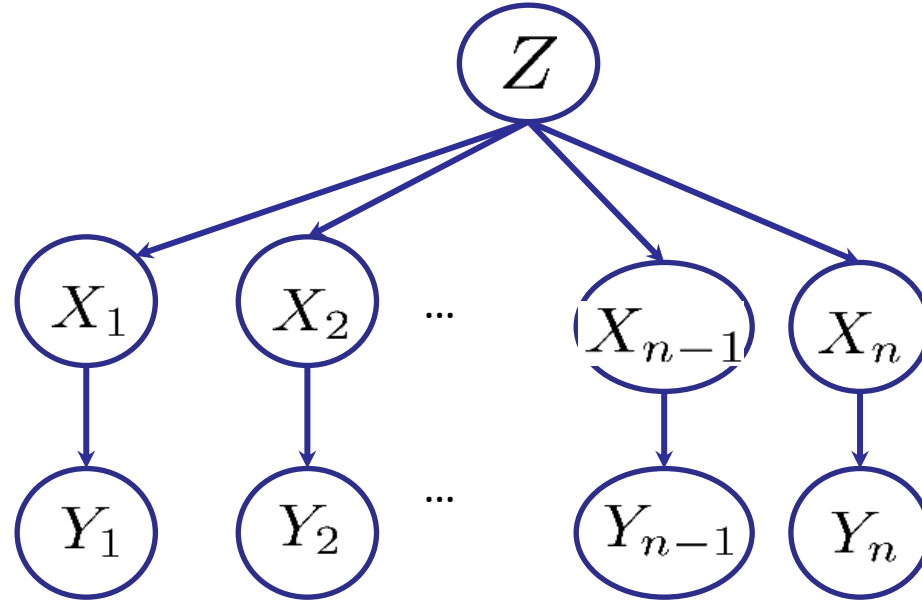
---

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



# 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:  $2^n$  versus 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!**



# “Easy” Structures: 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!!