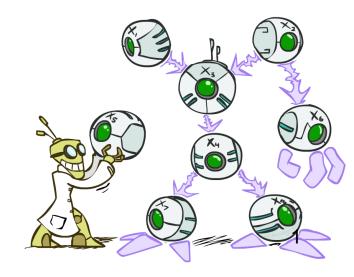
# CSE 573: Artificial Intelligence

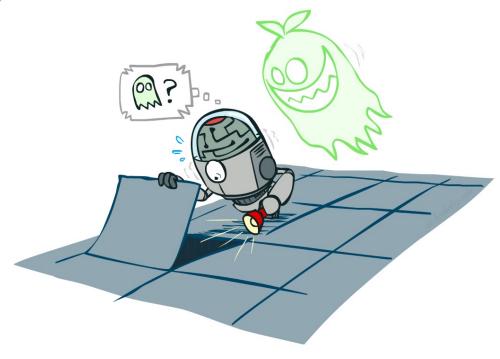
Hanna Hajishirzi Uncertainty, Bayes Nets

slides adapted from
Dan Klein, Pieter Abbeel ai.berkeley.edu
And Dan Weld, Luke Zettlemoyer



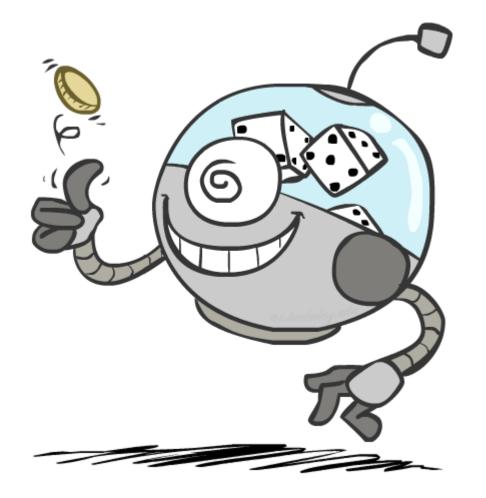
#### Our Status in CSE573

- We're done with Search and planning
- We are done with learning to make decisions
- Probabilistic Reasoning and Machine Learning
  - Diagnosis
  - Speech recognition
  - Tracking objects
  - Robot mapping
  - Genetics
  - Error correcting codes
  - ... lots more!



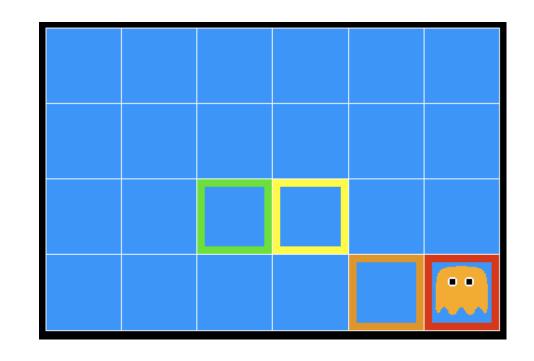
#### Outline

- Probability
- Bayes Nets
- You'll need all this stuff for the next few weeks, so make sure you go over it now!



### Inference in Ghostbusters

- A ghost is in the grid somewhere
- Sensor readings tell how close a square is to the ghost
  - On the ghost: red
  - 1 or 2 away: orange
  - 3 or 4 away: yellow
  - 5+ away: green

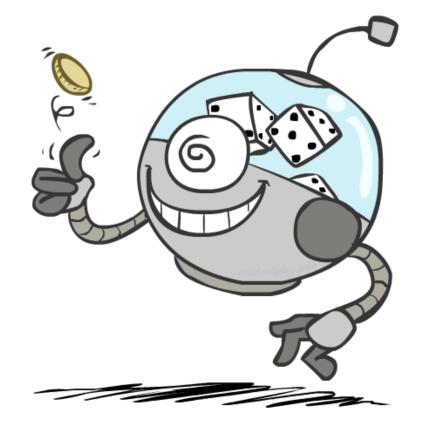


Sensors are noisy, but we know P(Color | Distance)

P(red   3)	P(orange   3)	P(yellow   3)	P(green   3)
0.05	0.15	0.5	0.3

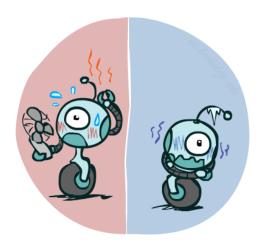
### Random Variables

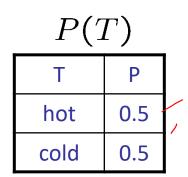
- A random variable is some aspect of the world about which we (may) have uncertainty
  - R = Is it raining?
  - T = Is it hot or cold?
  - D = How long will it take to drive to work?
  - L = Where is the ghost?
- We denote random variables with capital letters
- Random variables have domains
  - R in {true, false} (often write as {+r, -r})
  - T in {hot, cold}
  - D in  $[0, \infty)$
  - L in possible locations, maybe {(0,0), (0,1), ...}



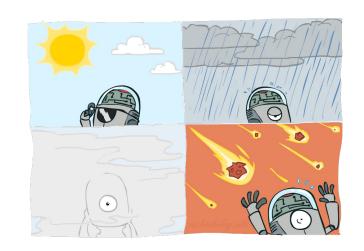
## **Probability Distributions**

- Associate a probability with each outcome
  - Temperature:





Weather:



P(W)

W	Р
sun	0.6
rain	0.1
fog	0.3
meteor	(0.0)

## **Probability Distributions**

Unobserved random variables have distributions

P(T)		
Т	Р	
hot	0.5	
cold	0.5	

D/D

1 (11)	
W	Р
sun	0.6
rain	0.1
fog	0.3
meteor	0.0

P(W)

- A distribution is a TABLE of probabilities of values
- A probability (lower case value) is a single number

$$P(W = rain) = 0.1$$

• Must have: 
$$\forall x \ P(X=x) \ge 0$$
 and  $\sum P(X=x) = 1$ 

$$P(hot) = P(T = hot),$$
  
 $P(cold) = P(T = cold),$   
 $P(rain) = P(W = rain),$   
...

OK if all domain entries are unique

#### Joint Distributions

• A *joint distribution* over a set of random variables:  $X_1, X_2, ... X_n$  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: 
$$P(x_1, x_2, \dots x_n) \ge 0$$

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

#### P(T,W)

Т	W	Р
hot	sun	0.4
hot	rain	0.1
cold	sun	0.2
cold	rain	0.3

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

#### **Events**

An event is a set E of outcomes

$$P(E) = \sum_{(x_1...x_n)\in E} P(x_1...x_n)$$

- From a joint distribution, we can calculate the probability of any event
  - Probability that it's hot AND sunny?
  - Probability that it's hot?
  - Probability that it's hot OR sunny?
- Typically, the events we care about are partial assignments, like P(T=hot)

P(T,W)

Т	W	Р
hot	sun	0.4
hot	rain	0.1
cold	sun	0.2
cold	rain	0.3

## Marginal Distributions

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

$D_{I}$	T	7	$\overline{W}$	1
1	(I)	,	VV	"

Т	W	Р
hot	sun	0.4
hot	rain	0.1
cold	sun	0.2
cold	rain	0.3

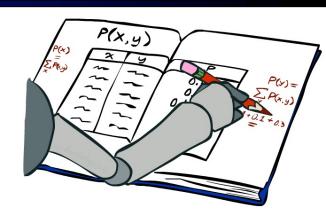
$$P(t) = \sum_{s} P(t, s)$$

$$P(s) = \sum_{t} P(t, s)$$

Т	Р
hot	0.5
cold	0.5



W	Р
sun	0.6
rain	0.4

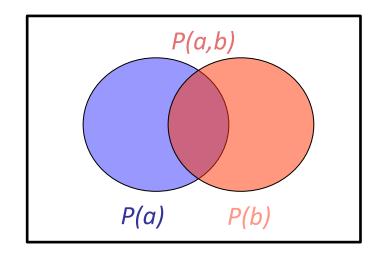


#### **Conditional Probabilities**

- A simple relation between joint and conditional probabilities
  - In fact, this is taken as the *definition* of a conditional probability

$$P(a|b) = \frac{P(a,b)}{P(b)}$$

Т	W	Р
hot	sun	0.4
hot	rain	0.1
cold	sun	0.2
cold	rain	0.3



$$P(W = s | T = c) = \frac{P(W = s, T = c)}{P(T = c)} = \frac{0.2}{0.5} = 0.4$$

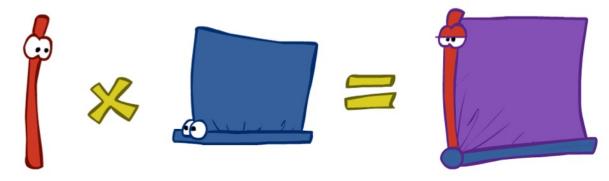
$$= P(W = s, T = c) + P(W = r, T = c)$$

$$= 0.2 + 0.3 = 0.5$$

#### The Product Rule

Sometimes have conditional distributions but want the joint

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



#### The Product Rule

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

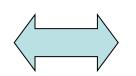
#### Example:

P(W)

R	Р
sun	0.8
rain	0.2

P(D|W)

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



P(D,W)

D	W	Р
wet	sun	
dry	sun	
wet	rain	
dry	rain	-

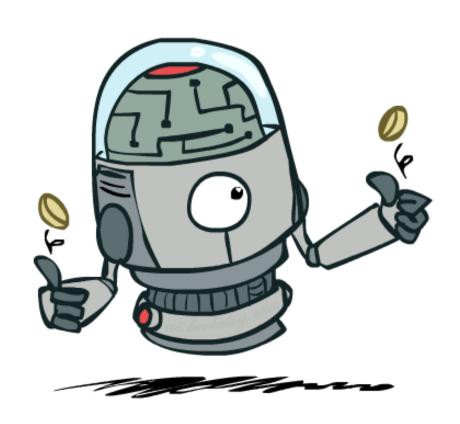
#### Probabilistic Models

- 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)

# Independence



## Independence

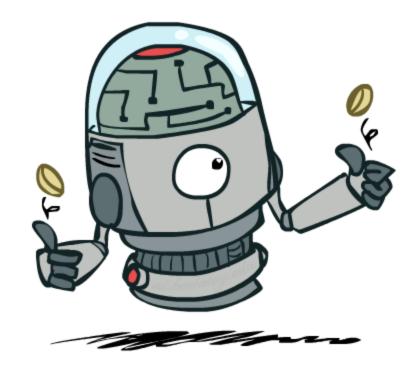
Two variables are independent if:

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

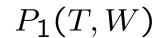
- This says that their joint distribution factors into a product two simpler distributions
- Another form:

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

- lacktriangle We write:  $X \coprod Y$
- 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?



Т	W	Р
hot	sun	0.4
hot	rain	0.1
cold	sun	0.2
cold	rain	0.3

#### P(T)

Т	Р
hot	0.5
cold	0.5

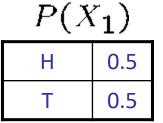
$$P_2(T,W)$$

Т	W	Р
hot	sun	0.3
hot	rain	0.2
cold	sun	0.3
cold	rain	0.2

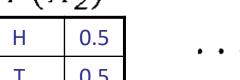
P(W)	
W	Р
sun	0.6
rain	0.4

# Example: Independence

N fair, independent coin flips:

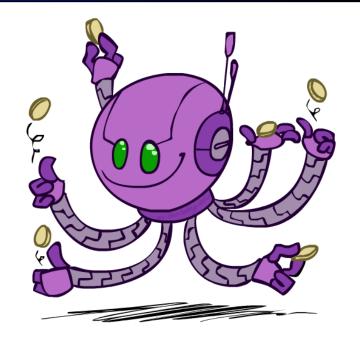


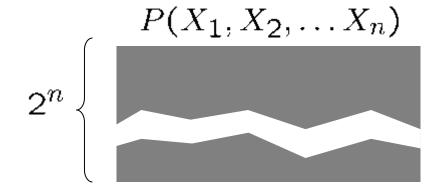
$P(X_2)$	
Н	0.5
Т	0.5

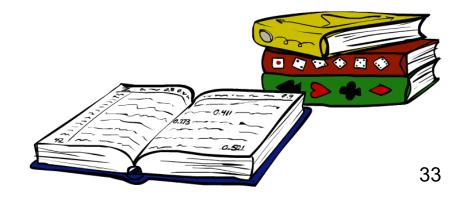


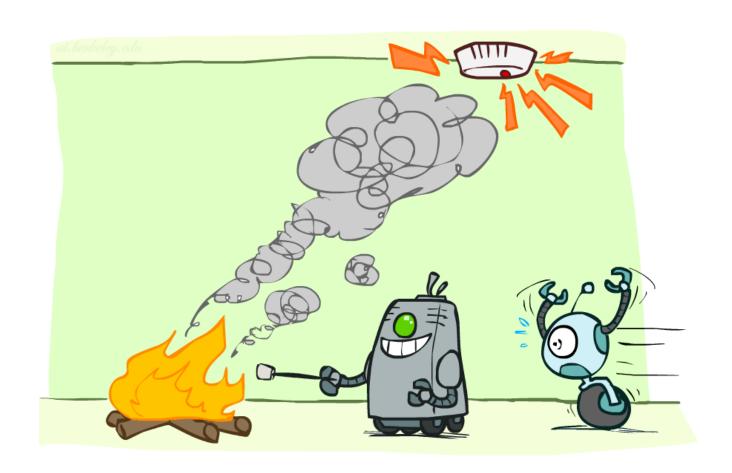
$-P(\Lambda_n)$	
Н	0.5
Т	0.5

D(V)

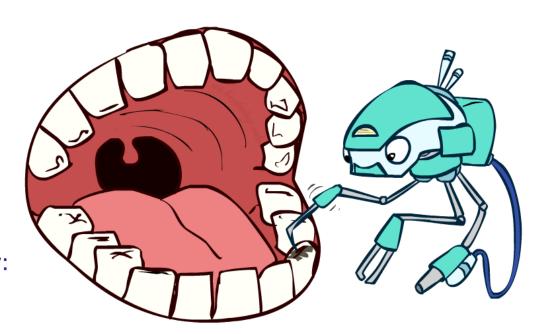








- P(Toothache, Cavity, Catch)
- If I have a cavity, the probability that the probe catches in it doesn't depend on whether I have a toothache:
  - P(+catch | +toothache, +cavity) = P(+catch | +cavity)
- The same independence holds if I don't have a cavity:
  - P(+catch | +toothache, -cavity) = P(+catch | -cavity)
- Catch is conditionally independent of Toothache given Cavity:
  - P(Catch | Toothache, Cavity) = P(Catch | Cavity)
- Equivalent statements:
  - P(Toothache | Catch , Cavity) = P(Toothache | Cavity)
  - P(Toothache, Catch | Cavity) = P(Toothache | Cavity) P(Catch | Cavity)
  - One can be derived from the other easily



- 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

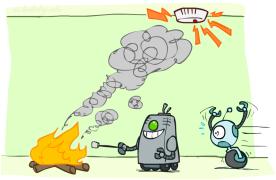
$$\forall x, y, z : P(x|z, y) = P(x|z)$$

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



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

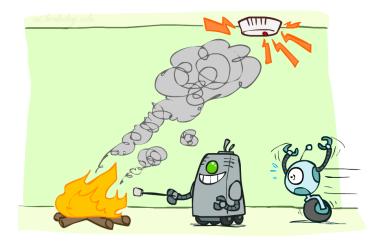


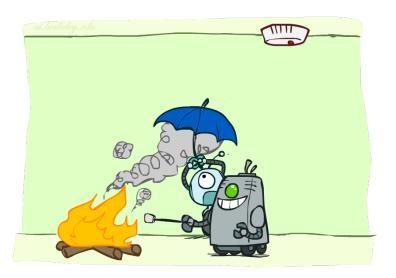


- What about this domain:
  - Traffic
  - Umbrella
  - Raining



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





## Conditional Independence and the Chain Rule

- Chain rule:  $P(X_1, X_2, ... X_n) = P(X_1)P(X_2|X_1)P(X_3|X_1, X_2)...$
- Trivial decomposition:

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

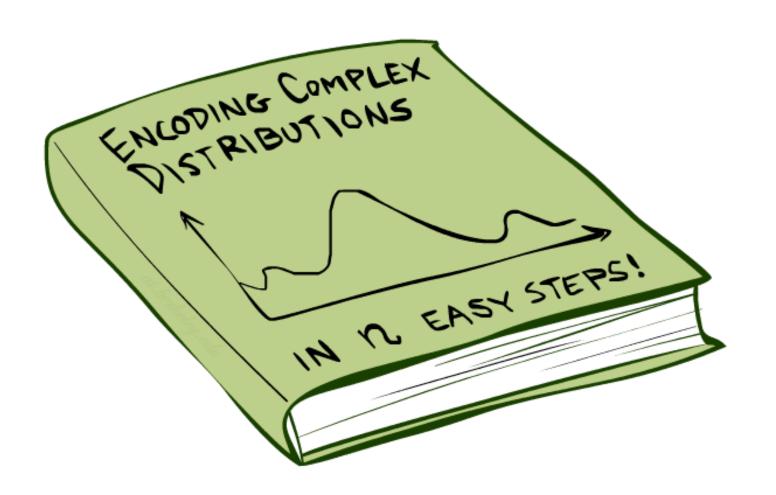
With assumption of conditional independence:

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



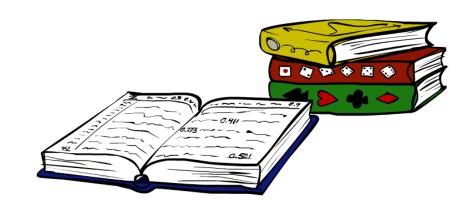
- We can represent joint distributions by multiplying these simpler local distributions.
- Bayes'nets / graphical models help us express conditional independence assumptions 41

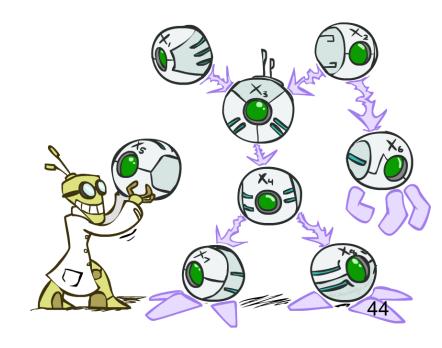
## Bayes'Nets: Big Picture



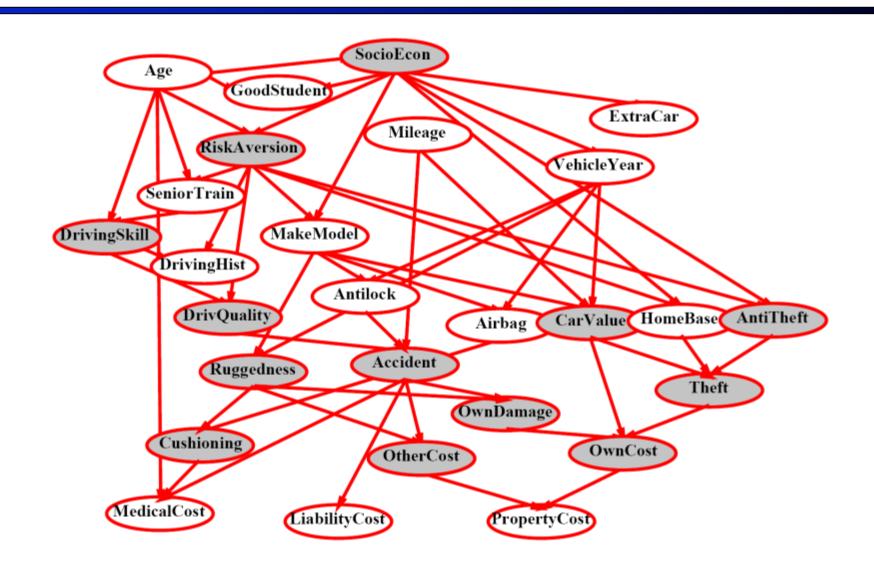
## Bayes' Nets: Big Picture

- 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





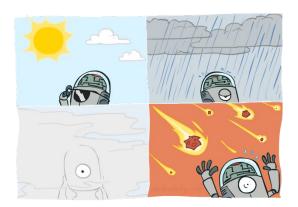
## Example Bayes' Net: Insurance



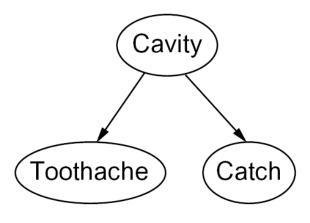
## **Graphical Model Notation**

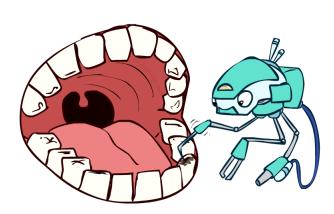
- Nodes: variables (with domains)
  - Can be assigned (observed) or unassigned (unobserved)





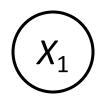
- Arcs: interactions
  - Indicate "direct influence" between variables
  - Formally: encode conditional independence (more later)
- For now: imagine that arrows mean direct causation (in general, they don't!)





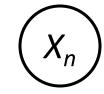
## Example: Coin Flips

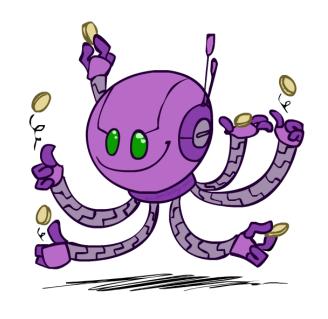
N independent coin flips





• •





No interactions between variables: absolute independence

## Example: Traffic

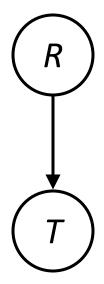
- Variables:
  - R: It rains
  - T: There is traffic
- Model 1: independence







Model 2: rain causes traffic



Why is an agent using model 2 better?

# Example: Traffic II

#### Variables

T: Traffic

R: It rains

L: Low pressure

■ D: Roof drips

B: Ballgame

C: Cavity



## Example: Alarm Network

#### Variables

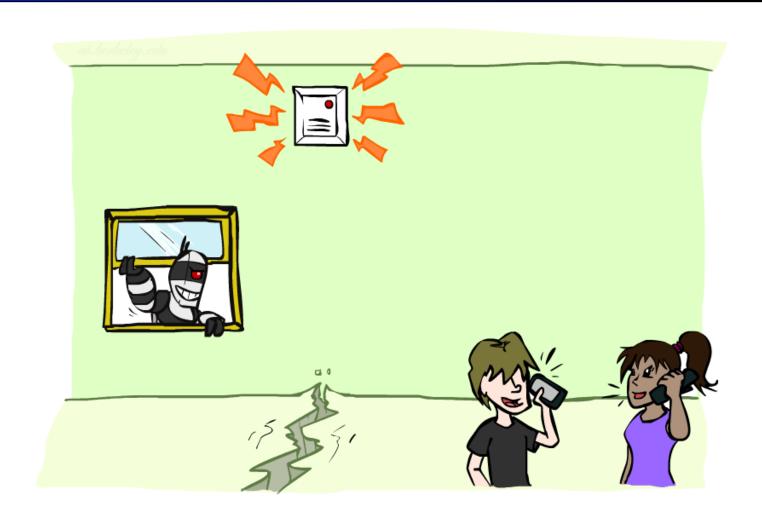
■ B: Burglary

A: Alarm goes off

M: Mary calls

■ J: John calls

■ E: Earthquake!



## Example: Alarm Network

#### Variables

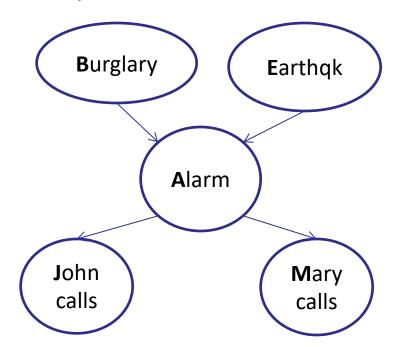
■ B: Burglary

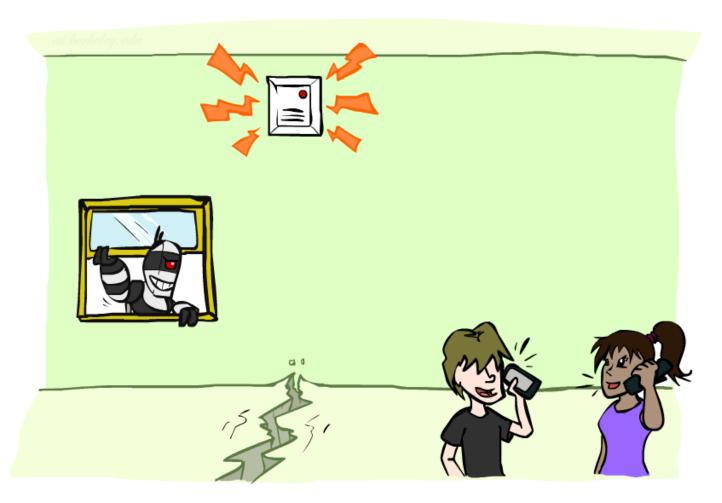
A: Alarm goes off

M: Mary calls

■ J: John calls

■ E: Earthquake!





# Bayes' Net Semantics



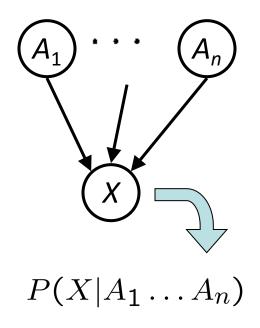
## Bayes' Net Semantics



- 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

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

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



A Bayes net = Topology (graph) + Local Conditional Probabilities<sub>55</sub>

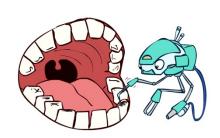
### Probabilities in BNs

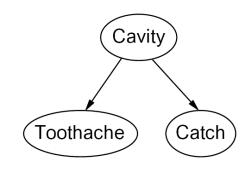


- Bayes' nets implicitly encode joint distributions
  - As a product of local conditional distributions
  - To see what probability a BN gives to a full assignment, multiply all the relevant conditionals together:

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

Example:





P(+cavity, +catch, -toothache)

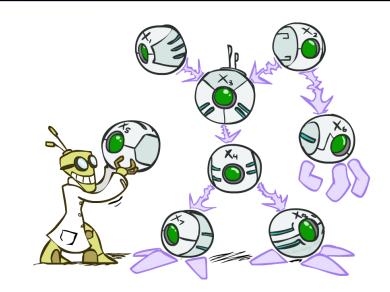
# Bayes' Net Representation

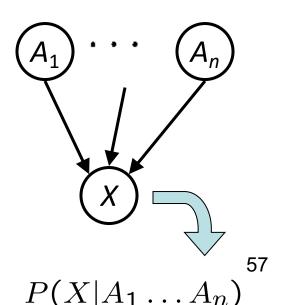
- A directed, acyclic graph, one node per random variable
- A conditional probability table (CPT) for each node
  - A collection of distributions over X, one for each combination of parents' values

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

- 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:

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





### Probabilities in BNs



Why are we guaranteed that setting

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

results in a proper joint distribution?

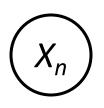
- Chain rule (valid for all distributions):  $P(x_1, x_2, \dots x_n) = \prod_{i=1}^n P(x_i | x_1 \dots x_{i-1})$
- Assume conditional independences:  $P(x_i|x_1,...x_{i-1}) = P(x_i|parents(X_i))$ 
  - $\rightarrow$  Consequence:  $P(x_1, x_2, \dots x_n) = \prod_{i=1}^n P(x_i | parents(X_i))$
- Not every BN can represent every joint distribution
  - The topology enforces certain conditional independencies

## Example: Coin Flips

$$X_1$$







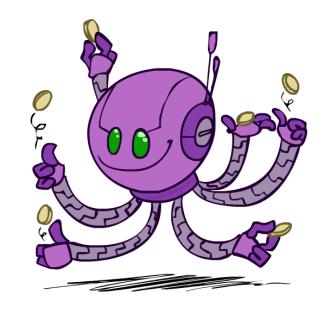
$$P(X_1)$$

h	0.5
t	0.5

P	(	X	2	)

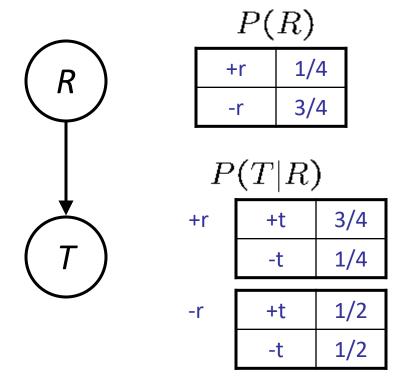
h	0.5
t	0.5

$$P(X_n)$$
h 0.5
t 0.5



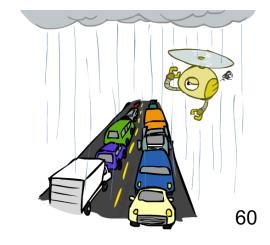
$$P(h, h, t, h) = P(h)P(h)P(t)P(h)$$

## Example: Traffic

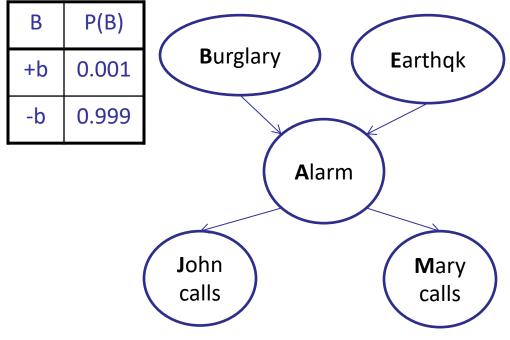


$$P(+r, -t) = P(+r)P(-t|+r) = \frac{1}{4} \cdot \frac{1}{4}$$





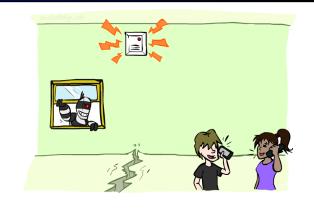
## Example: Alarm Network



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

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

ш	P(E)	
+e	0.002	
<del>-</del> e	0.998	

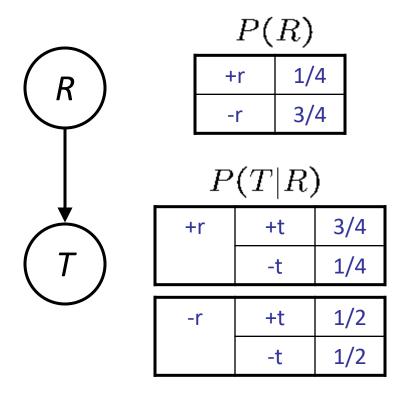


В	E	Α	P(A B,E)
+b	+e	+a	0.95
+b	+e	-a	0.05
+b	-е	+a	0.94
+b	-е	-a	0.06
-b	+e	+a	0.29
-b	+e	-a	0.71
-b	-е	+a	0.001
-b	-е	-a	0.999

P(M|A)P(J|A)P(A|B,E)

# Example: Traffic

#### Causal direction





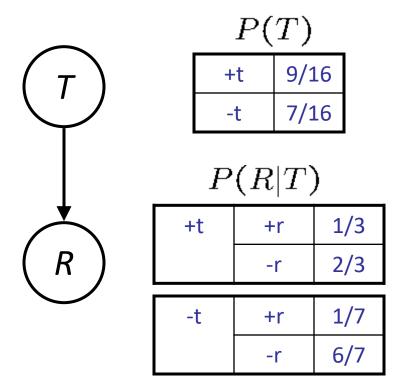


P	T	٦	Į	?)
•	(Τ	7	1	v

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

## Example: Reverse Traffic

Reverse causality?





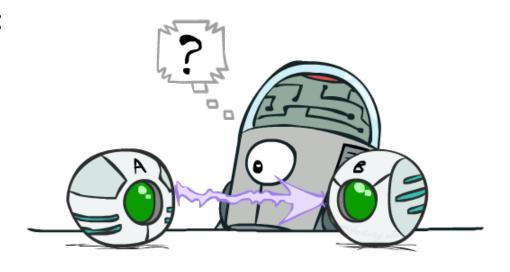
P(T,R)

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

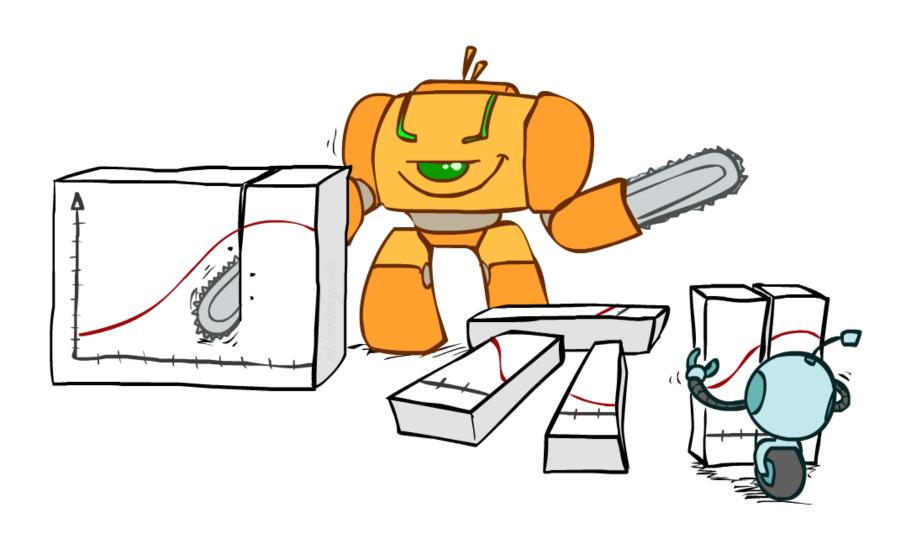
## Causality?

- 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

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



# Bayes Rule



## Bayes' Rule

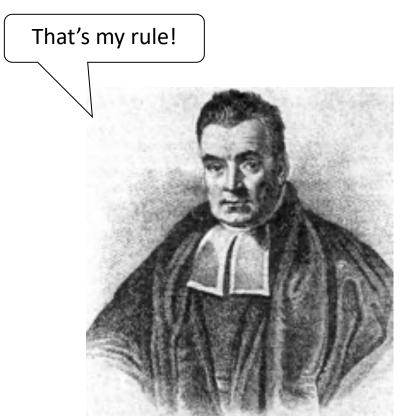
Two ways to factor a joint distribution over two variables:

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

Dividing, we get:

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

- Why is this at all helpful?
  - Lets us build one conditional from its reverse
  - Often one conditional is tricky but the other one is simple
  - Foundation of many systems we'll see later (e.g. ASR, MT)



In the running for most important AI equation!

## Inference with Bayes' Rule

Example: Diagnostic probability from causal probability:

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

- Example:
  - M: meningitis, S: stiff neck

$$P(+m) = 0.0001$$
 
$$P(+s|+m) = 0.8$$
 Example givens 
$$P(+s|-m) = 0.01$$

$$P(+m|+s) = \frac{P(+s|+m)P(+m)}{P(+s)} = \frac{P(+s|+m)P(+m)}{P(+s|+m)P(+m) + P(+s|-m)P(-m)} = \frac{0.8 \times 0.0001}{0.8 \times 0.0001 + 0.01 \times 0.999}$$

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

## Quiz: Bayes' Rule

Given:

#### P(W)

R	Р
sun	0.8
rain	0.2

#### P(D|W)

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

What is P(W | dry)?

## Quiz: Bayes' Rule

#### Given:

#### P(W)

R	Р
sun	0.8
rain	0.2

#### P(D|W)

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

### What is P(W | dry)?

 $P(sun|dry) \sim P(dry|sun)P(sun) = .9*.8 = .72$  $P(rain|dry) \sim P(dry|rain)P(rain) = .3*.2 = .06$ 

P(sun|dry)=12/13

P(rain|dry)=1/13

### **Uncertainty Summary**

• Conditional probability 
$$P(x|y) = \frac{P(x,y)}{P(y)}$$

Product rule

- P(x,y) = P(x|y)P(y)
- Chain rule  $P(X_1, X_2, ... X_n) = P(X_1)P(X_2|X_1)P(X_3|X_1, X_2)...$ =  $\prod_{i=1}^n P(X_i|X_1, ..., X_{i-1})$
- **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:  $X \perp\!\!\!\perp Y \mid Z$  $\forall x, y, z : P(x, y \mid z) = P(x \mid z) P(y \mid z)$

**BN** lecture

# 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

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

- 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:

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

