

# CSE 473: Artificial Intelligence

## Probability Review... HMMs



Daniel Weld

University of Washington

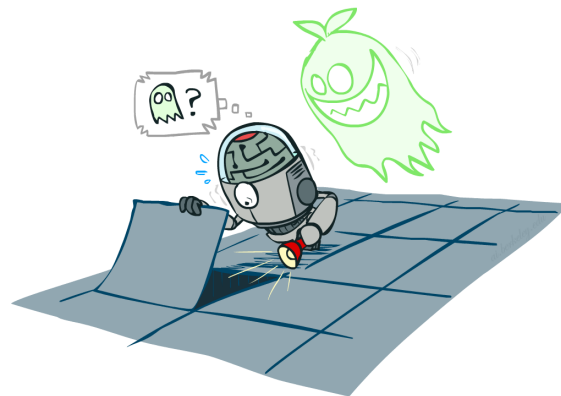
[These slides were created by Dan Klein and Pieter Abbeel for CS188 Intro to AI at UC Berkeley. All CS188 materials are available at <http://ai.berkeley.edu>.]

## Topics from 30,000'

- We're done with Part I Search and Planning!

- Part II: Probabilistic Reasoning

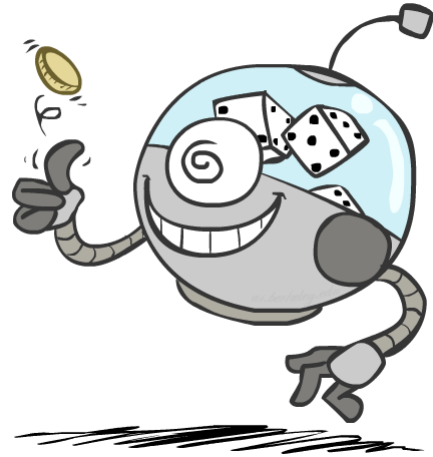
- Diagnosis
- Speech recognition
- Tracking objects
- Robot mapping
- Genetics
- Error correcting codes
- ... lots more!



- Part III: Machine Learning

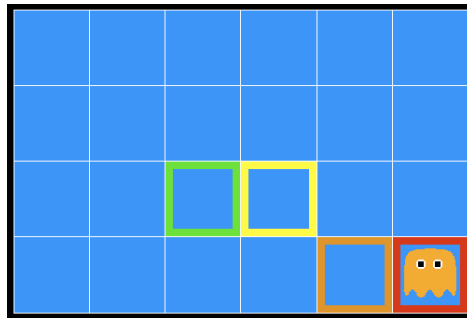
# Outline

- Probability
  - Random Variables
  - Joint and Marginal Distributions
  - Conditional Distribution
  - Product Rule, Chain Rule, Bayes' Rule
  - Inference
  - Independence
- You'll need all this stuff A LOT 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  $P(\text{Color} \mid \text{Distance})$



$P(\text{red} \mid 3)$	$P(\text{orange} \mid 3)$	$P(\text{yellow} \mid 3)$	$P(\text{green} \mid 3)$
0.05	0.15	0.5	0.3

[Demo: Ghostbuster – no probability (L12D1)]

## Video of Demo Ghostbuster – No probability



## Uncertainty

- **General situation:**

- **Observed variables (evidence):** Agent knows certain things about the state of the world (e.g., sensor readings or symptoms)
- **Unobserved variables:** Agent needs to reason about other aspects (e.g. where an object is or what disease is present)
- **Model:** Agent knows something about how the known variables relate to the unknown variables

0.11	0.11	0.11
0.11	0.11	0.11
0.11	0.11	0.11

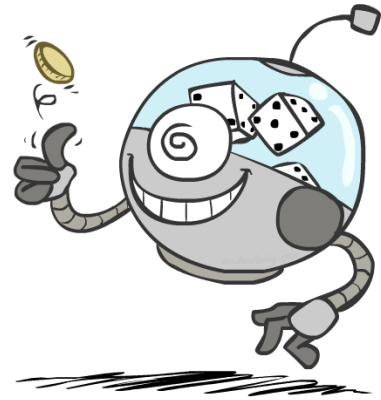
0.17	0.10	0.10
0.09	0.17	0.10
<-0.01	0.09	0.17

- Probabilistic reasoning gives us a framework for managing our beliefs and knowledge

<-0.01	<-0.01	0.03
<-0.01	0.05	0.05
<-0.01	0.05	0.81

# Random Variables

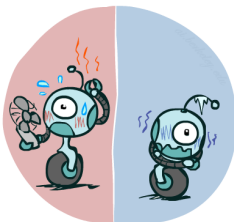
- A random variable is some aspect of the world about which we (may) have uncertainty
  - R = Is it raining?
  - T = Is it hot or cold?
  - D = How long will it take to drive to work?
  - L = Where is the ghost?
- We denote random variables with capital letters
- Like variables in a CSP, random variables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



# Probability Distributions

- Associate a probability with each value

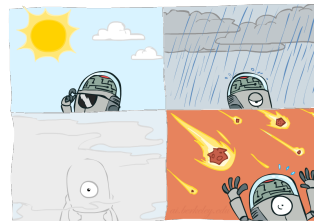
- Temperature:



$P(T)$

T	P
hot	0.5
cold	0.5

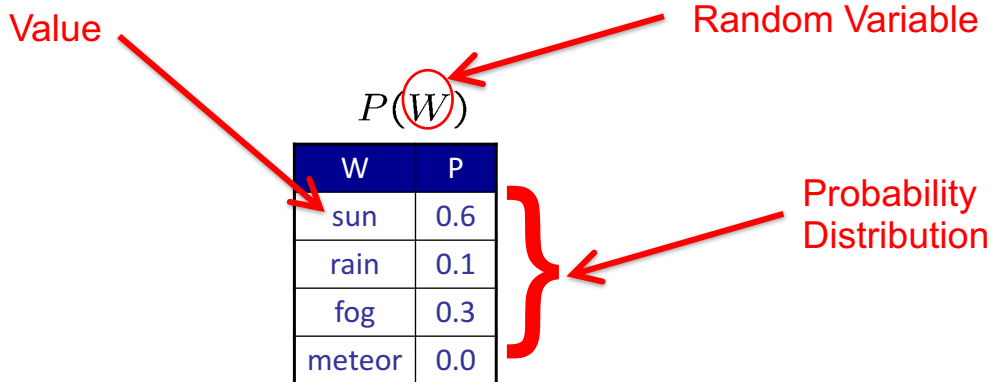
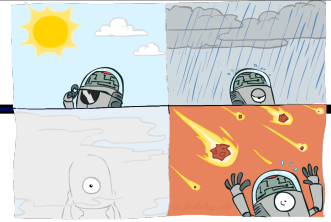
- Weather:



$P(W)$

W	P
sun	0.6
rain	0.1
fog	0.3
meteor	0.0

# What is....?



## Probability Distributions

- Unobserved random variables have distributions

$P(T)$

T	P
hot	0.5
cold	0.5

$P(W)$

W	P
sun	0.6
rain	0.1
fog	0.3
meteor	0.0

Shorthand notation:

$$P(\text{hot}) = P(T = \text{hot}),$$

$$P(\text{cold}) = P(T = \text{cold}),$$

$$P(\text{rain}) = P(W = \text{rain}),$$

...

OK if all domain entries are unique

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

$$P(W = \text{rain}) = 0.1$$

- Must have:  $\forall x P(X = x) \geq 0$  and  $\sum_x P(X = x) = 1$

# Joint Distributions

- A *joint distribution* over a set of random variables:  $X_1, X_2, \dots, X_n$  specifies a probability 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:  $P(x_1, x_2, \dots, x_n) \geq 0$

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

$P(T, W)$

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

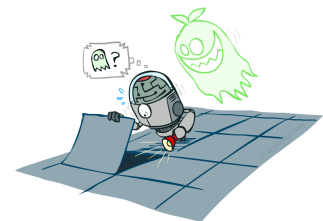
- Size of joint distribution if n variables with domain sizes d?
  - For all but the smallest distributions, impractical to write out!

# Probabilistic Models

- A *probabilistic model* is a joint distribution over a set of random variables

Distribution over T,W

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

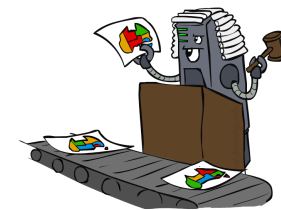


- Probabilistic models:

- (Random) variables with domains
- Joint distributions: say whether assignments (called "*outcomes*") are likely
- Normalized: sum to 1.0
- Ideally: only certain variables directly interact

Constraint over T,W

T	W	P
hot	sun	T
hot	rain	F
cold	sun	F
cold	rain	T



- Constraint satisfaction problems:

- Variables with domains
- Constraints: state whether assignments are possible
- Ideally: only certain variables directly interact

## Events

- An *event* is a set  $E$  of outcomes

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

- 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  $P(T=\text{hot})$

$P(T, W)$

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

## Quiz: Events

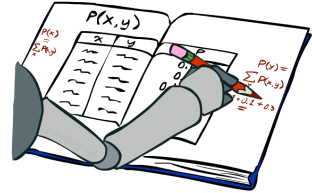
- $P(+x, +y)$  ?
- $P(+x)$  ?
- $P(-y \text{ OR } +x)$  ?

$P(X, Y)$

X	Y	P
+x	+y	0.2
+x	-y	0.3
-x	+y	0.4
-x	-y	0.1

# 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

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

$$P(T)$$

T	P
hot	0.5
cold	0.5

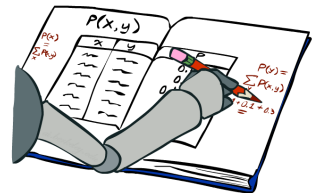
$$P(s) = \sum_t P(t, s)$$

$$P(W)$$

W	P
sun	0.6
rain	0.4

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

# Quiz: Marginal Distributions



$$P(X, Y)$$

X	Y	P
+x	+y	0.2
+x	-y	0.3
-x	+y	0.4
-x	-y	0.1

$$P(x) = \sum_y P(x, y)$$

$$P(X)$$

X	P
+x	
-x	

$$P(y) = \sum_x P(x, y)$$

$$P(Y)$$

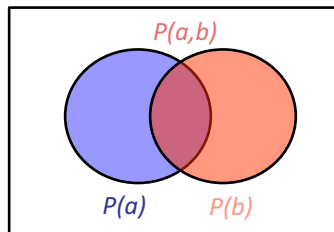
Y	P
+y	
-y	



## Conditional Probabilities

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

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



$P(T, W)$

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

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

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

$$= 0.2 + 0.3 = 0.5$$

## Quiz: Conditional Probabilities

$P(X, Y)$

X	Y	P
+x	+y	0.2
+x	-y	0.3
-x	+y	0.4
-x	-y	0.1

- $P(+x | +y) ?$

- $P(-x | +y) ?$

- $P(-y | +x) ?$

# Conditional Distributions

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

Conditional Distributions

$P(W|T)$

W	P
sun	0.8
rain	0.2

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

## Conditional Distributions - The Slow Way...

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

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

$$= \frac{P(W = s, T = c)}{P(W = s, T = c) + P(W = r, T = c)}$$

$$= \frac{0.2}{0.2 + 0.3} = 0.4$$

W	P
sun	0.4
rain	0.6

$\longrightarrow$

$$P(W = r|T = c) = \frac{P(W = r, T = c)}{P(T = c)}$$

$$= \frac{P(W = r, T = c)}{P(W = s, T = c) + P(W = r, T = c)}$$

$$= \frac{0.3}{0.2 + 0.3} = 0.6$$

## Normalization Trick

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

$P(T, W)$

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

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

## 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

- Why does this work? Sum of selection is P(evidence)! (P(T=c), here)

$$P(x_1|x_2) = \frac{P(x_1, x_2)}{P(x_2)} = \frac{P(x_1, x_2)}{\sum_{x_1} P(x_1, x_2)}$$

## Quiz: Normalization Trick

- $P(X | Y=-y)$  ?

$P(X, Y)$

X	Y	P
+x	+y	0.2
+x	-y	0.3
-x	+y	0.4
-x	-y	0.1

**SELECT** the joint probabilities matching the evidence



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



## To Normalize

- Dictionary: "To bring or restore to a normal condition"

All entries sum to ONE

- Procedure:

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

- Example 1

W	P
sun	0.2
rain	0.3

Normalize  
Z = 0.5

W	P
sun	0.4
rain	0.6

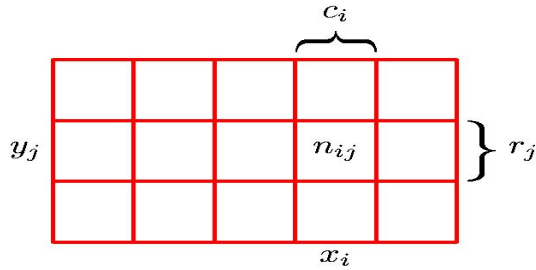
- Example 2

T	W	P
hot	sun	20
hot	rain	5
cold	sun	10
cold	rain	15

Normalize  
Z = 50

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

# Terminology



## Marginal Probability

$$p(X = x_i) = \frac{c_i}{N}$$

## Joint Probability

$$p(X = x_i, Y = y_j) = \frac{n_{ij}}{N}$$

## Conditional Probability

$$p(Y = y_j | X = x_i) = \frac{n_{ij}}{c_i}$$

↑  
X value is given




# Probabilistic Inference

- Probabilistic inference =  
*“compute a desired probability from other known probabilities (e.g. conditional from joint)”*
- We generally compute conditional probabilities
  - $P(\text{on time} \mid \text{no reported accidents}) = 0.90$
  - These represent the agent’s *beliefs* given the evidence
- 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*



# Probabilistic Inference in Ghostbusters

- A ghost is in the grid somewhere
- Noisy Sensor readings tell approx how close a square is to the ghost
  - 1 or 2 away: orange
  - Etc.

.05	.05	.05	.05	.05
.05	.05	.05	.05	.05
.05	.05	.05	.05	.05
.05	.05	.05	.05	

- Sensors are noisy, but we know  $P(\text{Color} \mid \text{Distance})$

$P(\text{red} \mid 3)$	$P(\text{orange} \mid 3)$	$P(\text{yellow} \mid 3)$	$P(\text{green} \mid 3)$
0.05	0.15	0.5	0.3

# Probabilistic Inference in Ghostbusters

- A ghost is in the grid somewhere
- Noisy Sensor readings tell approx how close a square is to the ghost
  - 1 or 2 away: orange
  - Etc.



*How update the probabilities?*

# Inference by Enumeration

- General case:

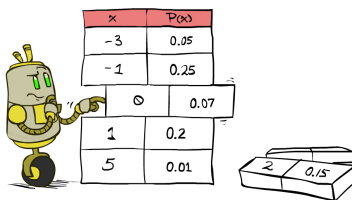
- Evidence variables:  $E_1 \dots E_k = e_1 \dots e_k$
  - Query\* variable:  $Q$
  - Hidden variables:  $H_1 \dots H_r$
- $\left. \begin{array}{l} X_1, X_2, \dots, X_n \\ \text{All variables} \end{array} \right\}$

- We want:

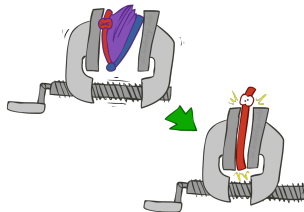
$$P(Q|e_1 \dots e_k)$$

*\* 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

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

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

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

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

## Inference by Enumeration

- $P(W=\text{sun})?$

- $P(W=\text{sun} \mid S=\text{winter})?$

- $P(W=\text{sun} \mid S=\text{winter}, T=\text{hot})?$

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

- Computational problems?
  - Worst-case time complexity  $O(d^n)$
  - Space complexity  $O(d^n)$  to store the joint distribution



## Don't be Fooled

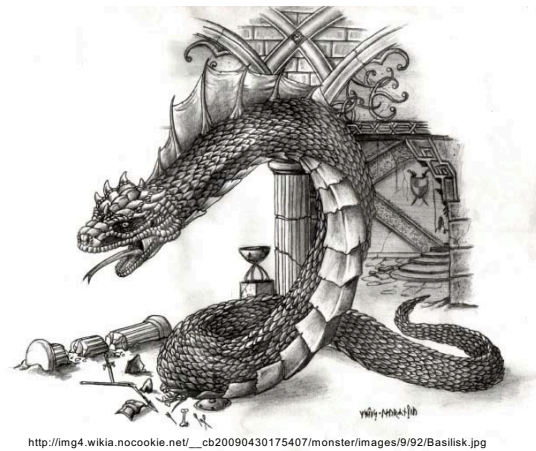
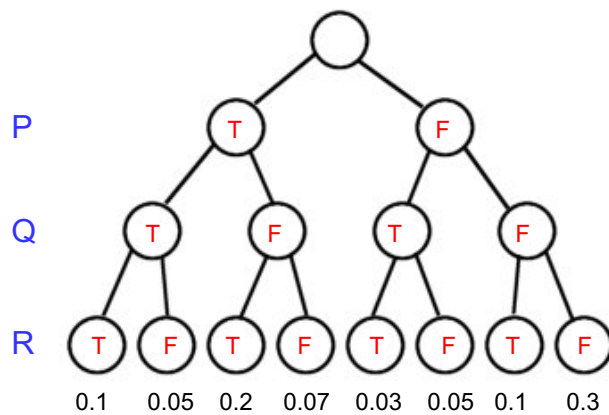
- It may look cute...



34

## Don't be Fooled

- It gets big...



35

# The Sword of Conditional Independence!



Slay  
the  
Basilisk!

I am a BIG joint  
distribution!



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

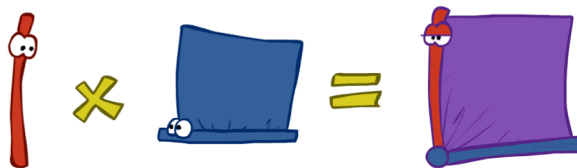
Or, equivalently:  $\forall x, y, z : P(x | z, y) = P(x | z)$

36

# The Product Rule

- Sometimes have conditional distributions but want the joint

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



## The Product Rule

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

- Example:

$P(W)$		$P(D W)$				$P(D, W)$		
R	P	D	W	P		D	W	P
sun	0.8	wet	sun	0.1	↔	wet	sun	
rain	0.2	dry	sun	0.9		dry	sun	
		wet	rain	0.7		wet	rain	
		dry	rain	0.3		dry	rain	

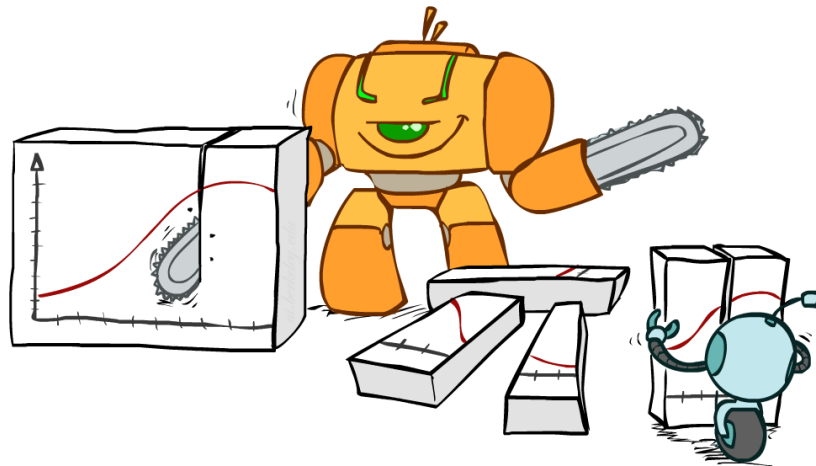
## The Chain Rule

- More generally, can always write any joint distribution as an incremental product of conditional distributions

$$P(x_1, x_2, x_3) = P(x_1)P(x_2|x_1)P(x_3|x_1, x_2)$$

$$P(x_1, x_2, \dots, x_n) = \prod_i P(x_i|x_1 \dots x_{i-1})$$

# 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!

That's my rule!



## 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

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

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

- Note: posterior probability of meningitis still very small =0.0079
- Note: you should still get stiff necks checked out! Why?

## Quiz: Bayes' Rule

- Given:

$P(W)$	
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

- What is  $P(W=\text{rain} | \text{dry})$  ?

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

# 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.  $P(R = \text{yellow} | G=(1,1)) = 0.1$

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

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

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

[Demo: Ghostbuster – with probability (L12D2) ]

## Video of Demo Ghostbusters with Probability



# Independence

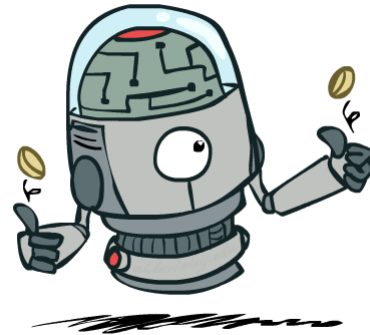
- Two variables are *independent* in a joint distribution if:

$$P(X, Y) = P(X)P(Y)$$

$$X \perp\!\!\!\perp Y$$

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

- Says the joint distribution *factors* into a product of two simple ones
- Usually variables aren't independent!
- Can use independence as a *modeling assumption*
  - Independence can be a simplifying assumption
  - Empirical* joint distributions: at best "close" to independent
  - What could we assume for {Weather, Traffic, Cavity}?
- Independence is like something from CSPs: what?



# Independence

$$P(A \wedge B) = P(A)P(B)$$

