

# CS 188: Artificial Intelligence

## Probability

Pieter Abbeel – UC Berkeley  
Many slides adapted from Dan Klein.

## Our Status in CS188

---

- We're done with Part I Search and Planning!
- Part II: Probabilistic Reasoning
  - Diagnosis
  - Tracking objects
  - Speech recognition
  - Robot mapping
  - Genetics
  - Error correcting codes
  - ... lots more!
- Part III: Machine Learning

2

## Part II: Probabilistic Reasoning

---

- Probability
- Distributions over LARGE Numbers of Random Variables
  - Representation
  - Independence
  - Inference
    - Variable Elimination
    - Sampling
  - Hidden Markov Models

3

## Probability

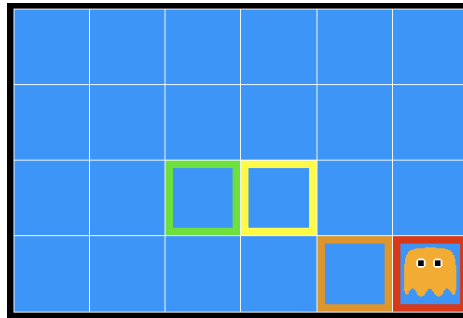
---

- Probability
  - Random Variables
  - Joint and Marginal Distributions
  - Conditional Distribution
  - Inference by Enumeration
  - Product Rule, Chain Rule, Bayes' Rule
  - Independence
- You' ll need all this stuff A LOT for the next few weeks, so make sure you go over it now and know it inside out! The next few weeks we will learn how to make these work computationally efficiently for LARGE numbers of random variables.

4

# 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  $P(\text{Color} | \text{Distance})$



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

# Uncertainty

- **General situation:**
  - **Evidence:** Agent knows certain things about the state of the world (e.g., sensor readings or symptoms)
  - **Hidden 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
- Probabilistic reasoning gives us a framework for managing our beliefs and knowledge

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

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

6

## Random Variables

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

7

## Probability Distributions

- Unobserved random variables have distributions

T	P
warm	0.5
cold	0.5

W	P
sun	0.6
rain	0.1
fog	0.3
meteor	0.0

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

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

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

8

# 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*):

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

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

$$P(T, W)$$

T	W	P
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?
- Must obey:

$$P(x_1, x_2, \dots, x_n) \geq 0$$

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

- For all but the smallest distributions, impractical to write out

9

# 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
  - Assignments are called *outcomes*
  - Joint distributions: say whether assignments (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 probs:
  - Variables with domains
  - Constraints: state whether assignments are possible
  - Ideally: only certain variables directly interact

10

# Events

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

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

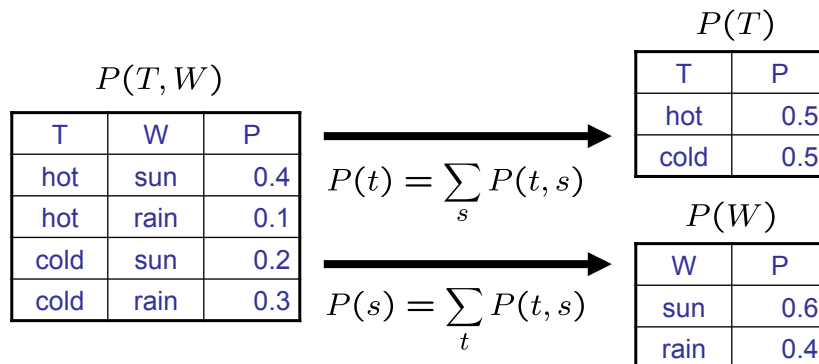
T	W	P
hot	sun	0.4
hot	rain	0.1
cold	sun	0.2
cold	rain	0.3

- 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  $P(T=hot)$

11

# Marginal Distributions

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

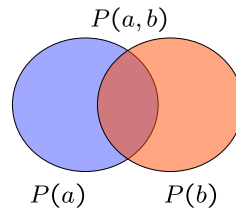


$$P(X_1 = x_1) = \sum_{x_2} P(X_1 = x_1, X_2 = x_2) \quad 12$$

# Conditional Probabilities

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

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



$P(T, W)$

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

$$P(W = r|T = c) = ???$$

13

# Conditional Distributions

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

Conditional Distributions

Joint Distribution

$P(W|T)$

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

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

$P(T, W)$

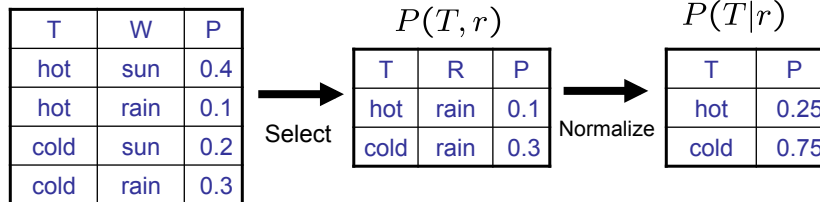
T	W	P
hot	sun	0.4
hot	rain	0.1
cold	sun	0.2
cold	rain	0.3

14

## Normalization Trick

- A trick to get a whole conditional distribution at once:
  - Select the joint probabilities matching the evidence
  - Normalize the selection (make it sum to one)

$P(T, W)$



- Why does this work? Sum of selection is  $P(\text{evidence})!$  ( $P(r)$ , here)

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

15

## 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*

16



## Inference by Enumeration

---

- P(sun)?
- P(sun | winter)?
- P(sun | winter, warm)?

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

17

## 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:  $P(Q|e_1 \dots e_k)$
- First, select the entries consistent with the evidence
- Second, sum out H to get joint of Query and evidence:

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

- Finally, normalize the remaining entries to conditionalize

- **Obvious problems:**
  - Worst-case time complexity  $O(d^n)$
  - Space complexity  $O(d^n)$  to store the joint distribution

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

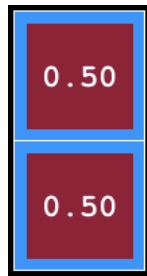
## Inference by Enumeration Example 2: Model for Ghostbusters

- Reminder: ghost is hidden, sensors are noisy

- T: Top sensor is red  
B: Bottom sensor is red  
G: Ghost is in the top

- Queries:  
 $P(+g) = ??$   
 $P(+g \mid +t) = ??$   
 $P(+g \mid +t, -b) = ??$

- Problem: joint distribution too large / complex



Joint Distribution

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

## The Product Rule

- Sometimes have conditional distributions but want the joint

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

- Example:

		$P(D W)$			$P(D, W)$		
$P(W)$		D	W	P	D	W	P
R	P	wet	sun	0.1	wet	sun	0.08
sun	0.8	dry	sun	0.9	dry	sun	0.72
rain	0.2	wet	rain	0.7	wet	rain	0.14
		dry	rain	0.3	dry	rain	0.06

## The Chain Rule

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

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

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

- Why is this always true?
- Can now build a joint distributions only specifying conditionals!
  - Bayesian networks essentially apply the chain rule plus make conditional independence assumptions.

21

## Bayes' Rule

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

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

That's my rule!

- Dividing, we get:

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



- 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!

22

## 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 is meningitis, s is stiff neck
 

$P(s m) = 0.8$	}	Example givens
$P(m) = 0.0001$		
$P(s) = 0.1$		

$$P(m|s) = \frac{P(s|m)P(m)}{P(s)} = \frac{0.8 \times 0.0001}{0.1} = 0.0008$$

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

23

## 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 = yellow | G=(1,1)) = 0.1

0.11	0.11	0.11
0.11	0.11	0.11
0.11	0.11	0.11

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

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

0.17	0.10	0.10
0.09	0.17	0.10
<0.01	0.09	0.17

24

# Independence

- Two variables are *independent* if:

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

- Says their joint distribution *factors* into a product two simpler ones.
- Usually variables are not independent!
- Equivalent definition of independence:

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

- 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 is like something from CSPs, what?

25

## Example: Independence?

$P_1(T, W)$			$P(T)$		$P_2(T, W)$		
T	W	P	T	P	T	W	P
warm	sun	0.4	warm	0.5	warm	sun	0.3
warm	rain	0.1	cold	0.5	warm	rain	0.2
cold	sun	0.2			cold	sun	0.3
cold	rain	0.3			cold	rain	0.2
			$P(W)$				
			W	P			
			sun	0.6			
			rain	0.4			

26

## Example: Independence

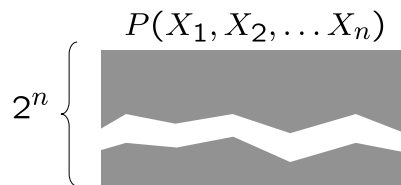
- N fair, independent coin flips:

H	0.5
T	0.5

H	0.5
T	0.5

...

H	0.5
T	0.5



27

## 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

28

## Conditional Independence

---

- Unconditional (absolute) independence very rare (why?)
- *Conditional independence* is our most basic and robust form of knowledge about uncertain environments:

$$\begin{aligned}\forall x, y, z : P(x, y|z) &= P(x|z)P(y|z) \\ \forall x, y, z : P(x|z, y) &= P(x|z)\end{aligned}\quad X \perp\!\!\!\perp Y|Z$$

- What about this domain:
  - Traffic
  - Umbrella
  - Raining
- What about fire, smoke, alarm?

29

## The Chain Rule Revisited

---

$$P(X_1, X_2, \dots, X_n) = P(X_1)P(X_2|X_1)P(X_3|X_1, X_2)\dots$$

- 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:
$$P(\text{Traffic, Rain, Umbrella}) = P(\text{Rain})P(\text{Traffic}|\text{Rain})P(\text{Umbrella}|\text{Rain})$$
- Representation size: 1 + 2 + 4 versus 1 + 2 + 2
- Bayes' nets / graphical models are concerned with distributions with conditional independences

30

# Ghostbusters Chain Rule

- Each sensor depends only on where the ghost is
- That means, the two sensors are conditionally independent, given the ghost position
- T: Top square is red  
B: Bottom square is red  
G: Ghost is in the top
- Givens:  
 $P(+g) = 0.5$   
 $P(+t \mid +g) = 0.8$   
 $P(+t \mid -g) = 0.4$   
 $P(+b \mid +g) = 0.4$   
 $P(+b \mid -g) = 0.8$

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

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