

Artificial Intelligence - INFOF311

Probability

Instructor : Tom Lenaerts

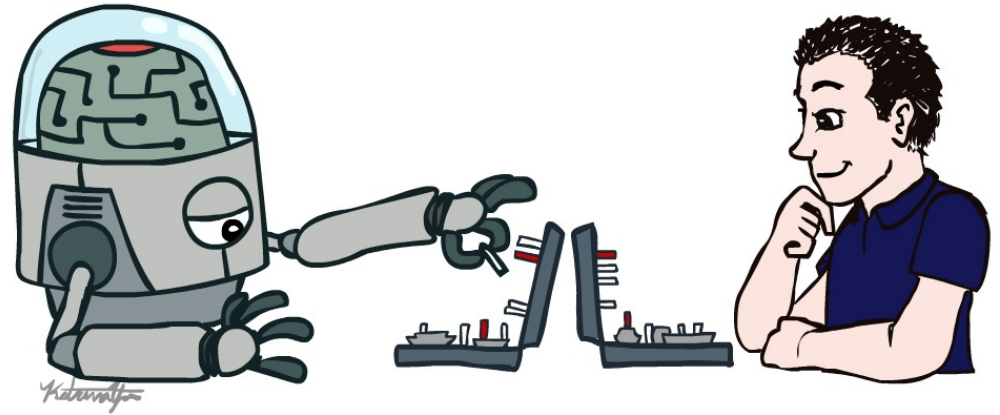


Acknowledgement

We thank Stuart Russell for his generosity in allowing us to use the slide set of the UC Berkeley Course CS188, Introduction to Artificial Intelligence. These slides were created by Dan Klein, Pieter Abbeel and Anca Dragan for CS188 Intro to AI at UC Berkeley. All CS188 materials are available at <http://ai.berkeley.edu>.



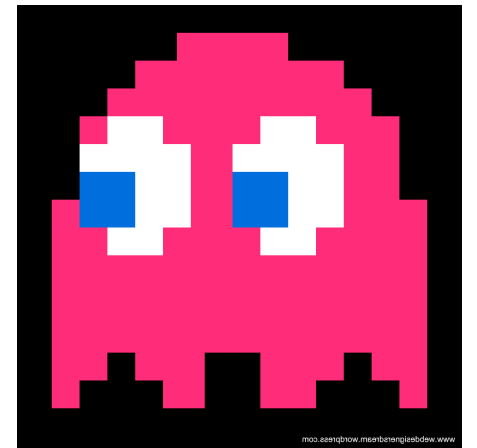
Center for
Human-Compatible
Artificial
Intelligence



The slides for INFOF311 are slightly modified versions of the slides of the spring and summer CS188 sessions in 2021 and 2022

Skipping logic agents

See course “*Informatique Fondamentale*”



Chapters 7-11, AI a modern approach

4 main themes :

Part 1 : Search and planning (uninformed and informed search, local search, game and adversarial search, ...)

Part 2: Probabilistic reasoning (Bayesian network, hidden Markov models,...)

Part 3: Decision making with uncertainty (MDP, reinforcement learning, ...)

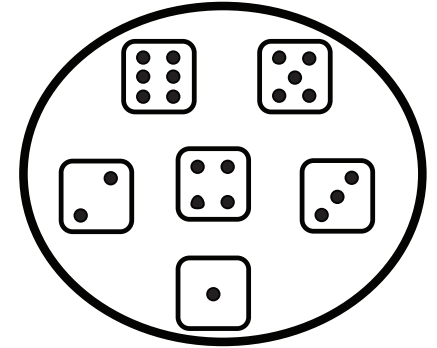
Part 4: Machine learning (naïve bayes, perceptrons, regression, neural networks, ...)

Uncertainty

- The real world is rife with uncertainty!
 - E.g., if I leave for the Brussels Airport 60 minutes before my flight, will I be there in time?
- Problems:
 - partial observability (road state, other drivers' plans, etc.)
 - noisy sensors (radio traffic reports, Google maps)
 - immense complexity of modelling and predicting traffic, security line, etc.
 - lack of knowledge of world dynamics (will tire burst? need COVID test?)
- Probabilistic assertions summarize effects of *ignorance* and *laziness*
- Combine probability theory + utility theory -> decision theory
 - **Maximize expected utility** : $a^* = \operatorname{argmax}_a \sum_s P(s \mid a) U(s)$

Basic laws of probability

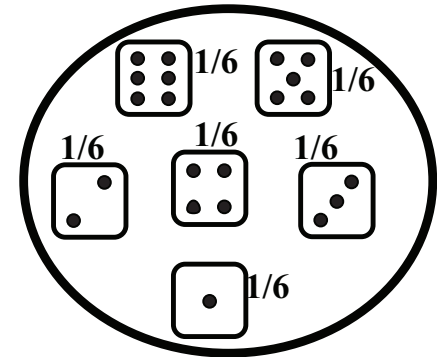
- Begin with a set Ω of possible worlds
 - E.g., 6 possible rolls of a die, $\{1, 2, 3, 4, 5, 6\}$



- A **probability model** assigns a number $P(\omega)$ to each world ω
 - E.g., $P(1) = P(2) = P(3) = P(4) = P(5) = P(6) = 1/6$.

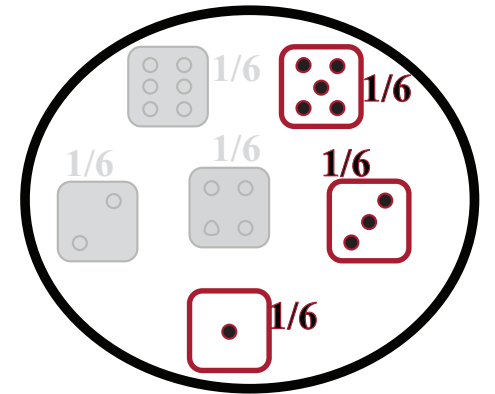
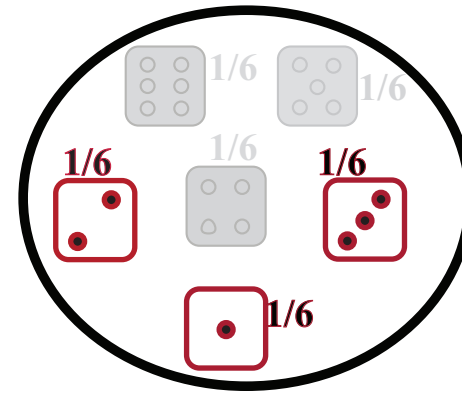
- These numbers must satisfy

- $0 \leq P(\omega) \leq 1$
- $\sum_{\omega \in \Omega} P(\omega) = 1$



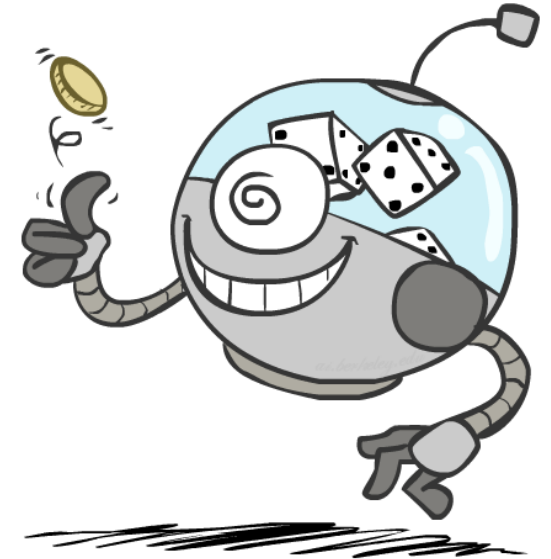
Basic laws contd.

- An **event** is any subset of Ω
 - E.g., “roll < 4” is the set {1,2,3}
 - E.g., “roll is odd” is the set {1,3,5}
- The probability of an event is the **sum** of probabilities over its worlds
 - $P(A) = \sum_{\omega \in A} P(\omega)$
 - E.g., $P(\text{roll} < 4) = P(1) + P(2) + P(3) = 1/2$



Random Variables

- A random variable (usually denoted by a capital letter) is some aspect of the world about which we (may) be uncertain
- Formally a **deterministic function** of ω
- The **range** of a random variable is the set of possible values
 - Odd = Is the dice roll an odd number? $\rightarrow \{true, false\}$
 - e.g. $Odd(1)=true$, $Odd(6) = false$
 - often write the event $Odd=true$ as odd , $Odd=false$ as $\neg odd$
 - T = Is it hot or cold? $\rightarrow \{hot, cold\}$
 - D = How long will it take to get to the airport? $\rightarrow [0, \infty)$
 - L_{Ghost} = Where is the ghost? $\rightarrow \{(0,0), (0,1), \dots\}$
- The **probability distribution** of a random variable X gives the probability for each value x in its range (probability of the event $X=x$)
 - $P(X=x) = \sum_{\{\omega: X(\omega)=x\}} P(\omega)$
 - $P(x)$ for short (when unambiguous)
 - $P(X)$ refers to the entire distribution (think of it as a vector or table)



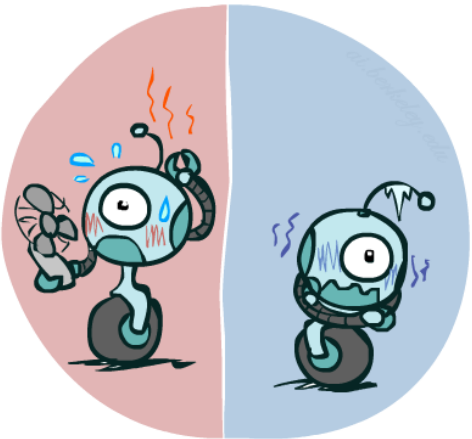
Probability Distributions

- Associate a probability with each value; sums to 1

- Temperature:

$P(T)$

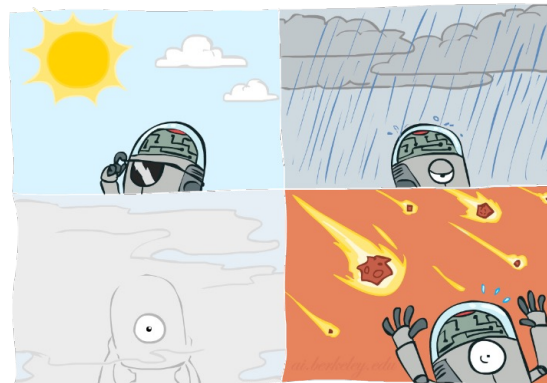
T	P
hot	0.5
cold	0.5



- Weather:

$P(W)$

W	P
sun	0.6
rain	0.1
fog	0.3
meteor	0.0



- Joint distribution*

$P(T,W)$

		Temperature	
		hot	cold
Weather	sun	0.45	0.15
	rain	0.02	0.08
	fog	0.03	0.27
	meteor	0.00	0.00

Making possible worlds

- In many cases we
 - begin with random variables and their domains
 - construct possible worlds as assignments of values to all variables
- E.g., two dice rolls $Roll_1$ and $Roll_2$
 - How many possible worlds?
 - What are their probabilities?
- Size of distribution for n variables with range size d ? d^n
- For all but the smallest distributions, cannot write out by hand!

Probabilities of events

- Recall that the probability of an event is the sum of probabilities of its worlds:
 - $P(A) = \sum_{\omega \in A} P(\omega)$
- So, given a joint distribution over all variables, can compute any event probability!
 - Probability that it's hot AND sunny?
 - Probability that it's hot?
 - Probability that it's hot OR not foggy?

- *Joint distribution*

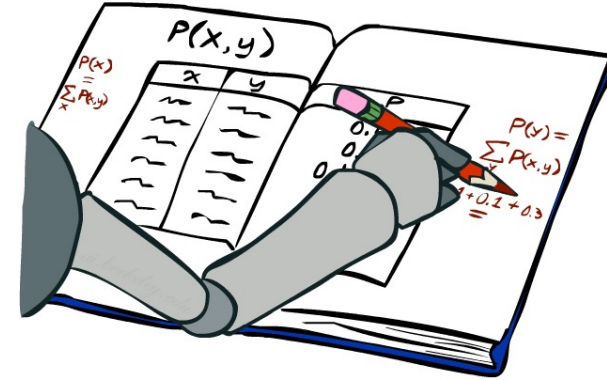
$P(T, W)$

		Temperature	
		hot	cold
Weather	sun	0.45	0.15
	rain	0.02	0.08
	fog	0.03	0.27
	meteor	0.00	0.00

Marginal Distributions

- Marginal distributions are sub-tables which eliminate variables
- **Marginalization** (*summing out*): Collapse a dimension by adding

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



		Temperature		
		hot	cold	
Weather	sun	0.45	0.15	0.60
	rain	0.02	0.08	0.10
	fog	0.03	0.27	0.30
	meteor	0.00	0.00	0.00
		0.50	0.50	

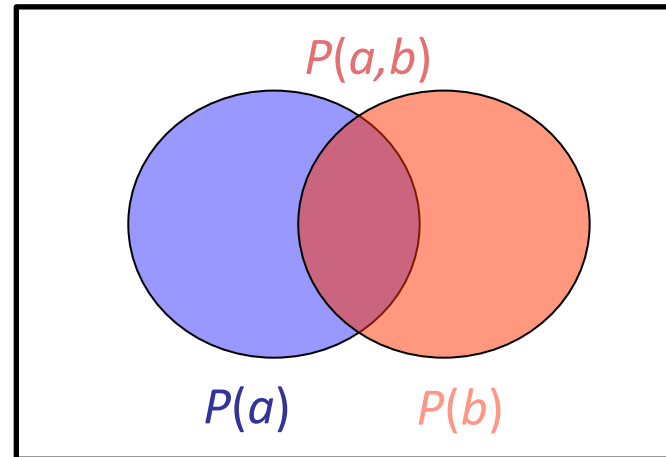
$P(W)$

$P(T)$

Conditional Probabilities

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

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



$P(T, W)$

		Temperature	
		hot	cold
Weather	sun	0.45	0.15
	rain	0.02	0.08
	fog	0.03	0.27
	meteor	0.00	0.00

$$P(W=s \mid T=c) = \frac{P(W=s, T=c)}{P(T=c)} = 0.15/0.50 = 0.3$$

$$\begin{aligned} &= P(W=s, T=c) + P(W=r, T=c) + P(W=f, T=c) + P(W=m, T=c) \\ &= 0.15 + 0.08 + 0.27 + 0.00 = 0.50 \end{aligned}$$

Conditional Distributions

- Distributions for one set of variables given another set

		Temperature	
		hot	cold
Weather	sun	0.45	0.15
	rain	0.02	0.08
	fog	0.03	0.27
	meteor	0.00	0.00

$P(W \mid T=h)$

hot

0.90
0.04
0.06
0.00

$P(W \mid T=c)$

cold

0.30
0.16
0.54
0.00

$P(W \mid T)$

hot

cold

0.90	0.30
0.04	0.16
0.06	0.54
0.00	0.00

Normalizing a distribution

- (Dictionary) To bring or restore to a normal condition

All entries sum to ONE

- Procedure:

- Multiply each entry by $\alpha = 1/(\text{sum over all entries})$

$P(W, T)$

		Temperature	
		hot	cold
Weather	sun	0.45	0.15
	rain	0.02	0.08
	fog	0.03	0.27
	meteor	0.00	0.00

$P(W, T=c)$

0.15
0.08
0.27
0.00

$$P(W | T=c) = P(W, T=c) / P(T=c) \\ = \alpha P(W, T=c)$$

Normalize

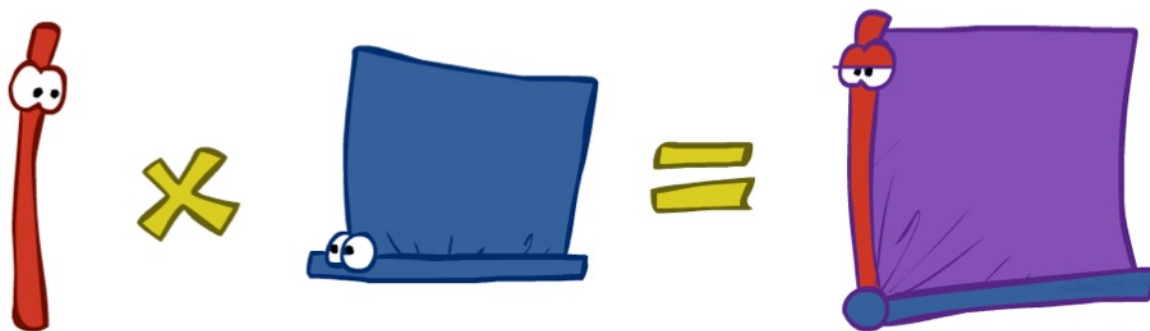
$$\alpha = 1/0.50 = 2$$

0.30
0.16
0.54
0.00

The Product Rule

- Sometimes have conditional distributions but want the joint

$$P(a \mid b) P(b) = P(a, b) \quad \longleftrightarrow \quad P(a \mid b) = \frac{P(a, b)}{P(b)}$$



The Product Rule: Example

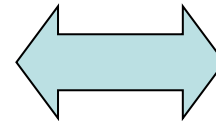
$$P(W \mid T) P(T) = P(W, T)$$

$P(W \mid T)$

	hot	cold
0.90	0.90	0.30
0.04	0.04	0.16
0.06	0.06	0.54
0.00	0.00	0.00

$P(T)$

T	P
hot	0.5
cold	0.5



$P(W, T)$

		Temperature	
		hot	cold
Weather	sun	0.45	0.15
	rain	0.02	0.08
	fog	0.03	0.27
	meteor	0.00	0.00

The Chain Rule

- A joint distribution can be written as a product of conditional distributions by repeated application of the product rule:
- $P(x_1, x_2, x_3) = P(x_3 \mid x_1, x_2) P(x_1, x_2) = P(x_3 \mid x_1, x_2) P(x_2 \mid x_1) P(x_1)$
- $P(x_1, x_2, \dots, x_n) = \prod_i P(x_i \mid x_1, \dots, x_{i-1})$

Probabilistic Inference

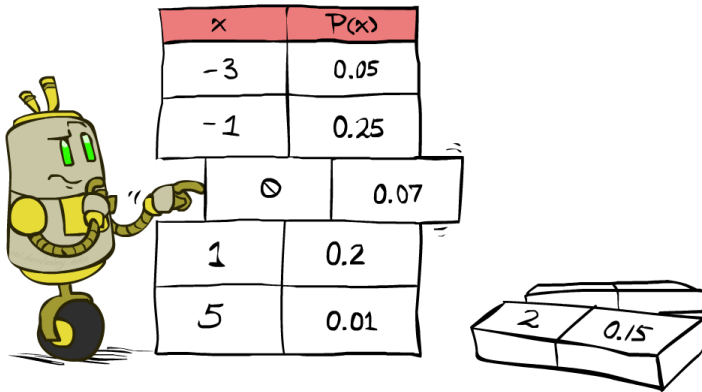
- Probabilistic inference: compute a desired probability from a probability model
 - Typically for a **query variable** given **evidence**
 - E.g., $P(\text{airport on time} \mid \text{no accidents}) = 0.90$
 - These represent the agent's *beliefs* given the evidence
- Probabilities change with new evidence:
 - $P(\text{airport on time} \mid \text{no accidents, 5 a.m.}) = 0.95$
 - $P(\text{airport on time} \mid \text{no accidents, 5 a.m., raining}) = 0.80$
 - Observing new evidence causes **beliefs to be updated**



Inference by Enumeration

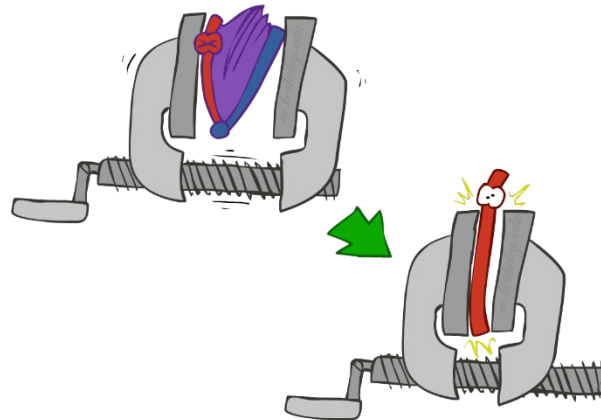
- Probability model $P(X_1, \dots, X_n)$ is given
- Partition the variables X_1, \dots, X_n into sets as follows:
 - Evidence variables: $E = e$
 - Query variables: Q
 - Hidden variables: H
- We want:
$$P(Q \mid e)$$

- Step 1: Select the entries consistent with the evidence



- Step 2: Sum out H from model to get joint of query and evidence

$$P(Q, e) = \sum_h \underbrace{P(Q, h, e)}_{X_1, \dots, X_n}$$



- Step 3: Normalize

$$P(Q \mid e) = \alpha P(Q, e)$$

Inference by Enumeration

- $P(W)$?

Season	Temp	Weather	P
summer	hot	sun	0.35
summer	hot	rain	0.01
summer	hot	fog	0.01
summer	hot	meteor	0.00
summer	cold	sun	0.10
summer	cold	rain	0.05
summer	cold	fog	0.09
summer	cold	meteor	0.00
winter	hot	sun	0.10
winter	hot	rain	0.01
winter	hot	fog	0.02
winter	hot	meteor	0.00
winter	cold	sun	0.15
winter	cold	rain	0.20
winter	cold	fog	0.18
winter	cold	meteor	0.00

Inference by Enumeration

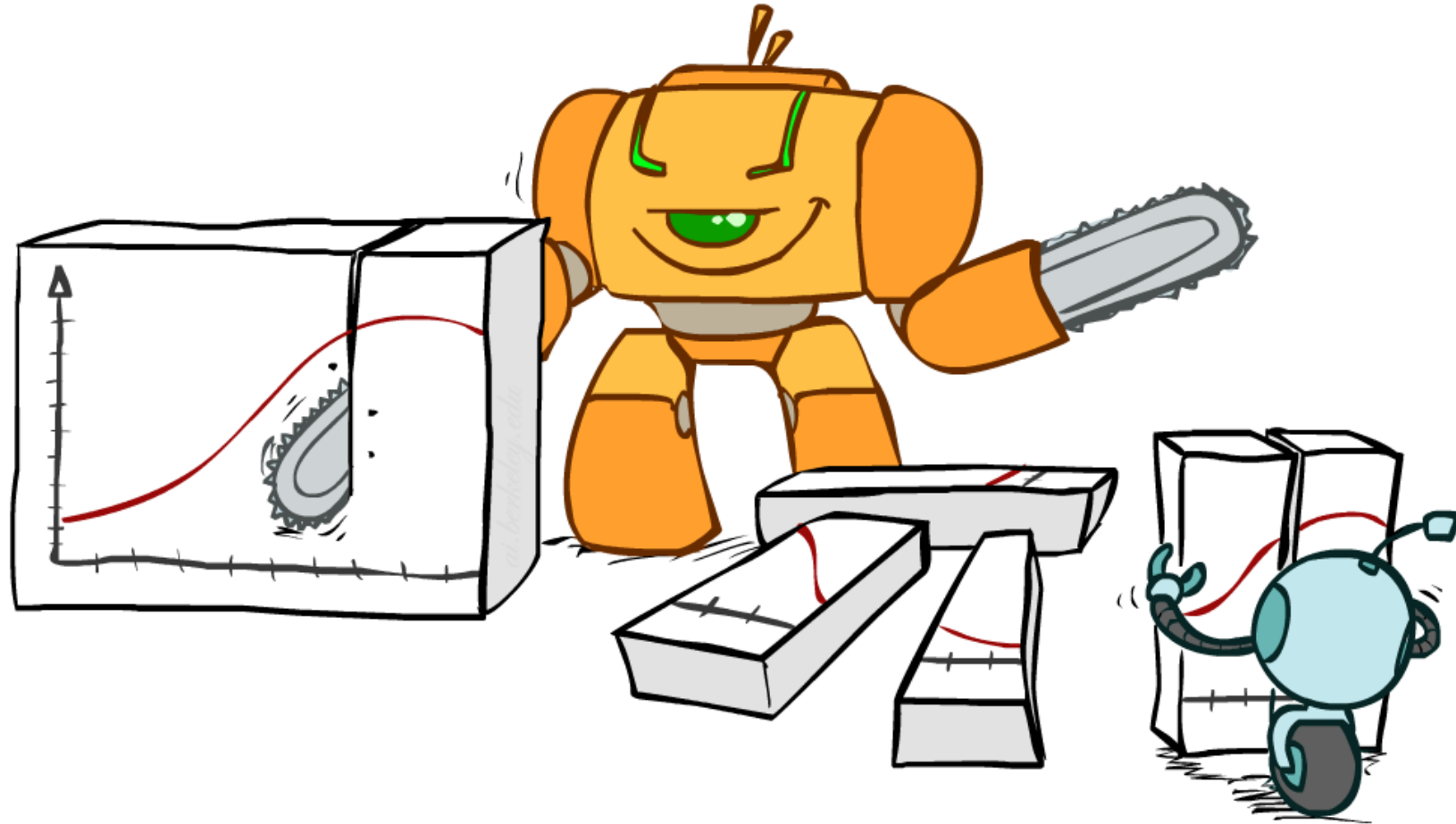
- $P(W)$?
- $P(W \mid \text{winter})$?

Season	Temp	Weather	P
summer	hot	sun	0.35
summer	hot	rain	0.01
summer	hot	fog	0.01
summer	hot	meteor	0.00
summer	cold	sun	0.10
summer	cold	rain	0.05
summer	cold	fog	0.09
summer	cold	meteor	0.00
winter	hot	sun	0.10
winter	hot	rain	0.01
winter	hot	fog	0.02
winter	hot	meteor	0.00
winter	cold	sun	0.15
winter	cold	rain	0.20
winter	cold	fog	0.18
winter	cold	meteor	0.00

Inference by Enumeration

- Obvious problems:
 - Worst-case time complexity $O(d^n)$ (exponential in #hidden variables)
 - Space complexity $O(d^n)$ to store the joint distribution
 - $O(d^n)$ data points to estimate the entries in the joint distribution

Bayes Rule



Bayes' Rule

- Write the product rule both ways:

$$P(a | b) P(b) = P(a, b) = P(b | a) P(a)$$

- Dividing left and right expressions, we get:

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

- 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
- Describes an “update” step from prior $P(a)$ to posterior $P(a | b)$
- 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:

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

- Example:

- M: meningitis, S: stiff neck

$$\left. \begin{array}{l} P(s \mid m) = 0.8 \\ P(m) = 0.0001 \\ P(s) = 0.01 \end{array} \right\} \begin{array}{l} \text{Example} \\ \text{gives} \end{array}$$

$$P(m \mid s) = \frac{P(s \mid m) P(m)}{P(s)} = \frac{0.8 \times 0.0001}{0.01}$$

- Note: posterior probability of meningitis still very small: 0.008 (80x bigger – why?)
- Note: you should still get stiff necks checked out! Why?

Independence

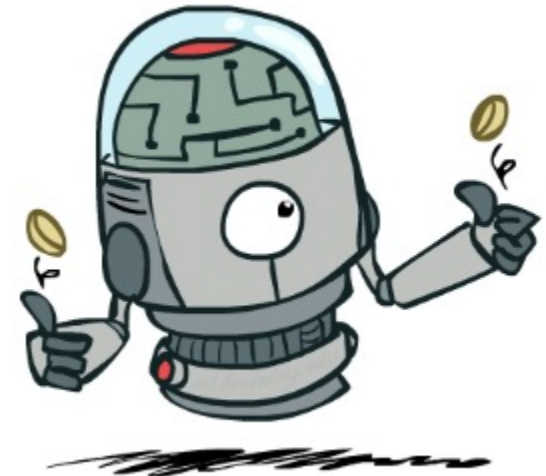
- Two variables X and Y are (absolutely) **independent** if

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

- I.e., the joint distribution **factors** into a product of two simpler distributions
- Equivalently, via the product rule $P(x, y) = P(x | y) P(y)$,

$$P(x | y) = P(x) \quad \text{or} \quad P(y | x) = P(y)$$

- Example: two dice rolls $Roll_1$ and $Roll_2$
 - $P(Roll_1=5, Roll_2=3) = P(Roll_1=5) P(Roll_2=3) = 1/6 \times 1/6 = 1/36$
 - $P(Roll_2=3 | Roll_1=5) = P(Roll_2=3)$



Example: Independence

- n fair, independent coin flips:

$P(X_1)$

H	0.5
T	0.5

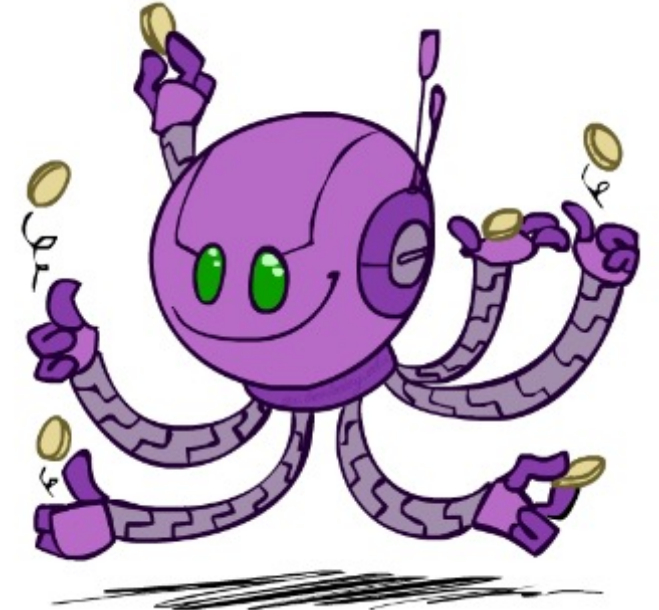
$P(X_2)$

H	0.5
T	0.5

...

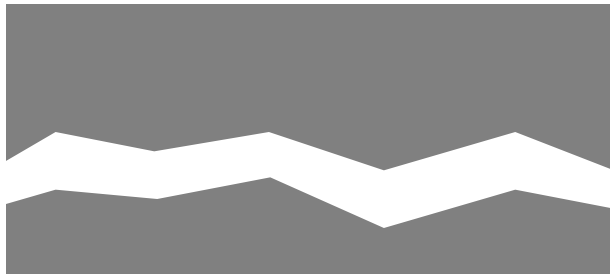
$P(X_n)$

H	0.5
T	0.5

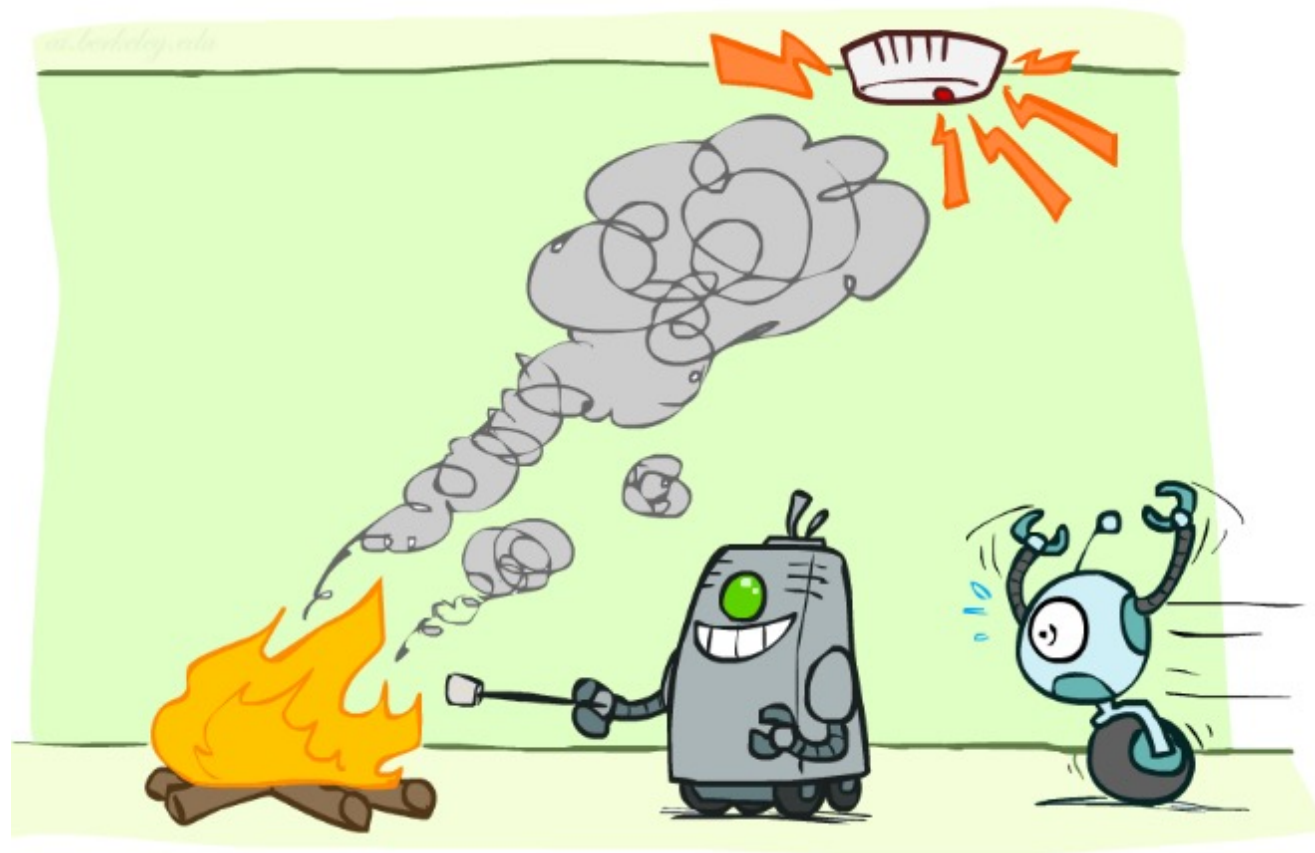


$P(X_1, X_2, \dots, X_n)$

2^n

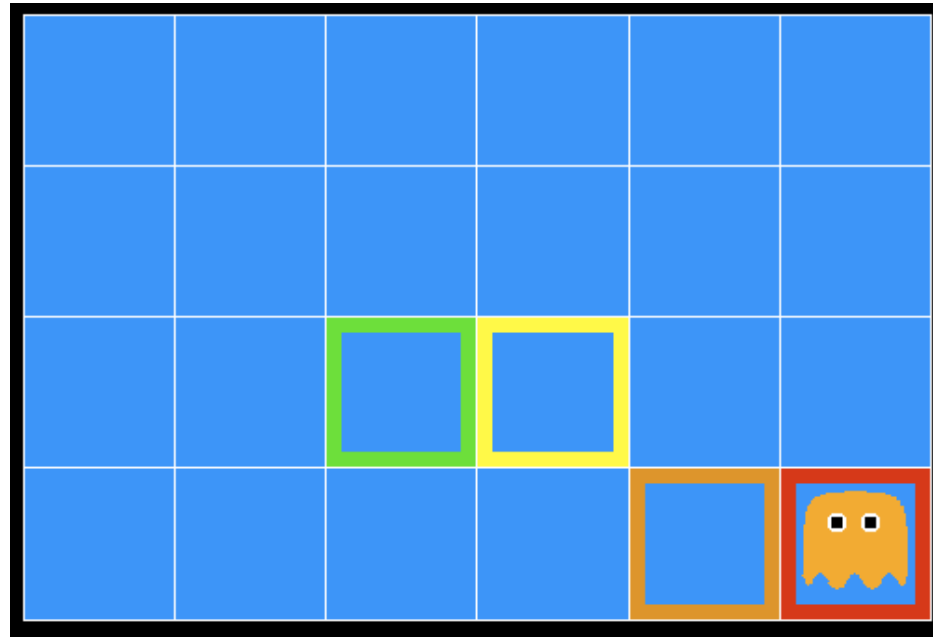


Conditional Independence



Ghostbusters

- A ghost is in the grid somewhere
- Sensor readings tell how close a square is to the ghost
 - On the ghost: usually red
 - 1 or 2 away: mostly orange
 - 3 or 4 away: typically yellow
 - 5+ away: often green
- Click on squares until confident of location, then “*bust*”



Video of Demo Ghostbusters with Probability



Ghostbusters model

- Variables and ranges:

- G (ghost location) in $\{(1,1),\dots,(3,3)\}$
- $C_{x,y}$ (color measured at square x,y) in $\{\text{red,orange,yellow,green}\}$


0.11	0.11	0.11
0.11	0.11	0.11
0.11	0.11	0.11

- Ghostbuster physics:

- **Uniform prior distribution** over ghost location: $P(G)$
- **Sensor model**: $P(C_{x,y} \mid G)$ (depends only on distance to G)
 - E.g. $P(C_{1,1} = \text{yellow} \mid G = (1,1)) = 0.1$

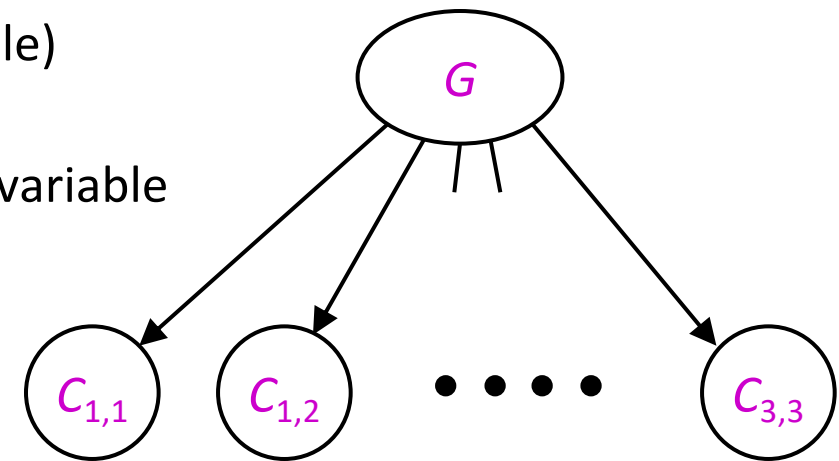
Ghostbusters model, contd.

- $P(G, C_{1,1}, \dots, C_{3,3})$ has $9 \times 4^9 = 2,359,296$ entries!!!
- Ghostbuster independence:
 - Are $C_{1,1}$ and $C_{1,2}$ independent?
 - E.g., does $P(C_{1,1} = \text{yellow}) = P(C_{1,1} = \text{yellow} \mid C_{1,2} = \text{orange})$?
- Ghostbuster physics again:
 - $P(C_{x,y} \mid G)$ ***depends only on distance to G***
 - So $P(C_{1,1} = \text{yellow} \mid \underline{G = (2,3)}) = P(C_{1,1} = \text{yellow} \mid \underline{G = (2,3)}, C_{1,2} = \text{orange})$
 - I.e., $C_{1,1}$ is ***conditionally independent*** of $C_{1,2}$ ***given G***

0.11		0.11
0.11	0.11	0.11
0.11	0.11	0.11

Ghostbusters model, contd.

- Apply the chain rule to decompose the joint probability model:
- $P(G, C_{1,1}, \dots, C_{3,3}) = P(G) P(C_{1,1} \mid G) P(C_{1,2} \mid G, C_{1,1}) P(C_{1,3} \mid G, C_{1,1}, C_{1,2}) \dots P(C_{3,3} \mid G, C_{1,1}, \dots, C_{3,2})$
- Now simplify using conditional independence:
- $P(G, C_{1,1}, \dots, C_{3,3}) = P(G) P(C_{1,1} \mid G) P(C_{1,2} \mid G) P(C_{1,3} \mid G) \dots P(C_{3,3} \mid G)$
- I.e., conditional independence properties of ghostbuster physics simplify the probability model from **exponential** to **quadratic** in the number of squares
- This is called a **Naïve Bayes** model:
 - One discrete query variable (often called the **class** or **category** variable)
 - All other variables are (potentially) evidence variables
 - Evidence variables are all conditionally independent given the query variable



Conditional Independence

- ***Conditional independence*** is our most basic and robust form of knowledge about uncertain environments.

- X is conditionally independent of Y given Z if and only if:

$$\forall x, y, z \quad P(x \mid y, z) = P(x \mid z)$$

or, equivalently, if and only if

$$\forall x, y, z \quad P(x, y \mid z) = P(x \mid z) P(y \mid z)$$

Conditional Independence

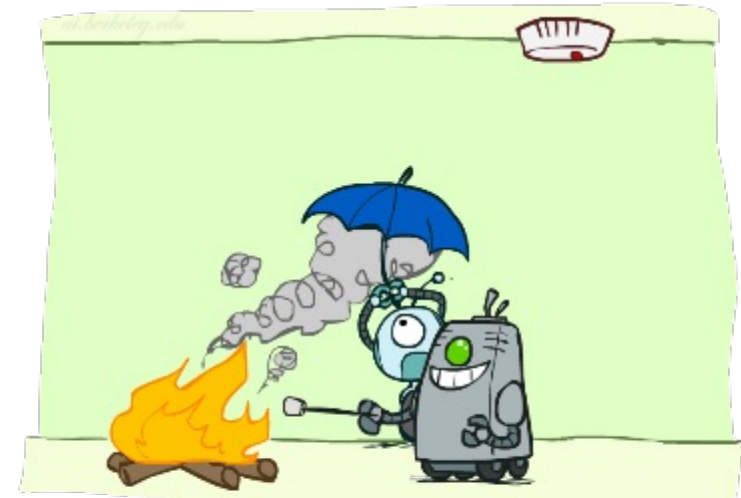
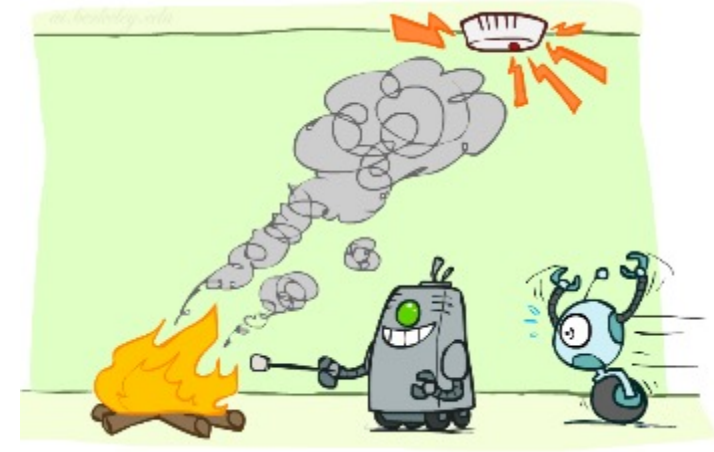
- What about this domain:
 - Traffic
 - Umbrella
 - Raining



Conditional Independence

- What about this domain:

- Fire
- Smoke
- Alarm



To Summarize ...

- Basic laws: $0 \leq P(\omega) \leq 1$, $\sum_{\omega \in \Omega} P(\omega) = 1$, $P(A) = \sum_{\omega \in A} P(\omega)$
- Random variable $X(\omega)$ has a value in each ω
 - Distribution $P(X)$ gives probability for each possible value x
 - Joint distribution $P(X,Y)$ gives total probability for each combination x,y
- Summing out/marginalization: $P(X=x) = \sum_y P(X=x,Y=y)$
- Conditional probability: $P(X|Y) = P(X,Y)/P(Y)$
- Chain rule: $P(X_1, \dots, X_n) = \prod_i P(X_i | X_1, \dots, X_{i-1})$
- Bayes Rule: $P(X|Y) = P(Y|X)P(X)/P(Y)$
- Independence: $P(X,Y) = P(X) P(Y)$ or $P(X|Y) = P(X)$ or $P(Y|X) = P(Y)$
- Conditional Independence: $P(X|Y,Z) = P(X|Z)$ or $P(X,Y|Z) = P(X|Z) P(Y|Z)$

Next time

- Bayes nets
- Elementary inference in Bayes nets