

Artificial Intelligence

12. Probability Review

Shashi Prabh

School of Engineering and Applied Science
Ahmedabad University

Uncertainty

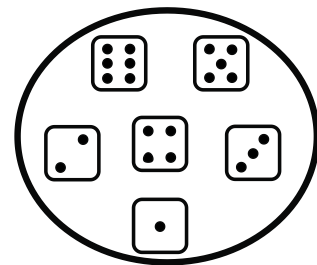
- The real world is rife with uncertainty!
 - If I leave for SFO 60 minutes before my flight, will I be there in time?
- Problems:
 - partial observability (road state, other drivers' plans, etc.)
 - noisy sensors (radio traffic reports, Google maps)
 - immense complexity of modelling & predicting traffic, security line, etc.
 - lack of knowledge of world dynamics (will tire burst? need COVID test?)

Uncertainty

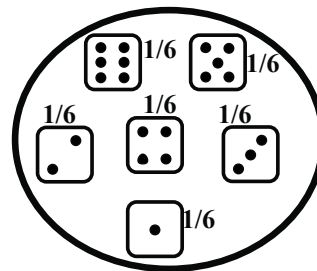
- Probabilistic assertions summarize effects of **ignorance** and **laziness**
- Probability theory + Utility theory \Rightarrow Decision theory
- **Maximize expected utility**
 - $a^* = \operatorname{argmax}_a \sum_s P(s \mid a) U(s)$

Basic laws of probability

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



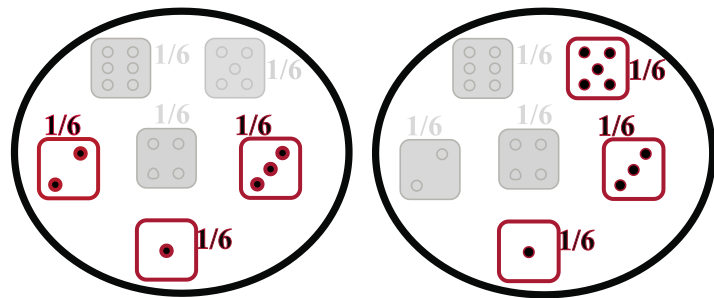
- A **probability model** assigns a number $P(\omega)$ to each world ω
 - E.g., $P(1) = P(2) = P(3) = P(4) = P(5) = P(6) = 1/6$.



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

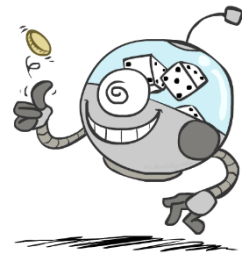
Basic laws contd.

- 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., $P(\text{roll} < 4) = P(1) + P(2) + P(3) = 1/2$
- De Finetti (1931): anyone who bets according to probabilities that violate these laws can be forced to lose money on every set of bets
 - No rational agent can have beliefs that violate probability axioms



Random Variables

- A random variable 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? $\rightarrow \{\text{true}, \text{false}\}$
 - e.g. **Odd**(1)=true, **Odd**(6) = false
 - often write the event **Odd**=true as **odd**, **Odd**=false as $\neg\text{odd}$
 - **T** = Is it hot or cold? $\rightarrow \{\text{hot}, \text{cold}\}$
 - **D** = How long will it take to get to the airport? $\rightarrow [0, \infty)$
 - **L_{Wumpus}** = Where is the wumpus? $\rightarrow \{(0,0), (0,1), \dots\}$



Random Variables

- 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)

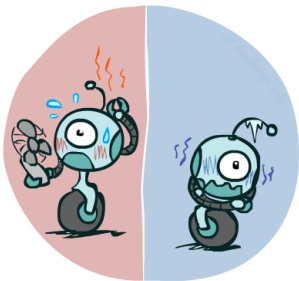
Probability Distributions

- Associate a probability with each value; sums to 1

- Temperature:

$P(T)$

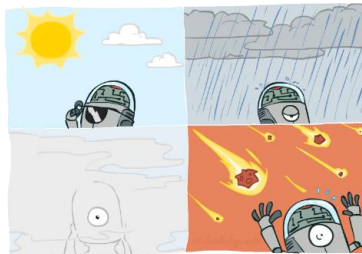
T	P
hot	0.5
cold	0.5



- 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!

- **Joint distribution**

$P(T, W)$

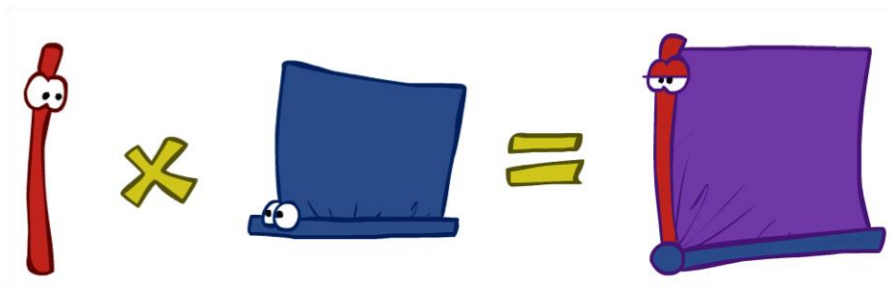
- Probability that it's hot AND sunny?
- Probability that it's hot?
- Probability that it's hot OR not foggy?

		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 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

$$P(x_1, x_2, x_3) = P(x_3 \mid x_1, x_2) P(x_1, x_2) = P(x_3 \mid x_1, x_2) P(x_2 \mid x_1) P(x_1)$$

or,

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

Conditional Probabilities

- A simple relation between joint and conditional probabilities

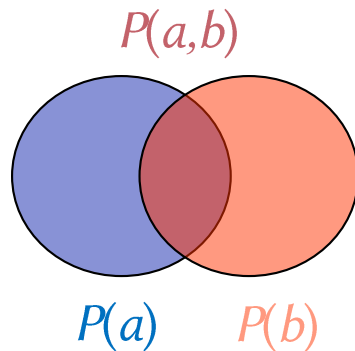
$$P(a | b) = \frac{P(a, b)}{P(b)}$$

- In fact, this is taken as the definition of conditional probability

$P(T, W)$

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

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



Conditional Probabilities

- A simple relation between joint and conditional probabilities

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

- In fact, this is taken as the definition of conditional probability

$P(T, W)$

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

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

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

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

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

Inference with Bayes' Rule

- Diagnostic probability from causal probability or likelihood:

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

- Likelihoods need not add to 1

Independence

- Two variables X and Y are **independent** if

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

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

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



Independence

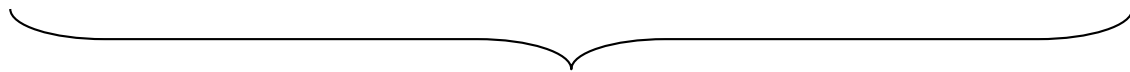
- Example: two dice rolls Roll_1 and Roll_2
 - $P(\text{Roll}_1=5, \text{Roll}_2=3) = P(\text{Roll}_1=5) P(\text{Roll}_2=3) = 1/6 \times 1/6 = 1/36$
 - $P(\text{Roll}_2=3 \mid \text{Roll}_1=5) = P(\text{Roll}_2=3)$



Example: Independence

- n fair, independent coin flips:

$P(X_1)$		$P(X_2)$		\dots		$P(X_n)$	
H	0.5	H	0.5			H	0.5
T	0.5	T	0.5			T	0.5



vs. full joint distribution 2^n

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



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:

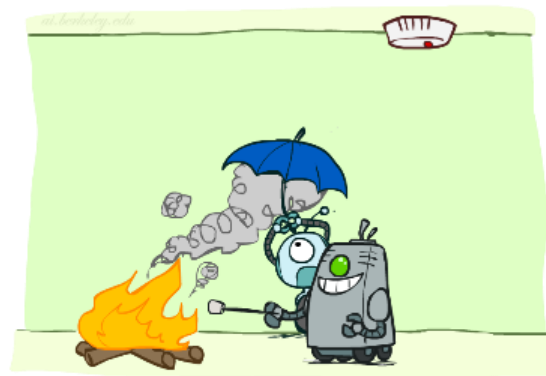
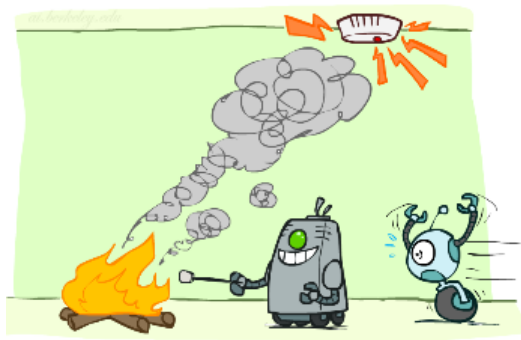
$$\forall x, y, z \quad P(x \mid y, z) = P(x \mid z)$$

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

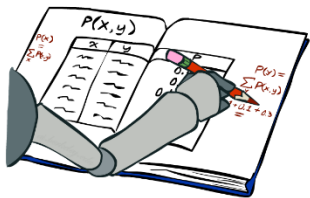
- What about this domain:
 - Fire
 - Smoke
 - Alarm



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	

$P(W)$

$P(T)$

Marginal Distributions

- $P(\text{cavity}) = ?$

	<i>toothache</i>		\neg <i>toothache</i>	
	<i>catch</i>	\neg <i>catch</i>	<i>catch</i>	\neg <i>catch</i>
<i>cavity</i>	0.108	0.012	0.072	0.008
\neg <i>cavity</i>	0.016	0.064	0.144	0.576

y

Bayes' Rule

- 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 the Bayes' Rule

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



- 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 AI systems
- In the running for the most important AI equation!

Inference with Bayes' Rule

- Diagnostic probability from causal probability or likelihood

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

- Example:

- M: meningitis
- S: stiff neck

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

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

- Posterior probability of meningitis still very small: **0.008**
 - You should still get stiff necks checked out! Why?

Probabilistic Inference

- Compute desired probability from a probability model
 - Typically for a **query variable** given **evidence**
 - $P(\text{airport on time} \mid \text{no accidents}) = 0.90$
 - These represent the agent's beliefs given the evidence
- **Probabilities can change with new evidence**
 - $P(\text{airport on time} \mid \text{no accidents, 5 a.m.}) = 0.95$
 - $P(\text{airport on time} \mid \text{no accidents, 5 a.m., raining}) = 0.80$
- **Beliefs are updated upon observing new evidences**




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: $\mathbf{E} = \mathbf{e}$
 - Query variables: \mathbf{Q}
 - Hidden variables: \mathbf{H}

■ We want: $P(\mathbf{Q} \mid \mathbf{e})$

- Step 1: Select the entries consistent with the evidence

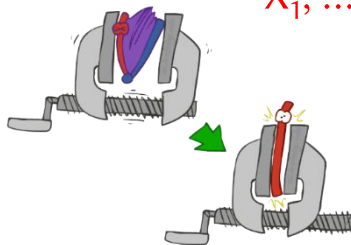


x	P(x)
-3	0.05
-1	0.25
0	0.07
1	0.2
5	0.01

2	0.15
---	------

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

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



- Step 3: Normalize

$$P(\mathbf{Q} \mid \mathbf{e}) = \alpha P(\mathbf{Q}, \mathbf{e})$$

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	0.01
summer	cold	rain	0.05
summer	cold	fog	0.10
summer	cold	meteor	0.00
winter	hot	sun	0.10
winter	hot	rain	0.01
winter	hot	fog	0.01
winter	hot	meteor	0.00
winter	cold	sun	0.10
winter	cold	rain	0.10
winter	cold	fog	0.15
winter	cold	meteor	0.00

Inference by Enumeration

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

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	0.01
summer	cold	rain	0.05
summer	cold	fog	0.10
summer	cold	meteor	0.00
winter	hot	sun	0.10
winter	hot	rain	0.01
winter	hot	fog	0.01
winter	hot	meteor	0.00
winter	cold	sun	0.10
winter	cold	rain	0.10
winter	cold	fog	0.15
winter	cold	meteor	0.00

Inference by Enumeration

- $P(\text{cavity} \mid \text{toothache}) = ?$
- $P(\neg\text{cavity} \mid \text{toothache}) = ?$

	<i>toothache</i>		$\neg\text{toothache}$	
	<i>catch</i>	$\neg\text{catch}$	<i>catch</i>	$\neg\text{catch}$
<i>cavity</i>	0.108	0.012	0.072	0.008
$\neg\text{cavity}$	0.016	0.064	0.144	0.576

Issues with Inference by Enumeration

- Worst-case time complexity $O(d^n)$
 - exponential in the number of hidden variables
- Space complexity $O(d^n)$ to store the joint distribution
- **All** the joint distribution entries must be estimated separately. That is $O(d^n)$ data points to estimate!
- We will use conditional independence to improve the inference complexity

- **Reading:** Chapter 12
- **Assignments:** PS 6
- Next:
 - Bayesian networks
 - Elementary inference in Bayesian networks