

Probability

In logic, every predicate is either true or false, and that's that.

- In the real world, that is still true:
 - Every alleged fact is either true or false
 - Every possible percept either appears or it doesn't
 - Every possible event either happens or doesn't
- But very often we just don't know which it is, true or false
- Standard logic can only handle certainty, true or false and nothing else
- Probability theory is one way of dealing with that
- The probability of a predicate, percept, or event
 - is a number on a continuous spectrum from 0 to 1
 - sometimes expressed as a percentage
 - $P(\text{event}) = 0$ means it is absolutely impossible for it to happen
 - $P(\text{event}) = 1$ means it is absolutely certain to happen
- But what do numbers *between* 0 and 1 mean?
 - e.g. weather forecast says 40% chance of rain
 - e.g. $P(\text{unfair-coin-lands-heads-up}) = 0.4$
 - Average result of an "infinite" number of identical experiments
 - Proportion of possible worlds in which it is true
 - 40¢ is a fair price to pay for a bet that would pay out \$1
- Stochastic \approx probability-based

Basic laws (Venn diagrams)

- $P(A \wedge B) = P(A) \times P(B)$
- $P(A \vee B) = P(A) + P(B) - P(A \wedge B)$
- $P(\neg A) = 1 - P(A)$
- Independent variables

Sometimes probabilities are Conditional

- The chances of an event happening depend on whether some other event has happened
- | means "given" or roughly "if we know that"
- For example, tooth aches do not happen much (if you clean your teeth)
 - $P(\text{toothache}) = 0.005$
 - but
 - $P(\text{toothache} \mid \text{cavity}) = 0.3$
 - This does not denote cause and effect, it could also be that
 - $P(\text{cavity} \mid \text{toothache}) = 0.75$
 - Useful for diagnostic purposes
 - It is still true that $P(\text{toothache}) = 0.005$
 - even after cavity has been observed
- $P(A \mid B) = P(A \wedge B) / P(B)$
 - Observing B rules out all cases where B is false
 - leaving the set of possibilities with total probability of just $P(B)$
 - within those remaining possibilities,
 - for $P(A \mid B)$ to be true, A must be true
 - so A and B is true

so the true probability of $A | B$ must be $P(A \wedge B)$
divided by $P(B)$, to make all the probabilities add up to 1

- $P(A \wedge B) = P(A | B) \times P(B)$ is sometimes a more convenient form
For A and B to be true, we need B to be true, then given that B is true we also need A to be true

Notation

- If the range of possible values for Weather is [sunny, rainy, cloudy, snowy] we might have
 $P(\text{Weather} = \text{sunny}) = 0.6$
 $P(\text{Weather} = \text{rainy}) = 0.1$
 $P(\text{Weather} = \text{cloudy}) = 0.29$
 $P(\text{Weather} = \text{snowy}) = 0.01$
 (these are not independent, so the \vee rule doesn't work)
 this is often written as
 $P(\text{Weather}) = \langle 0.6, 0.1, 0.29, 0.01 \rangle$
- Sometimes $P(\text{sunny})$ is written as an abbreviation for $P(\text{Weather} = \text{sunny})$
- For continuous variables, a Probability Density Function is used
- Probability = the area under the curve
- Numeric probabilities only really make sense for ranges of possible values

Joint distributions

- If we have three Boolean variables, Toothache, Cavity, and Catch the joint distribution is a $2 \times 2 \times 2$ table, all eight add up to exactly 1
- Use a bold **P** to represent that
- $P(\text{some possibility})$ can be found by adding up the relevant entries
- Marginalisation - eliminate variable(s) by summing entries
 $P(\text{some possibility}) = \sum \text{for all possible } x\text{'s of } P(\text{that possibility} \wedge X=x)$
- Conditioning
 $P(\text{some possibility}) = \sum \text{for all possible } x\text{'s of } P(\text{that possibility} | X=x) \times P(X=x)$

Bayes' rule

- From $P(A \wedge B) = P(A | B) \times P(B)$ and $P(A \wedge B) = P(B | A) \times P(A)$ we easily get Bayes' rule:
 $P(B | A) = P(A | B) \times P(B) / P(A)$
- When diagnosing an illness, a doctor might know all of
 $P(\text{effect} | \text{cause})$
 $P(\text{effect})$
 $P(\text{cause})$
 for a vast collection of different causes and effects,
 just from hundreds of years of observations and studies
 And from them $P(\text{cause} | \text{effect})$ can be calculated,
 the probability of a particular disease being the cause of the
 observed effects
- In an epidemic, $P(\text{cause})$ for one particular cause will increase a lot
 previously observed values for $P(\text{cause} | \text{effect})$ will become invalid

Combining evidence

- If we observe both toothache and catch, what is the probability distribution for Cavity?
- Assume this joint distribution:

	Toothache		¬Toothache	
	Catch	¬Catch	Catch	¬Catch
Cavity	.108	.012	.072	.008
¬Cavity	.016	.064	.144	.576

(they all add up to 1)

- If we have a joint distribution we can just add up the right entries
 $\mathbf{P}(\text{Cavity} \mid \text{toothache} \wedge \text{catch}) = \alpha \langle 0.108, 0.016 \rangle$ perhaps
 α stands for Normalise:
 multiply by something to make the probabilities add up to 1
 so $\alpha \langle 0.108, 0.016 \rangle = \langle 0.871, 0.129 \rangle$
- Does not scale up if there are a lot of variables
- Using Bayes' rule
 $\mathbf{P}(\text{Cavity} \mid \text{toothache} \wedge \text{catch}) =$
 $\alpha \mathbf{P}(\text{toothache} \wedge \text{catch} \mid \text{Cavity}) \times \mathbf{P}(\text{Cavity})$
- Again, there are probably too many variables for this to be practical
- If the variables were independent, we'd be better off
 but they aren't. Toothache and Catch are not unrelated
- But if, given knowledge of Cavity, they become independent
- From $\mathbf{P}(\text{toothache} \wedge \text{catch} \mid \text{Cavity}) =$
 $\mathbf{P}(\text{toothache} \mid \text{Cavity}) \times \mathbf{P}(\text{catch} \mid \text{Cavity})$
 we get $\mathbf{P}(\text{toothache} \wedge \text{catch} \mid \text{Cavity}) =$
 $\alpha \mathbf{P}(\text{toothache} \mid \text{Cavity}) \times \mathbf{P}(\text{catch} \mid \text{Cavity}) \times \mathbf{P}(\text{Cavity})$
- Not so many combinations to worry about now
- Conditional independence allows scalability

Naive Bayes Models

- If all the effects are conditionally independent given Cause, then
 $\mathbf{P}(\text{Cause} \wedge \text{Effect}_1 \wedge \text{Effect}_2 \wedge \dots) = \mathbf{P}(\text{Cause}) \times \prod \mathbf{P}(\text{Effect}_i \mid \text{Cause})$
- Naive because it relies on independence, but is sometimes used when there is none
- $\mathbf{P}(\text{Cause} \mid \text{Effect}=e) = \alpha \mathbf{P}(\text{Cause}) \times \prod \mathbf{P}(\text{Effect}=e_i \mid \text{Cause})$
- Example: text classification

Given some text, work out which section of the newspaper it came from
 We can know from prior observations these distributions

$\mathbf{P}(\text{Section})$

$\mathbf{P}(\text{Hasword}_w \mid \text{Section})$

If 9% of all articles are in the weather section then

$\mathbf{P}(\text{Section}=\text{weather}) = 0.09$

If 23% of all articles in the weather section use the word "rain" then

$\mathbf{P}(\text{Hasword}_{\text{rain}}=\text{true} \mid \text{Section}=\text{weather}) = 0.23$

From that and the current text itself, we can work out the probabilities of the text being from any given section