Statistical learning in the nervous system/ Intelligent systems Probabilistic models

Gergő Orbán golab.wigner.mta.hu





Course layout

Introduction of the mathematical framework Probabilistic models Probabilistic models in cognition: Bayesian behaviour Approximate methods Approximate methods in cognition: Sampling Discovering the probabilistic models used by humans Using learned statistical regularities: coding and compression Theory of coding and compression in memory processes Effective probabilistic models: nonparametric models Nonparametric models in cognition: memory dynamics Learning the structure of probabilistic models Structure learning in cognition: causal learning

RECAP: rational analysis



Horace Barlow Trinity College, Cambridge

A wing would be the most mystifying structure if one did not know that birds flew. One might observe that it could be extended a considerable distance, that it had smooth covering of feathers with conspicuous markings, that it was operated by powerful muscles, and that strength and lightness were prominent features of its constructions. These are important facts, but in themselves they do not tell us that birds fly. Yet, without knowing this and without understanding something of the principles of flight, a more detailed examination of the wing itself would probably be unrewarding.

RECAP: rational analysis



Horace Barlow Trinity College, Cambridge

A wing would be the most mystifying structure if one did not know that birds flew. One might observe that it could be extended a considerable distance, that it had smooth covering of feathers with conspicuous markings, that it was operated by powerful muscles, and that strength and lightness were prominent features of its constructions. These are important facts, but in themselves they do not tell us that birds fly. Yet, without knowing this and without understanding something of the principles of flight, a more detailed examination of the wing itself would probably be unrewarding.

Sensory communication (1962) MIT Press

 propositional logic: reasoning can be systematically performed

logic function: new propositions using basic operations (conjunction, disjunction, negation, implication)

 $C = f(A, B) = (A + \overline{B})(\overline{A} + A\overline{B}) + \overline{A}B(A + B)$

- propositional logic: reasoning can be systematically performed
- a finite set of operations is sufficient to cover the full proposition space



- propositional logic: reasoning can be systematically performed
- a finite set of operations is sufficient to cover the full proposition space
- rational analysis can be performed in a world where truth values of statements can be established

Turing machine an agent can be constructed that can perform an arbitrary function $x \rightarrow TM \rightarrow f(x)$ Universal Turing machine irrespective of the language used one can construct such computing machines $x \rightarrow TM \rightarrow f(x)$

- propositional logic: reasoning can be systematically performed
- a finite set of operations is sufficient to cover the full proposition space
- rational analysis can be performed in a world where truth values of statements can be established



- propositional logic: reasoning can be systematically performed
- a finite set of operations is sufficient to cover the full proposition space
- rational analysis can be performed in a world where truth values of statements can be established



- propositional logic: reasoning can be systematically performed
- a finite set of operations is sufficient to cover the full proposition space
- rational analysis can be performed in a world where truth values of statements can be established



- propositional logic: reasoning can be systematically performed
- a finite set of operations is sufficient to cover the full proposition space
- rational analysis can be performed in a world where truth values of statements can be established



- propositional logic: reasoning can be systematically performed
- a finite set of operations is sufficient to cover the full proposition space
- rational analysis can be performed in a world where truth values of statements can be established









• Scope of deductive reasoning is limited: we need to make inferences in uncertain situations

"The girl saw the boy with the telescope"



 Scope of deductive reasoning is limited: we need to make inferences in uncertain situations



"The girl saw the boy with the telescope"

• Scope of deductive reasoning is limited: we need to make inferences in uncertain situations



"The girl saw the boy with the telescope"

- Scope of deductive reasoning is limited: we need to make inferences in uncertain situations
- It is unclear how to construct adaptive forms of reasoning: training is hard, hand engineering is required

• Are there ways to formalise plausible reasoning?

- Are there ways to formalise plausible reasoning?
- Is there a rational way to approach this problem?

- Are there ways to formalise plausible reasoning?
- Is there a rational way to approach this problem?
- What exactly determines how much the observation of Aristotle's reading skills alters the plausibility of him being a philosopher?

- Are there ways to formalise plausible reasoning?
- Is there a rational way to approach this problem?
- What exactly determines how much the observation of Aristotle's reading skills alters the plausibility of him being a philosopher?
- In deductive reasoning, the veracity of statements can be identified after exhaustive chains of inferences. In plausible reasoning uncertainty grows with subsequent iterations. Can this be characterised?

- Are there ways to formalise plausible reasoning?
- Is there a rational way to approach this problem?
- What exactly determines how much the observation of Aristotle's reading skills alters the plausibility of him being a philosopher?
- In deductive reasoning, the veracity of statements can be identified after exhaustive chains of inferences. In plausible reasoning uncertainty grows with subsequent iterations. Can this be characterised?
- Discovering an efficient math behind plausible reasoning can provide a useful tool to model cognitive processes

- Are there ways to formalise plausible reasoning?
- Is there a rational way to approach this problem?
- What exactly determines how much the observation of Aristotle's reading skills alters the plausibility of him being a philosopher?
- In deductive reasoning, the veracity of statements can be identified after exhaustive chains of inferences. In plausible reasoning uncertainty grows with subsequent iterations. Can this be characterised?
- Discovering an efficient math behind plausible reasoning can provide a useful tool to model cognitive processes
- What does this do with 'weakly' rational behaviour at Wason's card selection task?

weak syllogisms

if A is true then B is true

if A is true then B is true

B is true

therefore, A becomes more plausible

A is false

weak syllogisms

if A is true then B is true

if A is true then B is true

B is true

A is false

therefore, A becomes more plausible

therefore, B becomes less plausible

• The truth value of many statements cannot be fully determined

weak syllogisms

if A is true then B is true

if A is true then B is true

B is true

A is false

therefore, A becomes more plausible

- The truth value of many statements cannot be fully determined
- Reasoning in uncertain situations is both useful and possible: we cannot rely on deductive reasoning

weak syllogisms

if A is true then B is true

if A is true then B is true

B is true

A is false

therefore, A becomes more plausible

- The truth value of many statements cannot be fully determined
- Reasoning in uncertain situations is both useful and possible: we cannot rely on deductive reasoning
- We do it automatically while making inferences, terming it *common sense*

weak syllogisms

if A is true then B is true

B is true

if A is true then B is true

A is false

therefore, A becomes more plausible

- The truth value of many statements cannot be fully determined
- Reasoning in uncertain situations is both useful and possible: we cannot rely on deductive reasoning
- We do it automatically while making inferences, terming it *common sense*
- Importantly, common sense is not an arbitrary relaxation of rules: George Pólya formalised a compact set of desiderata for reasoning

Artificial agent revisited

- We are pursuing common sense reasoning to understand thinking
- This provides better tools to build 'machines that think', i.e. better AI
- The rational agent developed provides a normative approach to human behaviour, in other words:
- 'rules will be deduced from simple desiderata which, it appears to us, would be desirable in human brains; i.e. we think that a rational person, on discovering that they were violating one of these desiderata, would wish to revise their thinking', *Jaynes*

• Aristotle wants to learn about the effectiveness of his teaching

what is the truth value of $S \rightarrow L$?

scribbling	excited look	learning	possible?
1	1	1	1
1	1	0	1
0	1	1	1
1	0	0	0
0	1	1	0
0	1	0	0
0	0	1	0
0	0	0	1

possible experiment outcomes

• Aristotle wants to learn about the effectiveness of his teaching

what is the truth value of $S \rightarrow L$?

scribbling	excited look	learning	possible?
1	1	1	1
1	1	0	1
0	1	1	1
1	0	0	0
0	1	1	0
0	1	0	0
0	0	1	0
0	0	0	1

possible experiment outcomes

How to establish numerical values for plausibilities?

Rational desiderata

- The proposition needs to be assigned a plausibility based on available evidence: the conditional plausibility that A is true, given that B is true, A | B
- The richness of real numbers is required by theory
- Convention: larger plausibility corresponds to higher numbers

combining predicaments:

the plausibility that A is true, given that both B and C are true:

A|BC

the plausibility that at least one of the propositions A and B is true, given that both C and D are true: A + B | CD

1. Level of certainty is expressed with real numbers
Rational desiderata

- Qualitative match of the formalism with common sense is required
- If (A|C') > (A|C)and (B|AC') = (B|AC)then $(AB|C') \ge (AB|C)$ & $(\overline{A}|C') < (\overline{A}|C)$
- Provides 'sense of direction'

1. Level of certainty is expressed with real numbers

Rational desiderata

- Qualitative match of the formalism with common sense is required
- If (A|C') > (A|C)and (B|AC') = (B|AC)then $(AB|C') \ge (AB|C)$ & $(\overline{A}|C') < (\overline{A}|C)$
- Provides 'sense of direction'

Level of certainty is expressed with real numbers
 Qualitative correspondence with common sense

Rational desiderata

- 3.a, Alternative ways to achieve conclusion lead to the same result
- 3.b, All evidence is taken into account by the agent
- 3.c, Equivalent state of knowledge is represented by equal level of plausibility

- 1. Level of certainty is expressed with real numbers
- 2. Qualitative correspondence with common sense
- 3. Consistency

Cox's theorem

Fundamental laws of probability (Kolmogorov's axioms) directly come from the three desiderata

- The three desiderata uniquely determine the mathematical form of representation of plausibility
- Plausibility has to behave as probability

Kolmogorov axioms

- 1. Probability is a real number between 0 and 1
- 2. Certainty is represented by P(A | B) = 1Certain falsehood: P(A | B) = 0
- 3. $P(A | B) + P(\overline{A} | B) = 1$
- +1 Conjunction: P(AB|C) = P(A|BC)P(B|C)

Cox's theorem

Fundamental laws of probability (Kolmogorov's axioms) directly come from the three desiderata

- The three desiderata uniquely determine the mathematical form of representation of plausibility
- Plausibility has to behave as probability

Kolmogorov axioms

- 1. Probability is a real number between 0 and 1
- 2. Certainty is represented by P(A | B) = 1Certain falsehood: P(A | B) = 0
- 3. $P(A | B) + P(\overline{A} | B) = 1$
- +1 Conjunction: P(AB|C) = P(A|BC)P(B|C)

We can use these to express operations in plausible reasoning

Cox's theorem

Fundamental laws of probability (Kolmogorov's axioms) directly come from the three desiderata

- The three desiderata uniquely determine the mathematical form of representation of plausibility
- Plausibility has to behave as probability

Kolmogorov axioms

- 1. Probability is a real number between 0 and 1
- 2. Certainty is represented by P(A | B) = 1Certain falsehood: P(A | B) = 0
- 3. $P(A | B) + P(\overline{A} | B) = 1$
- +1 Conjunction: P(AB|C) = P(A|BC)P(B|C)

We can use these to express operations in plausible reasoning

It will have similar power than deductive reasoning, but extends it

Dutch book argument

- Assume that we are willing to take bets with odds proportional to the probability of the occurrence of events
- If the plausibilities are not behaving according to the rules of probabilities then we will certainly loose money against someone using that strategy

Dutch book argument

- Assume that we are willing to take bets with odds proportional to the probability of the occurrence of events
- If the plausibilities are not behaving according to the rules of probabilities then we will certainly loose money against someone using that strategy

In an environment where truth values of statements cannot fully determined, a probabilistic description is lucrative

possible experiment outcoines					
scribbling	excited look	learning	P(S,E,L)		
1	1	1	0.2		
1	1	0	0.1		
0	1	1	0.06		
1	0	0	0.17		
0	0	0	0.43		

nassible avpariment autoomes

 $\sum P(A, E, L) = 1$

pos	nent outco	mes	
scribbling	excited look	learning	P(S,E,L)
1	1	1	0.2
1	1	0	0.1
0	1	1	0.06
1	0	0	0.17
0	0	0	0.43

- -

The probability table fully characterises the system

pos	sible experin	nent outco	mes	
scribbling	excited look	learning	P(S,E,L)	
1	1	1	0.2	
1	1	0	0.1	$\sum D(A - 1) = 1$
0	1	1	0.06	$\sum_{i=1}^{n} P(A, E, L) = 1$
1	0	0	0.17	
0	0	0	0.43	

The probability table fully characterises the system

pos	sible experin	nent outco	mes	
scribbling	excited look	learning	P(S,E,L)	
1	1	1	0.2	
1	1	0	0.1	$\sum D(A \in I) = 1$
0	1	1	0.06	$\sum F(A, E, L) = 1$
1	0	0	0.17	
0	0	0	0.43	

The probability table fully characterises the system

Instead of truth tables, we extended the scope of propositions

Goal: infer the probability of learning given the disciple wrote feverishly

T

scribbling	excited look	learning	P(S,E,L)
1	1	1	0.2
1	1	0	0.1
0	1	1	0.06
1	0	0	0.17
0	0	0	0.43

Infe	erence	& cond	litional	nroha	hility	
Goal: infer the probability of I $P(L \mid s, \overline{e}) = ?$ Stochastic variable S taking value s P(L \mid s, \overline{e}) = ?						
	scribbling	excited look	learning	P(S,E,L)		
	1	1	1	0.2		
	1	1	0	0.1		
	0	1	1	0.06		
	1	0	0	0.17		
	0	0	0	0.43		

Goal: infer the probability of learning given the disciple wrote feverishly

T

scribbling	excited look	learning	P(S,E,L)
1	1	1	0.2
1	1	0	0.1
0	1	1	0.06
1	0	0	0.17
0	0	0	0.43

Goal: infer the probability of learning given the disciple wrote feverishly



Goal: infer the probability of learning given the disciple wrote feverishly



Goal: infer the probability of learning given the disciple wrote feverishly





Algebra on probabilities?



George Boole, 1815-1864

Propositions: A, B, C

A B — *logical product* or *conjunction*'both of the propositions, A, B are true'

A + B – logical sum or disjunction

'at least one of the propositions, A, B is true'

scribbling	learning	P(S,L)
1	1	0.45
1	0	0.05
0	1	0.15
0	0	0.35

 $\mathsf{P}(l) = \mathsf{P}(s, l) + \mathsf{P}(\overline{s}, l)$

scribbling	learning	P(S,L)
1	1	0.45
1	0	0.05
0	1	0.15
0	0	0.35

Stochastic variable L taking value *l*

 $\mathsf{P}(l) = \mathsf{P}(s, l) + \mathsf{P}(\overline{s}, l)$

scribbling	learning	P(S,L)
1	1	0.45
1	0	0.05
0	1	0.15
0	0	0.35

Stochastic variable L taking value *l*

 $\mathsf{P}(l) = \mathsf{P}(s, l) + \mathsf{P}(\overline{s}, l)$

P the disciple has learned the course material

P she learned and was scribbling

OR P she learned *and* was not scribbling



	scribbling	learning	P(S,L)	Stochastic variable L taking value l
	1	1	0.45	$P(l) = P(s, l) + P(\overline{s}, l)$
/ ' + \	1	0	0.05	P the disciple has learned the course material $P(l) = 0.6$
	0	1	0.15	$P(\bar{l}) = 0.4$
	0	0	0.35	OR
				P she learned and was not scribbling

	scribbling	learning	P(S,L)	Stochastic variable L taking value <i>l</i>
	1	1	0.45	$P(l) = P(s, l) + P(\overline{s}, l)$
	1	0	0.05	P the disciple has learned the course material $P(l) = 0.6$
	0	1	0.15	$P(\bar{l}) = 0.4$
	0	0	0.35	OR
				P she learned and was not scribbling

in general cases when stochastic variables can take more values:

$$P(x) = \sum_{y \in Y} P(x, y)$$

also known as marginalisation

 $\mathsf{P}(s,l) = \mathsf{P}(l \,|\, s)\mathsf{P}(s)$

This turned up at the Cox theorem, being a consequence of the *Consistency* desideratum

 $\mathsf{P}(s,l) = \mathsf{P}(l \,|\, s)\mathsf{P}(s)$

This turned up at the Cox theorem, being a consequence of the *Consistency* desideratum

 $\mathsf{P}(s,l) = \mathsf{P}(l \,|\, s)\mathsf{P}(s)$

P the disciple has scribbled and she learned the course material

P the disciple has learned the course material if she was scribbling AND P she was scribbling

This turned up at the Cox theorem, being a consequence of the *Consistency* desideratum

 $\mathsf{P}(s,l) = \mathsf{P}(l \,|\, s)\mathsf{P}(s)$

P the disciple has scribbled and she learned the course material

P the disciple has learned the course material if she was scribbling AND P she was scribbling

chain rule

Goal: infer the probability of learning given the disciple wrote feverishly

T

scribbling	excited look	learning	P(S,E,L)
1	1	1	0.2
1	1	0	0.1
0	1	1	0.06
1	0	0	0.17
0	0	0	0.43

Goal: infer the probability of learning given the disciple wrote feverishly



Goal: infer the probability of learning given the disciple wrote feverishly



Goal: infer the probability of learning given the disciple wrote feverishly





Goal: infer the probability of learning given the disciple wrote feverishly

 $\mathsf{P}(\mathsf{L}\,|\,s,\overline{e})=?$



 $\mathsf{P}(s, e, l) = \mathsf{P}(l \,|\, s, \overline{e}) \mathsf{P}(s, \overline{e})$

Goal: infer the probability of learning given the disciple wrote feverishly

 $\mathsf{P}(\mathsf{L}\,|\,s,\overline{e})=?$



• Use the chain rule to construct the conditional probability:

$$P(s, e, l) = P(l \mid s, \overline{e})P(s, \overline{e})$$

$$P(l \mid s, \overline{e}) = \frac{P(s, e, l)}{P(s, \overline{e})}$$

Goal: infer the probability of learning given the disciple wrote feverishly



- Use the chain rule to construct the conditional probability:
- Probability of A assuming B is known ('given B'):

$$P(A | B) = \frac{P(A, B)}{P(B)}$$
Computational Cognitive Science

$$P(s, e, l) = P(l \mid s, \overline{e})P(s, \overline{e})$$

$$P(l \mid s, \overline{e}) = \frac{P(s, e, l)}{P(s, \overline{e})}$$
scribbling	learning	P(S,L)
1	1	0.45
1	0	0.05
0	1	0.15
0	0	0.35

$$P(l \mid s) = \frac{\mathsf{P}(s, l)}{\mathsf{P}(s)}$$

	scribbling	learning	P(S,L)
	1	1	0.45
/	1	0	0.05
	0	1	0.15
	0	0	0.35

$$P(l \mid s) = \frac{\mathsf{P}(s, l)}{\mathsf{P}(s)}$$

	scribbling	learning	P(S,L)	
	1	1	0.45	
+	1	0	0.05	
	0	1	0.15	
	0	0	0.35	

$P(1 \mid c)$	 P(s,l)
I(l s)	 P(s)

	scribbling	learning	P(S,L)
	1	1	0.45
+	1	0	0.05
	0	1	0.15
	0	0	0.35

$$P(l \mid s) = \frac{\mathsf{P}(s, l)}{\mathsf{P}(s)} = 0.75$$

	scribbling	learning	P(S,L)
	1	1	0.45
+	1	0	0.05
	0	1	0.15
	0	0	0.35

 $P(l \mid s) = \frac{\mathsf{P}(s, l)}{\mathsf{P}(s)} = 0.75$

By observing that the disciple was writing, the reasoning that she learned the lesson became more plausible

	scribbling	learning	P(S,L)
	1	1	0.45
+	1	0	0.05
	0	1	0.15
	0	0	0.35

$$P(l \mid s) = \frac{\mathsf{P}(s, l)}{\mathsf{P}(s)} = 0.75$$

By observing that the disciple was writing, the reasoning that she learned the lesson became more plausible

 $P(\overline{s} \,|\, l) = ?$

	scribbling	learning	P(S,L)
	1	1	0.45
+	1	0	0.05
	0	1	0.15
	0	0	0.35

$$P(l \mid s) = \frac{\mathsf{P}(s, l)}{\mathsf{P}(s)} = 0.75$$

By observing that the disciple was writing, the reasoning that she learned the lesson became more plausible

 $P(\overline{s} \,|\, l) = ?$

Knowing the probability table rich reasoning can be performed



- Probabilities can be used as a basis of plausible reasoning
- Instead of the truth table, the probability table describes the system
- Rich inferences can be made
- Using the two fundamental rules of probability theory, the sum and the product rules we can express these inferences

Challenge

- Complete characterisation requires filling the probability table
- Number of lines in the probability table increases drastically with the number of variables
- How many parameters? (how many numbers need to be given to specify the probabilistic model?)

Challenge

- Complete characterisation requires filling the probability table
- Number of lines in the probability table increases drastically with the number of variables
- How many parameters? (how many numbers need to be given to specify the probabilistic model?)



 $\mathsf{P}(\mathsf{S},\mathsf{E},\mathsf{R},\mathsf{T},\mathsf{F},\mathsf{L})$

P(S, E, R, T, F, L)

Scribbling Excited look Room where lecture was held Textbook usage Fast talking Learning

P(S, E, R, T, F, L)

Scribbling Excited look Room where lecture was held Textbook usage Fast talking Learning

= P(L | S, E, R, T, F)P(S, E, R, T, F)

Scribbling Excited look Room where lecture was held Textbook usage Fast talking Learning

Chain rule

P(S, E, R, T, F, L)

= P(L|S, E, R, T, F)P(S, E, R, T, F)

Scribbling Excited look Room where lecture was held Textbook usage Fast talking Learning

Chain rule

P(S, E, R, T, F, L)

= P(L|S, E, R, T, F)P(S, E, R, T, F)

Textbook usage is independent of other factors P(A, B) = P(A | B)P(B) = P(A)P(B)

Scribbling Excited look Room where lecture was held Textbook usage Fast talking Learning

Chain rule

P(S, E, R, T, F, L)

= P(L|S, E, R, T, F)P(S, E, R, T, F) = P(L|S, E, R, T, F)P(S, E, R, F)P(T)

Textbook usage is independent of other factors P(A, B) = P(A | B)P(B) = P(A)P(B)







discovering structure in the problem reduces complexity



discovering structure in the problem reduces complexity

finding independencies splits up the probability table into smaller subtables

Graphical models

P(S, E, R, T, F, L)

- Stochastic variables are represented by nodes
- Dependencies (conditional probabilities) are represented by links



Graphical models

$\mathsf{P}(\mathsf{S},\mathsf{E},\mathsf{R},\mathsf{T},\mathsf{F},\mathsf{L})$

Scribbling Excited look Room where lecture was held Textbook usage Fast talking Learning

- Stochastic variables are represented by nodes
- Dependencies (conditional probabilities) are represented by links



deductive reasoning

```
propositional logic
first order logic
λ-calculus
Universal Turing machine
```

expressive power

propositional logic first order logic λ-calculus Universal Turing machine

deductive reasoning

plausible reasoning

expressive power

deductive reasoning

propositional logic first order logic ↓ λ-calculus ↓ Universal Turing machine expressive power

plausible reasoning

probability table





Benefits of a graphical model representation

- helps in visualising the structure of the probabilistic model
- provides insights into the properties of the model (e.g. conditional independence)
- complex calculations for inference and learning can be expressed in graphical terms

Independence:

$$X \perp Y$$
 $P(X) = P(X \mid Y)$
 $P(X, Y) = P(X)P(Y)$



Independence:

 $X \perp Y$ P(X) = P(X | Y)P(X, Y) = P(X)P(Y)

P(X, Y) = P(X)

Conditional independence:



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



Independence:

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

P(X) = P(X | Y)P(X, Y) = P(X)P(Y)

Conditional independence:



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












Computational insight











observed variable: Aristotle's mood





hidden variable: headache

observed variable: Aristotle's mood

- The conditional dependency describes our knowledge on how variables depend on each other
- Based on observation we might want to perform inferences (reason about latent variables)





hidden variable: headache

observed variable: Aristotle's mood

The full model is specified by the joint distribution:

P(O, H) = P(O | H)P(H)

- The conditional dependency describes our knowledge on how variables depend on each other
- Based on observation we might want to perform inferences (reason about latent variables)





hidden variable: headache

observed variable: Aristotle's mood

The full model is specified by the joint distribution:

P(O, H) = P(O | H)P(H)

The full model is specified by the joint distribution:

P(O, H) = P(O | H)P(H)

- The conditional dependency describes our knowledge on how variables depend on each other
- Based on observation we might want to perform inferences (reason about latent variables)





hidden variable: headache

- The conditional dependency describes our knowledge on how variables depend on each other
- Based on observation we might want to perform inferences (reason about latent variables)

The full model is specified by the joint distribution:

```
P(O, H) = P(O | H)P(H)
```

however, we want to reason about the hidden state given the observation: P(H|O)

The full model is specified by the joint distribution:

P(O, H) = P(O | H)P(H)





hidden variable: headache

- The conditional dependency describes our knowledge on how variables depend on each other
- Based on observation we might want to perform inferences (reason about latent variables)

The full model is specified by the joint distribution:

however, we want to reason about the P(O, H) = P(O | H)P(H)hidden state given the observation:

P(H|O)

The full model is specified by the joint distribution:

P(O, H) = P(O | H)P(H) = P(H | O)P(O)





hidden variable: headache

- The conditional dependency describes our knowledge on how variables depend on each other
- Based on observation we might want to perform inferences (reason about latent variables)

The full model is specified by the joint distribution:

however, we want to reason about the P(O, H) = P(O | H)P(H)hidden state given the observation:

P(H|O)

The full model is specified by the joint distribution:

$$P(O, H) = P(O | H)P(H) = P(H | O)P(O)$$



Bayes rule:

$$P(H \mid O) = \frac{P(O \mid H)P(H)}{P(O)}$$





hidden variable: headache

- The conditional dependency describes our knowledge on how variables depend on each other
- Based on observation we might want to perform inferences (reason about latent variables)

The full model is specified by the joint distribution:

```
however, we want to reason about the
P(O, H) = P(O | H)P(H)
                        hidden state given the observation:
                                    P(H|O)
```

The full model is specified by the joint distribution:

posterior probability O, H = P(O|H)P(H) = P(H|O)P(O) $P(H|O) = \frac{P(O|H)P(H)}{P(O)}$ Bayes rule:





hidden variable: headache

• The conditional dependency describes our knowledge on how variables depend on each other

P(H|O)

 Based on observation we might want to perform inferences (reason about latent variables)

The full model is specified by the joint distribution:

```
however, we want to reason about the
P(O, H) = P(O | H)P(H)
                         hidden state given the observation:
```

The full model is specified by the joint distribution:

posterior probability O, H = P(O|H)Plikelihood function $P(H|O) = \frac{P(O|H)P(H)}{P(O)}$ Bayes rule:





hidden variable: headache

- The conditional dependency describes our knowledge on how variables depend on each other
- Based on observation we might want to perform inferences (reason about latent variables)

The full model is specified by the joint distribution:

```
however, we want to reason about the
P(O, H) = P(O | H)P(H)
                        hidden state given the observation:
                                    P(H|O)
```

The full model is specified by the joint distribution:

posterior probability O, H = P(O|H)Plikelihood function $P(H \mid O) = \frac{P(O \mid H)P(H)}{P(O)}$ prior distribution Bayes rule:





hidden variable: headache

- The conditional dependency describes our knowledge on how variables depend on each other
- Based on observation we might want to perform inferences (reason about latent variables)

The full model is specified by the joint distribution:

```
however, we want to reason about the
P(O, H) = P(O | H)P(H)
                        hidden state given the observation:
                                    P(H|O)
```

The full model is specified by the joint distribution:







hidden variable: headache

- The conditional dependency describes our knowledge on how variables depend on each other
- Based on observation we might want to perform inferences (reason about latent variables)

The full model is specified by the joint distribution:

```
however, we want to reason about the
P(O, H) = P(O | H)P(H)
                        hidden state given the observation:
                                    P(H|O)
```

The full model is specified by the joint distribution:







$P(H\,|\,O)$

the model how the world works: how observations are determined by the state of the environment and possibly noise as well



$P(H \,|\, O)$

the model how the world works: how observations are determined by the state of the environment and possibly noise as well



generative process

$P(H \,|\, O)$

the model how the world works: how observations are determined by the state of the environment and possibly noise as well



generative process



the model how the world works: how observations are determined by the state of the environment and possibly noise as well



generative process



the model how the world works: how observations are determined by the state of the environment and possibly noise as well



generative process



recognition model

the model how the world works: how observations are determined by the state of the environment and possibly noise as well



generative process

inverting the generative model: inferences / reasoning about observed data





measurement and inference



measurement and inference

measurement:

measuring the probability of coughing when having a flu or not

measurement and inference

measurement:

measuring the probability of coughing when having a flu or not

P(coughing = 1 | flu = 1)

measurement and inference

measurement:

measuring the probability of coughing when having a flu or not

```
P(\text{coughing} = 1 | \text{flu} = 1)
```

inference: infer the probability of a hypothesis under different conditions

measurement and inference

measurement:

measuring the probability of coughing when having a flu or not

P(coughing = 1 | flu = 1)

inference:

infer the probability of a hypothesis under different conditions

P(flu = 1 | coughing = 1)

measurement and inference

measurement:

measuring the probability of coughing when having a flu or not

P(coughing = 1 | flu = 1)

inference: infer the probability of a hypothesis under different conditions $\frac{P(\text{flu} = 1 | \text{coughing} = 1)}{P(\text{flu} = 1 | \text{coughing} = 1)}$

what is the connection between the two quantities?

measurement and inference

measurement:

measuring the probability of coughing when having a flu or not

P(coughing = 1 | flu = 1)

inference:

infer the probability of a hypothesis under different conditions

P(flu = 1 | coughing = 1)

what is the connection between the two quantities?

remember: **multiplication rule:** P(x, y) = P(x | y)P(y)

measurement and inference

measurement:

measuring the probability of coughing when having a flu or not

P(coughing = 1 | flu = 1)

inference:

infer the probability of a hypothesis under different conditions

P(flu = 1 | coughing = 1)

what is the connection between the two quantities?

remember: **multiplication rule:** P(x, y) = P(x | y)P(y)

or equivalently: P(x, y) = P(y | x)P(x)

measurement and inference

measurement:

measuring the probability of coughing when having a flu or not

P(coughing = 1 | flu = 1)

inference:

infer the probability of a hypothesis under different conditions

P(flu = 1 | coughing = 1)

what is the connection between the two quantities?

remember: multiplication rule: P(x, y) = P(x | y)P(y)

or equivalently:
$$P(x, y) = P(y | x)P(x)$$

 $P(\text{flu} = 1 | \text{coughing} = 1) = \frac{P(\text{coughing} = 1 | \text{flu} = 1)P(\text{flu} = 1)}{P(\text{coughing} = 1)}$

measurement and inference

measurement:

measuring the probability of coughing when having a flu or not

P(coughing = 1 | flu = 1)

inference:

infer the probability of a hypothesis under different conditions

P(flu = 1 | coughing = 1)

what is the connection between the two quantities?

remember: multiplication rule: P(x, y) = P(x | y)P(y)

or equivalently: P(x,y) = P(y | x)P(x)

 $P(\text{flu} = 1 | \text{coughing} = 1) = \frac{P(\text{coughing} = 1 | \text{flu} = 1)P(\text{flu} = 1)}{P(\text{coughing} = 1)}$ **Bayes rule:** 'inverts' a probabilistic relationship

Probability mass function / probability density

Probability mass function / probability density

P(x) summarises all the possible values of a variable in a single function



Probability mass function / probability density

P(x) summarises all the possible values of a variable in a single function

- discrete variables
 - value of the function is listed at all possible values of x
 - the function is the probability mass function


Probability mass function / probability density

P(x) summarises all the possible values of a variable in a single function

- discrete variables
 - value of the function is listed at all possible values of x
 - the function is the probability mass function
- continuous variables
 - the probability of any particular value is equal to 0
 - a finite interval on the variable can have a probability different from 0

_

• the function is a probability density function *p*





Examples of probability distributions

- Discrete valued
 - Bernoulli coin toss

 $Ber(x;p) = p^{x} \cdot (1-p)^{(1-x)}$

- multinomial dice throwing
- Continuous valued
 - uniform
 - Gaussian (normal) $Normal(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$ Gamma







Expectations and uncertainty





Reason #1

limited information

uncertainty in inferences

Formally:

- multiple hypotheses
- alternative hypotheses characterised by different probabilities





🗱 © 1998 Nature America Inc. • http://neurosci.nature.com

Where is the sun?

Jennifer Sun¹ and Pietro Perona^{1,2}

 ¹ California Institute of Technology 136-93, Pasadena, California 91125, USA
² Universita di Padova, Via Ognissanti 72, 35131 Padova, Italy Correspondence should be addressed to P.P. (perona@vision.caltech.edu)

nature *neuroscience* • volume 1 no 3 • july 1998





*





















$P(\text{feature} | \text{stimulus}) \propto P(\text{stimulus} | \text{feature}) \times P(\text{feature})$



 $P(\text{feature} | \text{stimulus}) \propto P(\text{stimulus} | \text{feature}) \times P(\text{feature})$

posterior: inference











Heider & Simmel, 1944



Heider & Simmel, 1944