

Probabilistic and Bayesian Analytics

Note to other teachers and users of these slides. Andrew would be delighted if you found this source material useful in giving your own lectures. Feel free to use these slides verbatim, or to modify them to fit your own needs. PowerPoint originals are available. If you make use of a significant portion of these slides in your own lecture, please include this message, or the following link to the source repository of Andrew's tutorials: <http://www.cs.cmu.edu/~awm/tutorials>. Comments and corrections gratefully received.

Andrew W. Moore
Professor
School of Computer Science
Carnegie Mellon University

www.cs.cmu.edu/~awm
awm@cs.cmu.edu
412-268-7599

Copyright © Andrew W. Moore

Slide 1

Probability

- The world is a very uncertain place
- 30 years of Artificial Intelligence and Database research danced around this fact
- And then a few AI researchers decided to use some ideas from the eighteenth century

Copyright © Andrew W. Moore

Slide 2

What we're going to do

- We will review the fundamentals of probability.
- It's really going to be worth it
- In this lecture, you'll see an example of probabilistic analytics in action: Bayes Classifiers

Copyright © Andrew W. Moore

Slide 3

Discrete Random Variables

- A is a Boolean-valued random variable if A denotes an event, and there is some degree of uncertainty as to whether A occurs.
- Examples
 - A = The US president in 2023 will be male
 - A = You wake up tomorrow with a headache
 - A = You have Ebola

Copyright © Andrew W. Moore

Slide 4

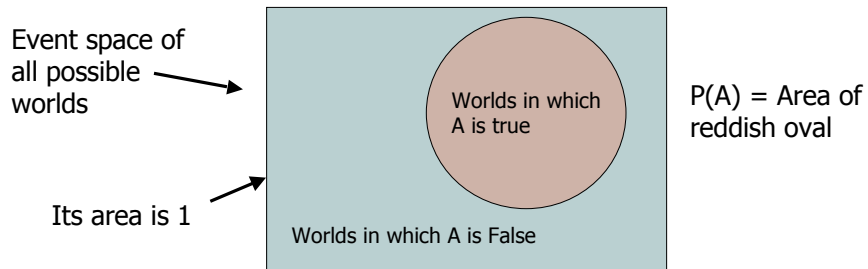
Probabilities

- We write $P(A)$ as “the fraction of possible worlds in which A is true”
- We could at this point spend 2 hours on the philosophy of this.
- But we won't.

Copyright © Andrew W. Moore

Slide 5

Visualizing A



**You can think of a probability as the *measurement* of a set.
In this case, the measurement is *area*.**

Copyright © Andrew W. Moore

Slide 6



The Axioms of Probability

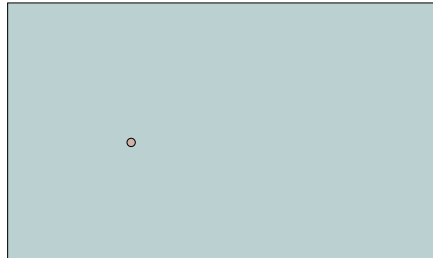
- $0 \leq P(A) \leq 1$
- $P(\text{True}) = 1$
- $P(\text{False}) = 0$
- $P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$

- These plus set theory are all you need

Where do these axioms come from? Were they "discovered"?
Answers coming up later.

Interpreting the axioms

- $0 \leq P(A) \leq 1$
- $P(\text{True}) = 1$
- $P(\text{False}) = 0$
- $P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$



The area of A can't get any smaller than 0

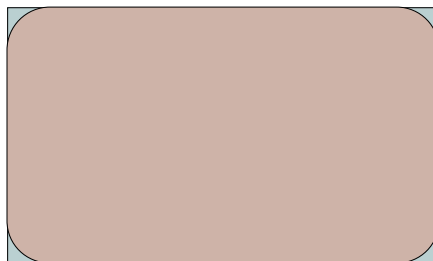
And a zero area would mean no world could ever have A true

Copyright © Andrew W. Moore

Slide 9

Interpreting the axioms

- $0 \leq P(A) \leq 1$
- $P(\text{True}) = 1$
- $P(\text{False}) = 0$
- $P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$



The area of A can't get any bigger than 1

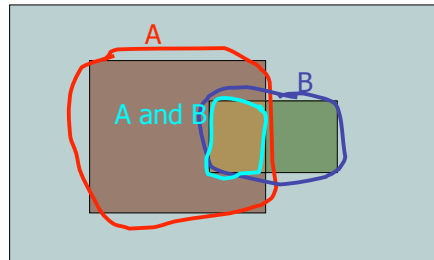
And an area of 1 would mean all worlds will have A true

Copyright © Andrew W. Moore

Slide 10

Interpreting the axioms

- $0 \leq P(A) \leq 1$
- $P(\text{True}) = 1$
- $P(\text{False}) = 0$
- $P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$



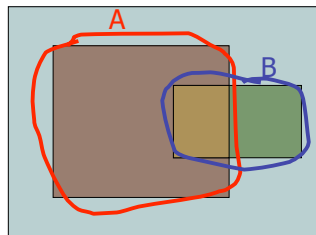
Copyright © Andrew W. Moore

Slide 11

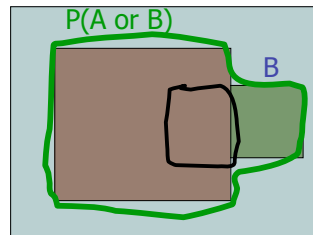
Interpreting the axioms

- $0 \leq P(A) \leq 1$
- $P(\text{True}) = 1$
- $P(\text{False}) = 0$
- $P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$

$P(A \text{ and } B)$ would be counted twice



Simple addition and subtraction



Copyright © Andrew W. Moore

Slide 12

These Axioms are Not to be Trifled With

- There have been attempts to do different methodologies for uncertainty
 - Fuzzy Logic
 - Three-valued logic
 - Dempster-Shafer
 - Non-monotonic reasoning
- But the axioms of probability are the only system with this property:
If you gamble using them you can't be unfairly exploited by an opponent using some other system [di Finetti 1931]

Copyright © Andrew W. Moore

Slide 13

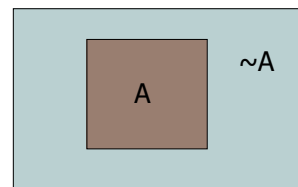
Theorems from the Axioms

- $0 \leq P(A) \leq 1$, $P(\text{True}) = 1$, $P(\text{False}) = 0$
- $P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$

From these we can prove:

$$P(\text{not } A) = P(\sim A) = 1 - P(A)$$

- How?
 - Substitute $B = \sim A$ above
 - Notice $P(A \text{ and } \sim A) = 0$
 - Notice $P(A \text{ or } \sim A) = 1$



Copyright © Andrew W. Moore

Slide 14

Side Note

- I am inflicting these proofs on you for two reasons:
 1. These kind of manipulations will need to be second nature to you if you use probabilistic analytics in depth
 1. (or even not in depth -- Dan)
 2. Suffering is good for you

Copyright © Andrew W. Moore

Slide 15

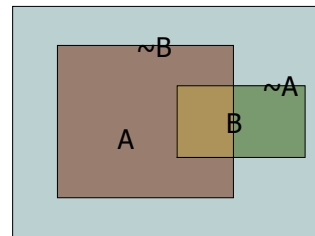
Another important theorem

- $0 \leq P(A) \leq 1$, $P(\text{True}) = 1$, $P(\text{False}) = 0$
- $P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$

From these we can prove:

$$P(A) = P(A \wedge B) + P(A \wedge \sim B)$$

- How?
 - One way:
 - Let $C = A \wedge B$, $D = A \wedge \sim B$
 - Substitute above
 - Notice $P(C \wedge D) = 0$
 - Notice $P(C \vee D) = P(A)$



Copyright © Andrew W. Moore

Slide 16

Multivalued Random Variables

- Suppose A can take on more than 2 values
- A is a *random variable with arity k* if it can take on exactly one value out of $\{v_1, v_2, \dots, v_k\}$
- Thus...

$$P(A = v_i \wedge A = v_j) = 0 \text{ if } i \neq j$$

$$P(A = v_1 \vee A = v_2 \vee \dots \vee A = v_k) = 1$$

Copyright © Andrew W. Moore

Slide 17

An easy fact about Multivalued Random Variables:

- Using the axioms of probability...
 $0 \leq P(A) \leq 1$, $P(\text{True}) = 1$, $P(\text{False}) = 0$
 $P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$

- And assuming that A obeys...

$$P(A = v_i \wedge A = v_j) = 0 \text{ if } i \neq j$$

$$P(A = v_1 \vee A = v_2 \vee \dots \vee A = v_k) = 1$$

- It's easy to prove that

$$P(A = v_1 \vee A = v_2 \vee \dots \vee A = v_i) = \sum_{j=1}^i P(A = v_j)$$

Copyright © Andrew W. Moore

Slide 18

An easy fact about Multivalued Random Variables:

- Using the axioms of probability...

$$0 \leq P(A) \leq 1, P(\text{True}) = 1, P(\text{False}) = 0$$

$$P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$$

- And assuming that A obeys...

$$P(A = v_i \wedge A = v_j) = 0 \text{ if } i \neq j$$

$$P(A = v_1 \vee A = v_2 \vee \dots \vee A = v_k) = 1$$

- It's easy to prove that

$$P(A = v_1 \vee A = v_2 \vee \dots \vee A = v_i) = \sum_{j=1}^i P(A = v_j)$$

- And thus we can prove

$$\sum_{j=1}^k P(A = v_j) = 1$$

Copyright © Andrew W. Moore

Slide 19

Another fact about Multivalued Random Variables:

- Using the axioms of probability...

$$0 \leq P(A) \leq 1, P(\text{True}) = 1, P(\text{False}) = 0$$

$$P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$$

- And assuming that A obeys...

$$P(A = v_i \wedge A = v_j) = 0 \text{ if } i \neq j$$

$$P(A = v_1 \vee A = v_2 \vee \dots \vee A = v_k) = 1$$

- It's easy to prove that

$$P(B \wedge [A = v_1 \vee A = v_2 \vee \dots \vee A = v_i]) = \sum_{j=1}^i P(B \wedge A = v_j)$$

Copyright © Andrew W. Moore

Slide 20

Another fact about Multivalued Random Variables:

- Using the axioms of probability...

$$0 \leq P(A) \leq 1, P(\text{True}) = 1, P(\text{False}) = 0$$

$$P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$$

- And assuming that A obeys...

$$P(A = v_i \wedge A = v_j) = 0 \text{ if } i \neq j$$

$$P(A = v_1 \vee A = v_2 \vee \dots \vee A = v_k) = 1$$

- It's easy to prove that

$$P(B \wedge [A = v_1 \vee A = v_2 \vee \dots \vee A = v_i]) = \sum_{j=1}^i P(B \wedge A = v_j)$$

- And thus we can prove

$$P(B) = \sum_{j=1}^k P(B \wedge A = v_j)$$

Copyright © Andrew W. Moore

Slide 21

Elementary Probability in Pictures

$$\sum_{j=1}^k P(A = v_j) = 1$$

Copyright © Andrew W. Moore

Slide 22

Elementary Probability in Pictures

$$P(B) = \sum_{j=1}^k P(B \wedge A = v_j)$$

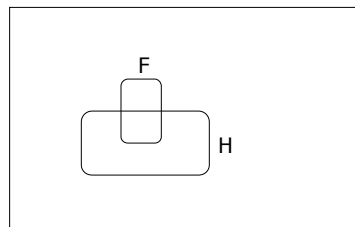
Copyright © Andrew W. Moore

Slide 23

Conditional Probability

- $P(A|B)$ = Fraction of worlds in which B is true that also have A true

H = "Have a headache"
F = "Coming down with Flu"



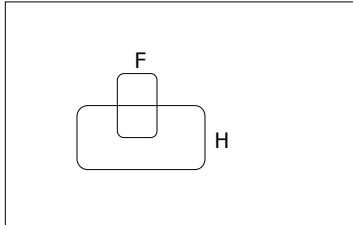
$P(H) = 1/10$
 $P(F) = 1/40$
 $P(H|F) = 1/2$

"Headaches are rare and flu is rarer, but if you're coming down with 'flu there's a 50-50 chance you'll have a headache."

Copyright © Andrew W. Moore

Slide 24

Conditional Probability



H = "Have a headache"
 F = "Coming down with Flu"

$P(H) = 1/10$
 $P(F) = 1/40$
 $P(H|F) = 1/2$

$P(H|F)$ = Fraction of flu-inflicted worlds in which you have a headache

= $\frac{\text{\#worlds with flu and headache}}{\text{\#worlds with flu}}$

= $\frac{\text{Area of "H and F" region}}{\text{Area of "F" region}}$

= $\frac{P(H \wedge F)}{P(F)}$

Copyright © Andrew W. Moore

Slide 25

Definition of Conditional Probability

$$P(A|B) = \frac{P(A \wedge B)}{P(B)}$$

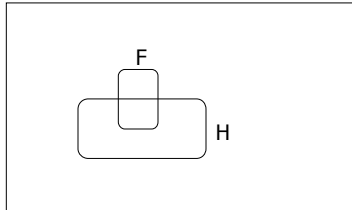
Corollary: The Chain Rule

$$P(A \wedge B) = P(A|B) P(B)$$

Copyright © Andrew W. Moore

Slide 26

Probabilistic Inference



H = "Have a headache"
F = "Coming down with Flu"

$$P(H) = 1/10$$
$$P(F) = 1/40$$
$$P(H|F) = 1/2$$

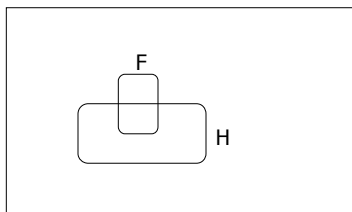
One day you wake up with a headache. You think: "Drat! 50% of flus are associated with headaches so I must have a 50-50 chance of coming down with flu"

Is this reasoning good?

Copyright © Andrew W. Moore

Slide 27

Probabilistic Inference



H = "Have a headache"
F = "Coming down with Flu"

$$P(H) = 1/10$$
$$P(F) = 1/40$$
$$P(H|F) = 1/2$$

$$P(F \wedge H) = \dots$$

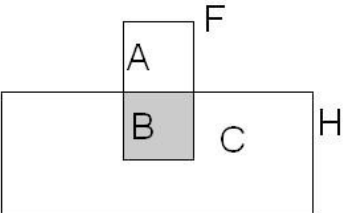
$$P(F|H) = \dots$$

Copyright © Andrew W. Moore

Slide 28

Another way to understand the intuition

Thanks to Jahanzeb Sherwani for contributing this explanation:



Let's say we have $P(F)$, $P(H)$, and $P(H|F)$, like in the example in class.

Area-wise, $P(F) = A + B$, $P(H) = B + C$,

Also, $P(H|F) = \frac{B}{A + B}$

Thus, to get the opposite conditional probability, ie, $P(F|H)$, we need to figure out $\frac{B}{B + C}$

Since we know $B / (A+B)$, we can get $B / (B+C)$ by multiplying by $(A+B)$ and dividing by $(B+C)$. But since we already calculated, $A+B = P(F)$, and $B+C = P(H)$, so we are actually multiplying by $P(F)$ and dividing by $P(H)$. Which is Bayes Rule:

$$P(F|H) = P(H|F) * \frac{P(F)}{P(H)}$$

Copyright © Andrew W. Moore

Slide 29

What we just did...

$$P(B|A) = \frac{P(A \wedge B)}{P(A)} = \frac{P(A|B) P(B)}{P(A)}$$

This is Bayes Rule

Bayes, Thomas (1763) An essay towards solving a problem in the doctrine of chances. *Philosophical Transactions of the Royal Society of London*, **53:370-418**



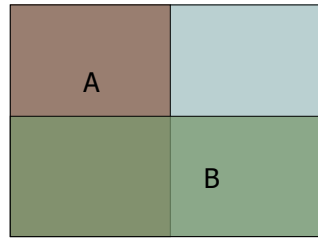
Copyright © Andrew W. Moore

Slide 30

Aside: Fixing my mistakes

“Independent”

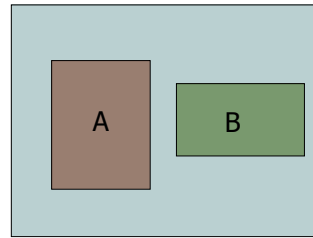
$$P(A \text{ and } B) = P(A) \cdot P(B)$$



Special case of
 $P(A \text{ and } B) = P(A) \cdot P(B|A)$

“Mutually Exclusive”

$$P(A \text{ or } B) = P(A) + P(B)$$



Special case of
 $P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$

These are not the same thing!

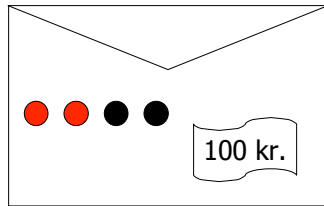
Copyright © Andrew W. Moore

Slide 31

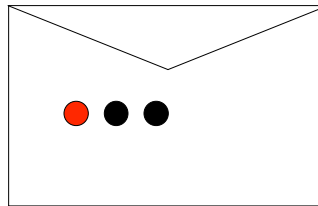
Copyright © Andrew W. Moore

Slide 32

Using Bayes Rule to Gamble



The "Win" envelope has 100 krónur and four beads in it



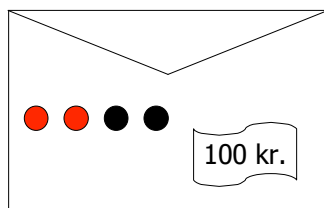
The "Lose" envelope has three beads and no money

Trivial question: someone draws an envelope at random (i.e. 50% chance of winning) and offers to sell it to you. How much should you pay?

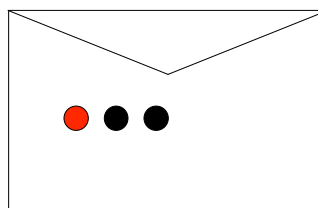
Copyright © Andrew W. Moore

Slide 33

Using Bayes Rule to Gamble



The "Win" envelope has 100 krónur and four beads in it



The "Lose" envelope has three beads and no money

Interesting question: before deciding, you are allowed to see one bead drawn from the envelope.

Suppose it's black: How much should you pay?

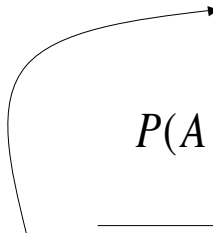
Suppose it's red: How much should you pay?

Copyright © Andrew W. Moore

Slide 34

More General Forms of Bayes Rule

$$P(A|B) = \frac{P(B|A)P(A)}{P(B|A)P(A) + P(B|\sim A)P(\sim A)}$$


$$P(A|B \wedge X) = \frac{P(B|A \wedge X)P(A \wedge X)}{P(B \wedge X)}$$

Here, the denominator just computes $P(B)$.

In the gambling example, we could have done:

$$P(\text{black}) = P(\text{black}|\text{win}) \cdot P(\text{win}) + P(\text{black}|\sim\text{win}) \cdot P(\sim\text{win})$$

Copyright © Andrew W. Moore

Slide 35

More General Forms of Bayes Rule

$$P(A = v_i | B) = \frac{P(B | A = v_i)P(A = v_i)}{\sum_{k=1}^{n_A} P(B | A = v_k)P(A = v_k)}$$

Copyright © Andrew W. Moore

Slide 36

Useful Easy-to-prove facts

$$P(A | B) + P(\sim A | B) = 1$$

$$\sum_{k=1}^{n_A} P(A = v_k | B) = 1$$

Implied by the definition of conditional probability

$$P(A|B) = \frac{P(A \wedge B)}{P(B)}$$

...and the "marginalization" formula

$$P(A) = P(A \wedge B) + P(A \wedge \sim B)$$

Copyright © Andrew W. Moore

Slide 37

The Joint Distribution

Example: Boolean variables A, B, C

Recipe for making a joint distribution
of M variables:

Copyright © Andrew W. Moore

Slide 38

The Joint Distribution

Example: Boolean variables A, B, C

Recipe for making a joint distribution of M variables:

1. Make a truth table listing all combinations of values of your variables (if there are M Boolean variables then the table will have 2^M rows).

A	B	C
0	0	0
0	0	1
0	1	0
0	1	1
1	0	0
1	0	1
1	1	0
1	1	1

Copyright © Andrew W. Moore

Slide 39

The Joint Distribution

Example: Boolean variables A, B, C

Recipe for making a joint distribution of M variables:

1. Make a truth table listing all combinations of values of your variables (if there are M Boolean variables then the table will have 2^M rows).
2. For each combination of values, say how probable it is.

A	B	C	Prob
0	0	0	0.30
0	0	1	0.05
0	1	0	0.10
0	1	1	0.05
1	0	0	0.05
1	0	1	0.10
1	1	0	0.25
1	1	1	0.10

Copyright © Andrew W. Moore

Slide 40

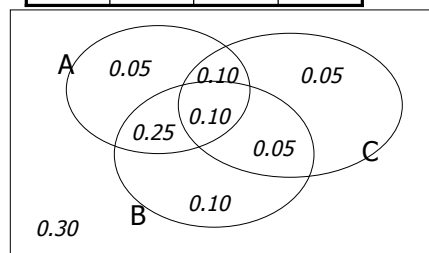
The Joint Distribution

Example: Boolean variables A, B, C

Recipe for making a joint distribution of M variables:

1. Make a truth table listing all combinations of values of your variables (if there are M Boolean variables then the table will have 2^M rows).
2. For each combination of values, say how probable it is.
3. If you subscribe to the axioms of probability, those numbers must sum to 1.

A	B	C	Prob
0	0	0	0.30
0	0	1	0.05
0	1	0	0.10
0	1	1	0.05
1	0	0	0.05
1	0	1	0.10
1	1	0	0.25
1	1	1	0.10



Copyright © Andrew W. Moore

Slide 41

Using the Joint

gender	hours_worked	wealth	prob
Female	v0:40.5-	poor	0.253122
		rich	0.0245895
	v1:40.5+	poor	0.0421768
		rich	0.0116293
Male	v0:40.5-	poor	0.331313
		rich	0.0971295
	v1:40.5+	poor	0.134106
		rich	0.105933

One you have the Joint Distribution, you can ask for the probability of any logical expression involving your attributes

$$P(E) = \sum_{\text{rows matching } E} P(\text{row})$$

Copyright © Andrew W. Moore

Slide 42

Using the Joint

gender	hours_worked	wealth	
Female	v0:40.5-	poor	0.253122
		rich	0.0245895
	v1:40.5+	poor	0.0421768
		rich	0.0116293
Male	v0:40.5-	poor	0.331313
		rich	0.0971295
	v1:40.5+	poor	0.134106
		rich	0.105933

$$P(\text{Poor} \wedge \text{Male}) = 0.4654$$

$$P(E) = \sum_{\text{rows matching } E} P(\text{row})$$

Copyright © Andrew W. Moore

Slide 43

Using the Joint

gender	hours_worked	wealth	
Female	v0:40.5-	poor	0.253122
		rich	0.0245895
	v1:40.5+	poor	0.0421768
		rich	0.0116293
Male	v0:40.5-	poor	0.331313
		rich	0.0971295
	v1:40.5+	poor	0.134106
		rich	0.105933

$$P(\text{Poor}) = 0.7604$$

$$P(E) = \sum_{\text{rows matching } E} P(\text{row})$$

Copyright © Andrew W. Moore

Slide 44

Inference with the Joint

gender	hours_worked	wealth	
Female	v0:40.5-	poor	0.253122
		rich	0.0245895
	v1:40.5+	poor	0.0421768
		rich	0.0116293
Male	v0:40.5-	poor	0.331313
		rich	0.0971295
	v1:40.5+	poor	0.134106
		rich	0.105933

$$P(E_1 | E_2) = \frac{P(E_1 \wedge E_2)}{P(E_2)} = \frac{\sum_{\text{rows matching } E_1 \text{ and } E_2} P(\text{row})}{\sum_{\text{rows matching } E_2} P(\text{row})}$$

Copyright © Andrew W. Moore

Slide 45

Inference with the Joint

gender	hours_worked	wealth	
Female	v0:40.5-	poor	0.253122
		rich	0.0245895
	v1:40.5+	poor	0.0421768
		rich	0.0116293
Male	v0:40.5-	poor	0.331313
		rich	0.0971295
	v1:40.5+	poor	0.134106
		rich	0.105933









$$P(E_1 | E_2) = \frac{P(E_1 \wedge E_2)}{P(E_2)} = \frac{\sum_{\text{rows matching } E_1 \text{ and } E_2} P(\text{row})}{\sum_{\text{rows matching } E_2} P(\text{row})}$$

$$P(\text{Male} | \text{Poor}) = 0.4654 / 0.7604 = 0.612$$

Copyright © Andrew W. Moore

Slide 46

Inference with the Joint

gender	hours_worked	wealth	
Female	v0:40.5-	poor	0.253122 
		rich	0.0245895 
	v1:40.5+	poor	0.0421768 
		rich	0.0116293 
Male	v0:40.5-	poor	0.331313 
		rich	0.0971295 
	v1:40.5+	poor	0.134106 
		rich	0.105933 

$$P(E_1 | E_2) = \frac{P(E_1 \wedge E_2)}{P(E_2)} = \frac{\sum_{\text{rows matching } E_1 \text{ and } E_2} P(\text{row})}{\sum_{\text{rows matching } E_2} P(\text{row})}$$

$$P(\text{Poor} | \text{Male}) = 0.4654 / 0.6685 = 0.696$$

Copyright © Andrew W. Moore

Slide 47

Inference is a big deal

- I've got this evidence. What's the chance that this conclusion is true?
 - I've got a sore neck: how likely am I to have meningitis?
 - I see my lights are out and it's 9pm. What's the chance my spouse is already asleep?

Copyright © Andrew W. Moore

Slide 48

Inference is a big deal

- I've got this evidence. What's the chance that this conclusion is true?
 - I've got a sore neck: how likely am I to have meningitis?
 - I see my lights are out and it's 9pm. What's the chance my spouse is already asleep?
- There's a thriving set of industries growing based around Bayesian Inference. Highlights are: Medicine, Pharma, Help Desk Support, Engine Fault Diagnosis

Copyright © Andrew W. Moore

Slide 49

Where do Joint Distributions come from?

- Idea One: Expert Humans
- Idea Two: Simpler probabilistic facts and some algebra (axioms, Bayes rule, ...)

Example: Suppose you had three binary variables A, B, C, and you knew

$$\begin{array}{ll} P(A) = 0.7 & P(C|A \wedge B) = 0.1 \\ & P(C|A \wedge \sim B) = 0.8 \\ P(B|A) = 0.2 & P(C|\sim A \wedge B) = 0.3 \\ P(B|\sim A) = 0.1 & P(C|\sim A \wedge \sim B) = 0.1 \end{array}$$

Then you can automatically compute the JD using the chain rule

$$P(A=x \wedge B=y \wedge C=z) = P(C=z|A=x \wedge B=y) P(B=y|A=x) P(A=x)$$

Once you have recovered JD, you can ask whatever you want!

Copyright © Andrew W. Moore

Slide 50

Where do Joint Distributions come from?

- Idea Three: Learn them from data!

Prepare to see one of the most impressive learning algorithms you'll come across in the entire course....

Learning a joint distribution

Build a JD table for your attributes in which the probabilities are unspecified

A	B	C	Prob
0	0	0	?
0	0	1	?
0	1	0	?
0	1	1	?
1	0	0	?
1	0	1	?
1	1	0	?
1	1	1	?

The fill in each row with









$$\hat{P}(\text{row}) = \frac{\text{records matching row}}{\text{total number of records}}$$

A	B	C	Prob
0	0	0	0.30
0	0	1	0.05
0	1	0	0.10
0	1	1	0.05
1	0	0	0.05
1	0	1	0.10
1	1	0	0.25
1	1	1	0.10

Fraction of all records in which A and B are True but C is False

Example of Learning a Joint

- This Joint was obtained by learning from three attributes in the UCI "Adult" Census Database [Kohavi 1995]

gender	hours_worked	wealth	
Female	v0:40.5-	poor	0.253122 
		rich	0.0245895 
	v1:40.5+	poor	0.0421768 
		rich	0.0116293 
Male	v0:40.5-	poor	0.331313 
		rich	0.0971295 
	v1:40.5+	poor	0.134106 
		rich	0.105933 

Copyright © Andrew W. Moore

Slide 53

Where are we?

- We have recalled the fundamentals of probability
- We have become content with what Joint Distributions are and how to use them
- And we even know how to learn Joint Distributions from data.

Copyright © Andrew W. Moore

Slide 54