
More probability

CS311
David Kauchak
Spring 2013

*Some material borrowed from:
Sara Owsley Sood and others*

Admin

- Assign 3 Tuesday at the beginning of class (in class)
- Should have looked at written 2 by now
- Written 3 out soon

- Mancala tournament: good news and bad news

Another example

Start with the joint probability distribution:

	toothache		\neg toothache	
	catch	\neg catch	catch	\neg catch
cavity	.108	.012	.072	.008
\neg cavity	.016	.064	.144	.576

$P(\text{toothache}) = ?$

Another example

Start with the joint probability distribution:

	toothache		\neg toothache	
	catch	\neg catch	catch	\neg catch
cavity	.108	.012	.072	.008
\neg cavity	.016	.064	.144	.576

$P(\text{toothache}) = 0.108 + 0.012 + 0.016 + 0.064 = 0.2$

Another example

Start with the joint probability distribution:

	toothache		¬toothache	
	catch	¬catch	catch	¬catch
cavity	.108	.012	.072	.008
¬cavity	.016	.064	.144	.576

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

Another example

Start with the joint probability distribution:

	toothache		¬toothache	
	catch	¬catch	catch	¬catch
cavity	.108	.012	.072	.008
¬cavity	.016	.064	.144	.576

$$\begin{aligned}
 P(\neg\text{cavity} \mid \text{toothache}) &= \frac{P(\neg\text{cavity}, \text{toothache})}{P(\text{toothache})} \\
 &= \frac{0.016 + 0.064}{0.108 + 0.012 + 0.016 + 0.064} \\
 &= 0.4
 \end{aligned}$$

Normalization

	toothache		¬toothache	
	catch	¬catch	catch	¬catch
cavity	.108	.012	.072	.008
¬cavity	.016	.064	.144	.576

Denominator can be viewed as a **normalization constant** α

$$\begin{aligned}
 P(\text{CAVITY} \mid \text{toothache}) &= \alpha P(\text{CAVITY}, \text{toothache}) \\
 &= \alpha [P(\text{CAVITY}, \text{toothache}, \text{catch}) + P(\text{CAVITY}, \text{toothache}, \neg\text{catch})] \\
 &= \alpha \langle 0.108, 0.012 \rangle + \langle 0.072, 0.008 \rangle \\
 &= \alpha \langle 0.12, 0.08 \rangle = \langle 0.6, 0.4 \rangle
 \end{aligned}$$

↑ unnormalized $p(\text{cavity} \mid \text{toothache})$ ↓ unnormalized $p(\neg\text{cavity} \mid \text{toothache})$

General idea: compute distribution on query variable by fixing **evidence variables** and summing over **hidden/unknown variables**

More Probability

In the United States, 55% of children get an allowance and 41% of children get an allowance and do household chores. What is the probability that a child does household chores given that the child gets an allowance?

$$\begin{aligned}
 p(\text{chores} \mid \text{allow}) &= p(\text{chores}, \text{allow}) / p(\text{allow}) \\
 &= 0.41 / 0.55 = 0.745
 \end{aligned}$$

Still more probability

- A math teacher gave her class two tests. 25% of the class passed both tests and 42% of the class passed the first test. What is the probability that a student who passed the first test also passed the second test?

Another Example

A patient takes a lab test and the result comes back positive. The test has a false negative rate of 2% and false positive rate of 2%. Furthermore, 0.5% of the entire population have this cancer.

What is the probability of cancer if we know the test result is positive?

Another Example

A patient takes a lab test and the result comes back positive. The test has a false negative rate of 2% and false positive rate of 2%. Furthermore, 0.5% of the entire population have this cancer.

What is the probability of cancer if we know the test result is positive?

$p(\text{cancer}) = 0.005$ false negative: negative result even though we have cancer
 $p(\text{false_neg}) = 0.02$
 $p(\text{false_pos}) = 0.02$ false positive: positive result even though we don't have cancer
 $p(\text{cancer} | \text{pos}) = ?$

Another Example

$p(\text{cancer}) = 0.005$ false negative: negative result even though we have cancer
 $p(\text{false_neg}) = 0.02$
 $p(\text{false_pos}) = 0.02$ false positive: positive result even though we don't have cancer
 $p(\text{cancer} | \text{pos}) = ?$

$$p(\text{cancer} | \text{pos}) = \frac{p(\text{cancer}, \text{pos})}{p(\text{pos})}$$

Another Example

$p(\text{cancer}) = 0.005$
 $p(\text{false_neg}) = 0.02$
 $p(\text{false_pos}) = 0.02$
 $p(\text{cancer} | \text{pos}) = ?$

false negative: negative result even though we have cancer

false positive: positive result even though we don't have cancer

$$\frac{p(\text{cancer, pos})}{p(\text{pos})} = \frac{p(\text{cancer})(1 - p(\text{false_neg}))}{p(\text{cancer})(1 - p(\text{false_neg})) + p(\neg\text{cancer})p(\text{false_pos})}$$

1-p(false_neg) gives us the probability of the test correctly identifying us with cancer

two ways to get a positive result: cancer with a correct positive and not cancer with a false positive

Another Example

$p(\text{cancer}) = 0.005$
 $p(\text{false_neg}) = 0.02$
 $p(\text{false_pos}) = 0.02$
 $p(\text{cancer} | \text{pos}) = ?$

false negative: negative result even though we have cancer

false positive: positive result even though we don't have cancer

$$p(\text{cancer} | \text{pos}) = 0.1975$$

Contrast this with $p(\text{pos} | \text{cancer}) = 0.98$

Obtaining probabilities



We've talked a lot about probabilities, but not where they come from

- intuition/guess
 - this can be very hard
 - people are not good at this for anything but the simplest problems
- estimate from data!

Estimating probabilities



Total Flips: 10
Number Heads: 5
Number Tails: 5

Probability of Heads:
Number Heads / Total Flips = 0.5

Probability of Tails:
Number Tails / Total Flips = 0.5 = 1.0 - Probability of Heads

The experiments, the sample space and the events must be defined clearly for probability to be meaningful

Theoretical Probability

Maximum entropy principle

- When one has only partial information about the possible outcomes one should choose the probabilities so as to maximize the uncertainty about the missing information
- Alternatives are always to be judged equally probable if we have no reason to expect or prefer one over the other

Maximum likelihood estimation

- set the probabilities so that we maximize how likely our data is

Turns out these approaches do the same thing!

Maximum Likelihood Estimation

Number of times an event occurs in the data

Total number of times experiment was run
(total number of data collected)

Maximum Likelihood Estimation

Number of times an event occurs in the data

Total number of times experiment was run
(total number of data collected)

[Rock/Paper/Scissors](http://www.nytimes.com/interactive/science/rock-paper-scissors.html)

<http://www.nytimes.com/interactive/science/rock-paper-scissors.html>

How is it done?

Maximum Likelihood Estimation

Number of times an event occurs in the data

Total number of times experiment was run
(total number of data collected)

[Rock/Paper/Scissors](http://www.nytimes.com/interactive/science/rock-paper-scissors.html)

<http://www.nytimes.com/interactive/science/rock-paper-scissors.html>



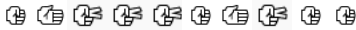
- Analyze the prior choices
- Select probability of next choice based on data

How?

Maximum Likelihood Estimation

Number of times an event occurs in the data

Total number of times experiment was run
(total number of data collected)



$$P(\text{rock}) =$$

$$P(\text{rock} \mid \text{scissors}) =$$

$$P(\text{rock} \mid \text{scissors, scissors, scissors}) =$$

Maximum Likelihood Estimation

Number of times an event occurs in the data

Total number of times experiment was run
(total number of data collected)



$$P(\text{rock}) = 4/10 = 0.4$$

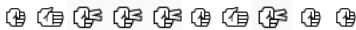
$$P(\text{rock} \mid \text{scissors}) = 2/4 = 0.5$$

$$P(\text{rock} \mid \text{scissors, scissors, scissors}) = 1/1 = 1.0$$

Maximum Likelihood Estimation

Number of times an event occurs in the data

Total number of times experiment was run
(total number of data collected)



$$P(\text{rock}) = 4/10 = 0.4$$

$$P(\text{rock} \mid \text{scissors}) = 2/4 = 0.5$$

$$P(\text{rock} \mid \text{scissors, scissors, scissors}) = 1/1 = 1.0$$

Which of these do you think is most accurate?

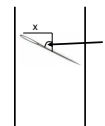
Law of Large Numbers

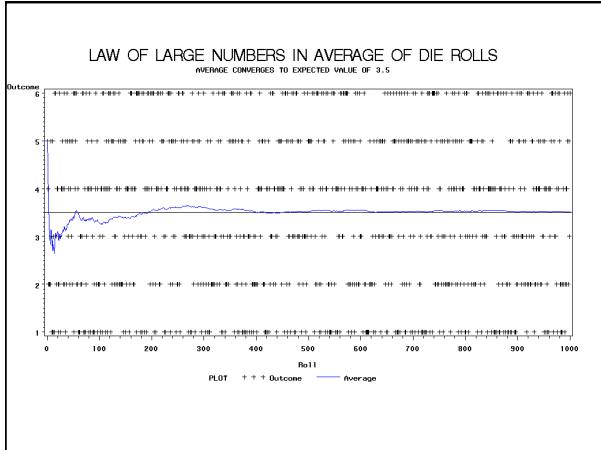
As the number of experiments increases the relative frequency of an event more closely approximates the actual probability of the event.

- if the theoretical assumptions hold

Buffon's Needle for Computing π

- <http://mste.illinois.edu/reese/buffon/buffon.html>






Large Numbers Reveal Problems in Assumptions

Results of 1,000,000 throws of a die

Number	1	2	3	4	5	6
Fraction	.155	.159	.164	.169	.174	.179

Probabilistic Reasoning



Evidence

- What we know about a situation


Hypothesis

- What we want to conclude

Compute

- $P(\text{Hypothesis} | \text{Evidence})$

Probabilistic Reasoning



Evidence

- What we know about a situation

Hypothesis

- What we want to conclude

Compute

- $P(\text{Hypothesis} | \text{Evidence})$

Credit card application?

Credit Card Application

E is the data about the applicant's age, job, education, income, credit history, etc,

H is the hypothesis that the credit card will provide positive return.

The decision of whether to issue the credit card to the applicant is based on the probability $P(H|E)$.

Probabilistic Reasoning



Evidence

- What we know about a situation

Hypothesis

- What we want to conclude

Compute

- $P(\text{Hypothesis} | \text{Evidence})$

Medical diagnosis?

Medical Diagnosis

E is a set of symptoms, such as, coughing, sneezing, headache, ...

H is a disorder, e.g., common cold, SARS, swine flu.

The diagnosis problem is to find an H (disorder) such that $P(H|E)$ is maximum.

Chain rule (aka product rule)

$$p(X|Y) = \frac{P(X,Y)}{P(Y)} \quad \Rightarrow \quad p(X,Y) = P(X|Y)P(Y)$$

We can view calculating the probability of X AND Y occurring as two steps:

1. Y occurs with some probability $P(Y)$
2. Then, X occurs, given that Y has occurred

or you can just trust the math... 😊

Chain rule (aka product rule)

$$p(X|Y) = \frac{P(X,Y)}{P(Y)} \quad \Rightarrow \quad p(X,Y) = P(Y|X)P(X)$$

We can view calculating the probability of X AND Y occurring as two steps:

1. X occurs with some probability P(X)
2. Then, Y occurs, given that X has occurred

or you can just trust the math... ☺

Chain rule

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

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

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

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

$$p(X_1, X_2, \dots, X_n) = ?$$

Bayes' rule (theorem)

$$p(X|Y) = \frac{P(X,Y)}{P(Y)} \quad \Rightarrow \quad p(X,Y) = P(X|Y)P(Y)$$

$$p(X|Y) = \frac{P(X,Y)}{P(Y)} \quad \Rightarrow \quad p(X,Y) = P(Y|X)P(X)$$

$$p(X|Y) = \frac{P(Y|X)P(X)}{P(Y)}$$

Bayes' rule

Allows us to talk about P(Y|X) rather than P(X|Y)

Sometimes this can be more intuitive

Why?

$$p(X|Y) = \frac{P(Y|X)P(X)}{P(Y)}$$

Bayes' rule

$p(\text{disease} \mid \text{symptoms})$

- For everyone who had those symptoms, how many had the disease?
- $p(\text{symptoms} \mid \text{disease})$
 - For everyone that had the disease, how many had this symptom?

$p(\text{good_lender} \mid \text{credit_features})$

- For everyone who had these credit features, how many were good lenders?
- $p(\text{credit_features} \mid \text{good_lender})$
 - For all the good lenders, how many had this feature

$p(\text{cause} \mid \text{effect})$ vs. $p(\text{effect} \mid \text{cause})$

$p(H \mid E)$ vs. $p(E \mid H)$

Bayes' rule

$$p(\text{good_lender} \mid \text{features}) = \frac{P(\text{features} \mid \text{good_lender})P(\text{good_lender})}{P(\text{features})}$$

We often already have data on good lenders, so $p(\text{features} \mid \text{good_lender})$ is straightforward

$p(\text{features})$ and $p(\text{good_lender})$ are often easier than $p(\text{good_lender} \mid \text{features})$

Allows us to properly handle changes in just the underlying distribution of good_lenders, etc.

Other benefits

Simple lender model:

- score: is credit score > 600
- debt: debt < income

$$p(\text{Good} \mid \text{Credit, Debt}) = \frac{P(\text{Credit, Debt} \mid \text{Good})P(\text{Good})}{P(\text{Credit, Debt})}$$

Other benefits

It's in the 1950s and you train your model "diagnostically" using just $p(\text{Good} \mid \text{Credit, Debt})$.

However, in the 1960s and 70s the population of people that are good lenders drastically increases (baby-boomers learned from their depression era parents and are better with their money)

$$p(\text{Good} \mid \text{Credit, Debt})$$

Intuitively what should happen?

Other benefits

It's in the 1950s and you train your model "diagnostically" using just $p(\text{Good} \mid \text{Credit}, \text{Debt})$.

However, in the 1960s and 70s the population of people that are good lendees drastically increases (baby-boomers learned from their depression era parents and are better with their money)

$$p(\text{Good} \mid \text{Credit}, \text{Debt})$$

Probability of "good" should increase, but that's hard to figure out from just this equation

Other benefits

$$p(\text{Good} \mid \text{Credit}, \text{Debt}) = \frac{P(\text{Credit}, \text{Debt} \mid \text{Good})P(\text{Good})}{P(\text{Credit}, \text{Debt})}$$

Modeled using Bayes' rule, it's clear how much the probability should change.

Measure what the new $P(\text{Good})$ is.

When it rains...

Marie is getting married tomorrow at an outdoor ceremony in the desert. In recent years, it has rained only 5 days each year. Unfortunately, the weatherman has predicted rain for tomorrow. When it actually rains, the weatherman correctly forecasts rain 90% of the time. When it doesn't rain, he incorrectly forecasts rain 5% of the time. What is the probability that it will rain on the day of Marie's wedding?

$$\begin{aligned} p(\text{rain}) &= 5/365 \\ p(\text{predicted} \mid \text{rain}) &= 0.9 \\ p(\text{predicted} \mid \neg \text{rain}) &= 0.05 \end{aligned}$$

When it rains...

$$\begin{aligned} p(\text{rain}) &= 5/365 \\ p(\text{predicted} \mid \text{rain}) &= 0.9 \\ p(\text{predicted} \mid \neg \text{rain}) &= 0.05 \end{aligned}$$

$$\begin{aligned} p(\text{rain} \mid \text{predicted}) &= \frac{p(\text{predicted} \mid \text{rain})p(\text{rain})}{p(\text{predicted})} \\ &= \frac{0.9 * 5 / 365}{p(\text{predicted})} \end{aligned}$$

When it rains...

$$p(\text{rain}) = 5/365$$

$$p(\text{predicted}|\text{rain}) = 0.9$$

$$p(\text{predicted}|\sim\text{rain}) = 0.05$$

$$p(\text{predicted}) = p(\text{predicted}|\text{rain})p(\text{rain}) + p(\text{predicted}|\sim\text{rain})p(\sim\text{rain})$$

$$p(\sim\text{rain}|\text{predicted}) = \frac{p(\text{predicted}|\sim\text{rain})p(\sim\text{rain})}{p(\text{predicted})}$$

$$= 0.05 * 360/365$$

Monty Hall

- 3 doors
 - behind two, something bad
 - behind one, something good
- You pick one door, but are not shown the contents
- Host opens one of the other two doors that has the bad thing behind it (he always opens one with the bad thing)
- You can now switch your door to the other unopened. Should you?



Monty Hall

$p(\text{win})$ initially?

- 3 doors, 1 with a winner, $p(\text{win}) = 1/3$

$p(\text{win} | \text{shown_other_door})?$

- One reasoning:
 - once you're shown one door, there are just two remaining doors
 - one of which has the winning prize
 - $1/2$

This is not correct!

Be careful! – Player picks door 1

	winning location	host opens
1/3	Door 1	1/2 Door 2 1/2 Door 3
1/3	Door 2	1 Door 3
1/3	Door 3	1 Door 2

In these two cases, switching will give you the correct answer.
Key: host knows where it is.

Another view

1000 doors

- behind 999, something bad
- behind one, something good

You pick one door, but are not shown the contents

Host opens 998 of the other 999 doors that have the bad thing behind it
(he always opens ones with the bad thing)

In essence, you're picking between it being behind your one door or
behind any one of the other doors (whether that be 2 or 999)

