

$P(Y_1 = 1 | B)$ Bayesians sometimes refer ⁽¹²⁾ to this as the annoying normalizing constant, because it's often the hardest ingredient in Bayes's Theorem to pin down. I will

now show you 3 different ways to estimate this probability, all of them

instructive: Method 1: 2x2 contingency table
(8 Jan 19)

Imagine randomly ^{blood} sampling 100,000 people from the population $P = \{ \text{all people similar to Bob in all relevant ways} \}$ (i.e., ^{people with the} same B) and running all the blood samples through "Determine HIV"; cross-tabulate (truth) against (what blood test says), by which I mean what we expect to happen:

Truth (really is)

		HIV+	HIV-	
test says (+)	(1)	(2)	(3)	(*)
test says (-)	(4)	(5)	(6)	
	(7)	(8)	100,000 = (9)	

what we know (13)

so far:

I $P(\text{really is HIV+} | \mathcal{B}) = .01$

II $P(\text{test says HIV+} | \text{really is HIV+}) = .999$

Fact I is called the prevalence of HIV in Bob's population \mathcal{P} ; let's use it first.

People call (*) a 2 by 2 contingency table; the counts (integers) in the body of the table (1, 2, 4, 5) record how many blood samples have the indicated row attribute (what the test says) and the indicated column attribute (the truth) (e.g.; 4 = # blood samples for which test says - and truth is HIV+). The counts

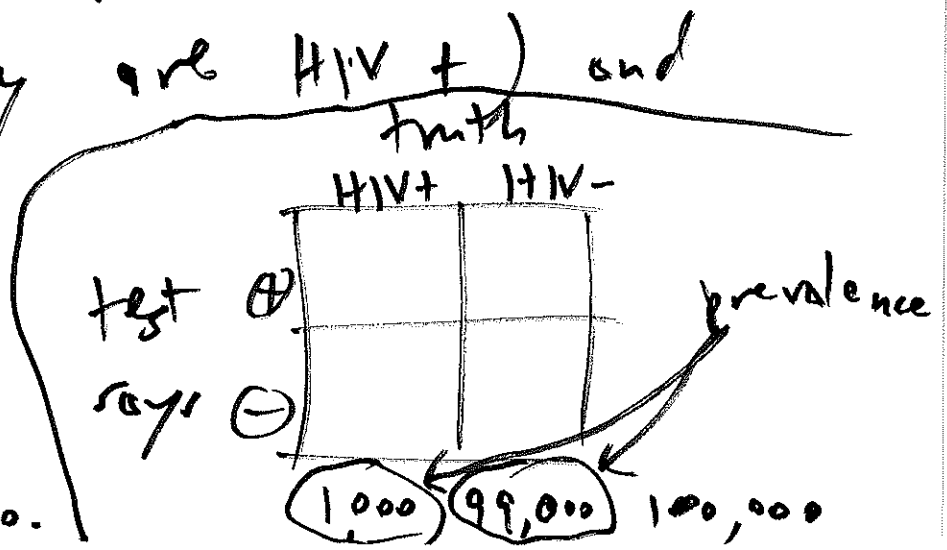
at the margins of the table (3, 6, 5, 14, 7, 14),
 record how many blood samples have the
 indicated row (or column) attribute,
 adding up over the column (or row)
 attribute, so that $1 + 2 = 3$, $4 + 10 = 14$,
 $1 + 4 = 5$, $2 + 5 = 7$ (e.g. 6 = # blood
 samples for which test says \ominus).

Q_i

How can we use the prevalence 1 to
 begin filling in the table? A_i 1% of
 100,000 is 1,000, so $7 = 1,000$ (because
 the first column represents all the blood
 samples that really are HIV+) and

Therefore 8
 has to be

$$100,000 - 1,000 = 99,000.$$



Q2: How can we bring the sensitivity to bear on the table!

A2: well,

sensitivity = P(test says HIV+ | really is HIV+) = .999;

Q3: Where are we in the table when we know that a blood sample really is HIV+?

In the first column.

So, in 1,000 truly HIV+ samples, the test gets it right

(0.999) * 1000 = 999 times, & only makes 1 mistake, leading to this table.

	truth		
	HIV+	HIV-	
test says (+)	999		Sensitivity
test says (-)	1		
	1,000	99,000	100,000

Q4: How can we fill in the rest of the table?

A4: we need

new information on how good the test is:

The missing analogue to sensitivity is called the specificity P(test says HIV- | really is HIV-)

It turns out that "Determine HIV"'s specificity is almost as good as its sensitivity:

$$P(\text{test says HIV-} \mid \text{really is HIV-}) = 0.994$$

Now we can

complete the table, by using the specificity on the second column:

$$(0.994) \cdot (99,000) = 98,406$$

truly HIV- samples are correctly marked as HIV-

by the test, leaving $(99,000 - 98,406) = 594$ mistakes;

and the row totals are now available & complete the table.

		truth		
		HIV+	HIV-	
test says	+	999	594	1,593
	-	1	98,406	98,407
		1,000	99,000	100,000

specificity

We can now compute $P(\theta = 1 \mid y_1 = 1, B)$, and

the answer is surprising: [given that the test says HIV+ ($y_1 = 1$)] puts us in the

first row of the table

so $P(\theta = 1 \mid y_1 = 1, B)$ disappointingly low?

$$= \frac{999}{1,593} = 0.627$$

sensitivity

specificity

$$\text{sensitivity} = 0.999, \text{ specificity} = 0.994$$

A related & also relevant probability: (15)

$$P(\theta=0 | \gamma_1=1, B) = \text{false positive rate}$$

$$= \frac{594}{1593} = 0.373 = 1 - P(\theta=1 | \gamma_1=1, B)$$

With sensitivity & specificity values so close to 1, a false positive rate of 37% is disappointingly high; why has this happened?

Q5 Let's look at $P(\theta=1 | \gamma_1=1, B) = 63\%$

Q6 When is a fraction small? A6 When its

numerator is small, or its denominator is large, or both. $63\% = \frac{999}{1,593}$ (9)

Given the low prevalence in Bob's population γ , 999 is almost as large as it possibly could be; its biggest possible value is only $999+1 = 1,000$

b) The denominator, 1,593, = 999 + 594, (18)
 got large in 2 ways: the terrific sensitivity
 yielded 999, and the prevalence was so low
 that the second column total was 99,000;

		Truth		
		HIV+	HIV-	
test + says -	+	999	594	1,593
	-	1	98,406	98,407
		1,000	99,000	100,000

even with a specificity
 of 0.994, $(0.006)(99,000)$
 = 594 is big. Note

that this test's false negative

rate is $P(\theta = 1 | \gamma_1 = 0, \beta) = \frac{1}{98,407} \approx 0.0001$,

a fantastically low value. A7

Evidently the
 developers of this blood test were much more
 worried about false negatives than false positives;

is this sensible?

A7 Answering this involves

two additional questions: Q8 The test will be
 used for what purpose? Q9 what are the

real-world consequences of each type of error? (19)

A₈ In CS1 the test will be used by You (the doctor) to advise your patient Bob on his HIV status (other uses would include screening blood donated to a blood bank).

Note: something interesting

just happened: before Q₅, we were doing science (the acquisition of knowledge for its own sake); starting with Q₅, we're thinking about decision

theory (using scientific knowledge to make a ^{behavioral} choice: what should ^(as Bob's doctor) you do if "Determine HIV" comes back positive?)

A₉ False positive = (test says ⊕ but really ⊖)

This is bad: if your treatment plan for Bob (20) is based only on the "Determine HIV" result, you will now begin treating Bob for HIV when in fact this treatment is unnecessary, & Bob will think he has a serious disease when in fact he doesn't.

False negative =
(test says \ominus but really \oplus)

This is also bad, but for a different reason:

You will tell Bob he's HIV negative when he's actually HIV positive; he will go untreated, & he may unintentionally infect other people.

We can't make any additional progress without making a value judgment: Q90 which type of bad is worse, or are they equally bad?

A90 Statisticians & economists have developed the concept of a utility function.

to make value judgments like this one ⁽²⁾

precise: $u_{\text{Bob}}(a, \theta^* | B) = \text{the numerical value}$
utility function

Bob places on action a if the unknown

θ has the value θ^* . This requires specifying

$(A|B)_{\text{Bob}} = \{ \text{all feasible actions judged by Bob to be relevant (worthy of consideration)} \}$
← the action space

Suppose (for example)

that Bob chooses to put his treatment plan completely in your hands; then $(A_{\text{Bob}} | B) = (A_{\text{you}} | B)$

Suppose further (for

example) that you choose $(A_{\text{you}} | B) = \{ \begin{matrix} \text{+} \\ \text{-} \end{matrix} \} = \{ \text{treat Bob for HIV if "determine HIV" returns a positive result \& prescribe no treatment for Bob if test result is } \ominus \}$

$u_{Bob}(a, \theta | B)$ truth
 $\theta=1 \quad \theta=0$

test (+)	u_{11}	u_{12}
says (-)	u_{21}	u_{22}

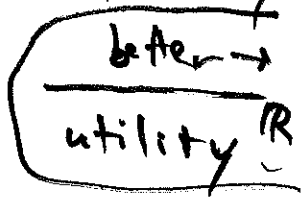
Now, to approach this formally, You would need to elicit the u_{ij} values from Bob (since You're

acting as his ^(proxy) agent, so to speak).

A simple

set of reasonable u_{ij} values can be obtained by (without loss of generality)

(a) adopting the convention that large utility values are preferred over small values



and (b) regarding utility as negative loss:

$U_{Bob}(a, \theta^* | B) = -L_{Bob}(a, \theta^* | B)$, where

$L_{Bob}(a, \theta^* | B)$ is the loss Bob incurs if action a is chosen & the unknown θ takes the value θ^* .

From this point of view

the diagonal entries in the table above are arguably 0, yielding the new table

$u_{Bob}(a, \theta | B)$ truth
 $\theta=1 \quad \theta=0$

test (+)	0	$-l_{12}$
says (-)	$-l_{21}$	0

(more subtle choices are possible)

$l_{12} > 0$ is the loss incurred by Bob if a false positive occurs, and $l_{21} > 0$ is Bob's loss under a false negative. It helps to set

$l_{12} = L, l_{21} = c \cdot L$ for some $c > 0$; then Bob only needs to think about c . Almost everyone

would agree that false negatives are (a lot) worse than false positives here, so that $c \gg 1$. Fact:

all tests that try to detect HIV attempt to estimate a quantity g

		truth	
		$\theta = 1$	$\theta = 0$
test says	\oplus	0	$-L$ (fake \oplus)
	\ominus	$-c \cdot L$ (false \ominus)	0

(e.g., concentration of antibodies against the virus) & use a decision rule of the form

if $g \leq T$ report HIV \ominus | if $g \geq T$ report HIV \oplus

T is a tuning constant that

trades off false ⊕ against false ⊖ errors: (24)
 as $T \uparrow$ it gets harder to report HIV ⊕, which
 lowers the false ⊕ rate but increases the
 false ⊖ rate, and vice versa as $T \downarrow$.

As You
(continued)

		HIV+	HIV-	
test	⊕	999	594	1,593
says	⊖	1	98,406	98,407
		1,000	99,000	100,000

can see that
 qualitatively the
 designers of
 "Determine HIV"
 have their utility function aligned with
 Bob's: the test is terrified of false
 negatives, & inevitably as a result it
 makes a lot of false positive mistakes
 (i.e., the designers aren't stupid).

creative laziness; Bayes's Theorem in
 odds form) Method 2

odds form) CS1

$$P(\theta=1 | \gamma_1=1, \mathcal{B}) = \frac{P(\theta=1 | \mathcal{B}) P(\gamma_1=1 | \theta=1, \mathcal{B})}{P(\gamma_1=1 | \mathcal{B})}$$

The denominator in the right-hand side (25) is annoying; ^{the} shall we cope with this? Advice

If solving a problem directly is difficult, try to be creatively lazy, by solving it indirectly. Application of this advice:

The opposite possibility to $P(\theta=1 | y_1=1, \mathcal{B})$ is $P(\theta=0 | y_1=1, \mathcal{B})$, and crucially this has the same denominator; take the ratio of these two probabilities and the annoying denominator will cancel:

$$P(\theta=1 | y_1=1, \mathcal{B}) = \frac{P(\theta=1 | \mathcal{B}) P(y_1=1 | \theta=1, \mathcal{B})}{P(y_1=1 | \mathcal{B})}$$

$$P(\theta=0 | y_1=1, \mathcal{B}) = \frac{P(\theta=0 | \mathcal{B}) P(y_1=1 | \theta=0, \mathcal{B})}{P(y_1=1 | \mathcal{B})}$$

Divide these two equations to obtain (26)
 (likelihood ratio)

$$\left[\frac{P(\theta=1 | Y_1=1, B)}{P(\theta=0 | Y_1=1, B)} \right] \textcircled{*} = \left[\frac{P(\theta=1 | B)}{P(\theta=0 | B)} \right] \cdot \left[\frac{P(Y_1=1 | \theta=1, B)}{P(Y_1=1 | \theta=0, B)} \right]$$

Definition The conditional odds ratio

in favor of A given B is $O_{A|B} = \frac{P(A|B)}{P(\text{not } A|B)}$

With this definition, $\textcircled{*}$ becomes

$$\left(\begin{array}{l} \text{posterior odds} \\ \text{in favor of} \\ \theta=1 \text{ given} \\ Y_1=1 \ \& \ B \end{array} \right) = \left(\begin{array}{l} \text{prior odds} \\ \text{in favor} \\ \text{of } \theta=1 \\ \text{given } B \end{array} \right) \cdot \left(\begin{array}{l} \text{data odds} \\ \text{Bayes} \\ \text{factor} \\ \text{in favor} \\ \text{of } \theta=1 \\ \text{given } Y_1=1 \ \& \ B \end{array} \right)$$

(Bayes's Theorem in odds form)

Note that we know (via the prevalence, sensitivity & specificity) all of the probabilities on the right-hand side of $\textcircled{*}$.

$$\left(\begin{array}{l} \text{prior odds in favor} \\ \text{of } \theta=1 \text{ given } B \end{array} \right) = \frac{.01}{.99} = 99 \text{ to } 1 \text{ odds against } \theta=1$$

Def. The conditional odds ratio against (27)

$$\text{A given B is } O_{A|B}^{\text{not}} = \frac{P(\text{not } A | B)}{P(A | B)} = \frac{1}{O_{A|B}}$$

The sensitivity gives us the numerator of the Bayes factor: $P(Y_1=1 | \theta=1, X) = 0.999$

And the denominator is just 1 minus the specificity: $P(Y_1=1 | \theta=0, X) = 1 - P(Y_1=0 | \theta=0) = 1 - 0.994 = 0.006$.

So the Bayes factor is

$$\frac{P(Y_1=1 | \theta=1)}{P(Y_1=1 | \theta=0)} = \frac{0.999}{0.006} = \frac{333}{2} = 166.5 \text{ to } 1$$

in favor of $\theta=1$

The prior information pushes away from $(\theta=1)$ with a 99 to 1 strength on the odds scale, but the data information pushes toward $(\theta=1)$ with a 166.5 to 1 strength on the same scale; Bayes's Theorem tells

us that the optimal way to combine these information sources is multiplicatively

the odds scale:

$$\frac{P(\theta=1 | \gamma_1=1, B)}{P(\theta=0 | \gamma_1=1, B)} = \left[\frac{P(\theta=1 | B)}{P(\theta=0 | B)} \right] \cdot \left[\frac{P(\gamma_1=1 | \theta=1, B)}{P(\gamma_1=1 | \theta=0, B)} \right]$$

So the posterior odds, in favor of

$$= \frac{1}{99} \cdot \frac{333}{2} = \frac{37}{22}$$

Bob being HIV+, given the blood test results and Bob's background information, is $\frac{37}{22}$.

Q: How are odds ratios and probabilities related?

$$O_{A|B} = \frac{P(A|B)}{P(\text{not } A|B)} = \frac{P_{A|B}}{1 - P_{A|B}}$$

Solve this for

$P_{A|B}$ in terms of $O_{A|B}$:

$$P_{A|B} = \frac{O_{A|B}}{1 + O_{A|B}}$$

Thus the posterior probability that Bob is HIV+, given the test results & Bob's background information,

$$P(\theta=1 | \gamma_1=1, B) = \frac{\frac{37}{22}}{1 + \frac{37}{22}} = \frac{37}{59} = \frac{999}{1593} = 0.627$$