

Method 3) (the hard way) ("extending the conversation") CS1 we know $P(\theta=1 | B)$,

$P(Y_1=1 | \theta=1)$ and $P(Y_1=0 | \theta=0)$;
 ← or put blood test information in \mathcal{B}

we want $P(\theta=1 | Y_1=1, B) = \frac{P(\theta=1 | B) P(Y_1=1 | \theta=1, B)}{P(Y_1=1 | B)}$

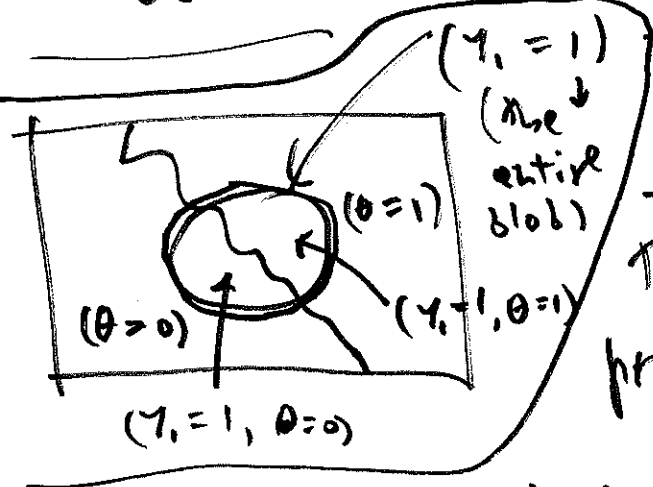
Both of the numerator probabilities are known; but what about the annoying denominator? Q: 13

^{can we} How get $P(Y_1=1 | B)$ from $P(\theta=1 | B)$, $P(Y_1=1 | \theta=1)$ and $P(Y_1=0 | \theta=0)$?

A: Notice that $P(Y_1=1 | B)$ is hard, but $P(Y_1=1 | \theta=1)$ and $P(Y_1=0 | \theta=0)$ are easy; in other words, we don't know how the data (the blood test) will come out, but we do know how the data will come out ^{probabilistically} if we knew the truth $\{P(Y_1=1 | \theta=1), P(Y_1=0 | \theta=0)\}$. So let's extend the conversation by bringing θ into (J.V. Lindley: 1923-2013)

the picture. Since the only two possibilities for the truth are $(\theta=1)$ and $(\theta=0)$, those two propositions form a partition: a collection of mutually exclusive possibilities that is exhaustive of all the possibilities.

Step 1



$$P(\gamma_1=1 | \mathcal{B}) = P(\gamma_1=1, \theta=0 | \mathcal{B}) + P(\gamma_1=1, \theta=1 | \mathcal{B})$$

This is progress (θ is now in the conversation), but more work is needed.

Step 2 θ is on the wrong side of the conditioning bar in $P(\gamma_1=1, \theta=0 | \mathcal{B})$ to be useful to us, so let's force it to move to the other side: as before,

$$P(\gamma_1=1, \theta=0 | \mathcal{B}) = P(\gamma_1=1 | \theta=0, \mathcal{B}) \cdot [?]$$

doesn't so why were useful, so let's try $P(\gamma_1=1, \theta=0 | \mathcal{B}) = P(\gamma_1=1 | \theta=0, \mathcal{B}) \cdot [?]$

$$P(Y_1=1, \theta=0 | B) = P(Y_1=1 | \theta=0, B) \cdot \boxed{?}$$

" (definition of conditional probability) "

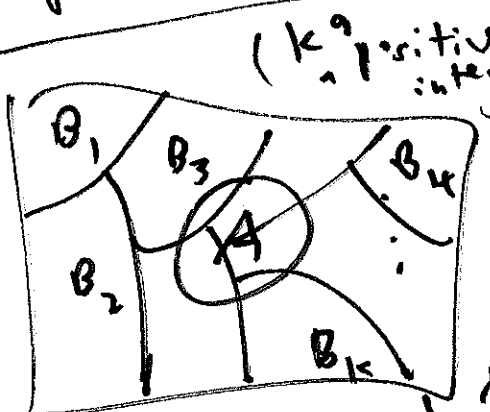
$$\frac{P(Y_1=1, \theta=0, B)}{P(B)} = \frac{P(Y_1=1, \theta=0, B)}{P(\theta=0, B)} \cdot \boxed{\frac{P(\theta=0, B)}{P(B)}}$$

So it works, and the answer is

$$P(Y_1=1, \theta=0 | B) = P(Y_1=1 | \theta=0, B) \cdot \boxed{P(\theta=0 | B)}$$

Thus $P(Y_1=1 | B) = P(\theta=0 | B) \cdot P(Y_1=1 | \theta=0, B) + P(\theta=1 | B) \cdot P(Y_1=1 | \theta=1, B)$

This is a special case of the Law of Total Probability.



(k a positive integer)

If $\{B_1, \dots, B_k\}$ form a partition of {all possibilities},

then $P(A) = \sum_{i=1}^k P(B_i) P(A | B_i)$

($P(B_i) > 0$ by definition of a partition)

We've just computed the annoying denominator by partitioning over the truth.

$$P(\gamma_1 = 1 | B) = P(\theta = 0 | B) \cdot P(\gamma_1 = 1 | \theta = 0, B) + P(\theta = 1 | B) \cdot P(\gamma_1 = 1 | \theta = 1, B)$$

$$= (0.99) [1 - P(\gamma_1 = 0 | \theta = 0, B)]$$

prevalence \rightarrow (0.99) specificity \rightarrow $[1 - P(\gamma_1 = 0 | \theta = 0, B)]$ sensitivity \leftarrow (0.999)

$$+ (0.01) \cdot (0.999)$$

$$= \frac{99}{100} \cdot \frac{6}{1000} + \frac{1}{100} \cdot \frac{999}{1000} = \frac{1593}{100000} \text{ and}$$

finally $P(\theta = 1 | \gamma_1 = 1, B) = \frac{\frac{1}{100} \cdot \frac{999}{100000}}{\frac{1593}{100000}} = \frac{999}{1593} \approx 0.63$

Extending the conversation to include the unknown θ , which in applications of Bayesian learning amounts to partitioning over the truth, is a powerful technique that will come up a number of times in what follows; Bayes's Theorem in odds form is also highly useful.

The Big Picture, $P = (Q, C) \rightarrow (\theta, D, B)$ ③
Revisited

① I can make the identification of Q and C from P unique by adopting the convention that if you, given P , have a different Q and/or C in mind, you're working on a different problem than I am.

(unfortunately) the mapping from (Q, C) to (θ, D, B) is not necessarily unique.

In CS1, (this will typically be true), for example, θ and D are uniquely specified, but different reasonable choices of B are possible:

to obtain $P(\theta=1 | B)$ from the medical literature, it was necessary to specify $P = \{ \text{all U.S. adults similar to Bob in all relevant ways} \}$;

I chose $B_1 = (\text{male})$, $B_2 = (\text{age } 28)$, $B_3 = (\text{gay})$, $B_4 = (\text{mostly safe sex})$, but it would be reasonable to also consider $B_5 = (\text{multiple partners})$ and $B_6 =$

(last tested Θ 11 months ago) if ^{population} data were available on these variables as well.

The next step

in the model-building was to go from (θ, D, B) to $\{p(\theta|B), p(D|\theta B)\}$

In CS1 $p(\theta|B) = \int p(\theta=1|B)$

and $p(D|\theta B) = \begin{cases} p(\gamma_1=1 | \theta=1, B) \\ 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{cases}$
this is ^{called} the sampling distribution ("sample")

$\begin{cases} p(\theta=0|B) \\ \uparrow \\ \text{this is} \\ \text{the } \underline{\text{prior}} \\ \text{information} \end{cases}$

because it specifies how the data $D=(\gamma_1)$ is likely to come out if θ were known; in CS2 we'll see how the sampling distribution is converted into the likelihood information (distribution)

$l(\theta|D B)$. So the general paradigm, ^{so far} is

$$P = (Q, C) \rightarrow (\theta, D, B) \rightarrow M_{IP} \stackrel{\Delta}{=} \{p(\theta|B), p(D|\theta B)\}$$

Here IP stands for Inference (drawing conclusions about the unknown θ) and Prediction (estimating

new (not yet observed) data values D^*). / If ⁽³⁵⁾

You need to go beyond science (inference, prediction) to make a choice (decision-making).
The paradigm becomes will concentrate on inference & prediction in this course

$P = (Q, C) + (\theta, D, B) + \Pi_{IPD} = \{ p(\theta|B), p(D|\theta B),$

Statistical data science $\{ (A|B), u(q, \theta|B) \}$
 $(A|C)$ (data wrangling)

encompasses 4 main activities: ① Description
(data curation)
(graphical and numerical) of existing data sets

D; ② inference (drawing conclusions about the
(θ)
underlying scientific process, that gave rise
to the data); ③ prediction (drawing inferences
about new data D^*); and ④ decision (choosing
the best action [because you have to make a
choice] even though you have uncertainty
about relevant quantities (θ)). Good data science

almost always begins with a (possibly ⁽³⁶⁾ extensive) graphical & numerical descriptive exploration of the data set D , with a particular focus on missing data; more on this toward the end of the class.

In CS4 we had uncertainty about how to specify the prior information $p(\theta|B)$; this will often be true in real-world applications.

In CS1 we didn't have any uncertainty about the sampling distribution $p(D|\theta, B)$ (because the sensitivity & specificity of the blood test were "known"); in ~~more~~ complicated problems, even there, the 0.999 \pm ? sensitivity & 0.994 \pm ? specificity were estimated from data. You will typically

also have uncertainty about $p(D|\theta, B)$.

Thus in general you will have 2 levels of uncertainty: You're uncertain about (level 1)

θ , but you're also uncertain about (level 2) B , (how to specify your uncertainty about θ , through B , $p(\theta|B)$, and $p(D|\theta, B)$).

ie., the mapping from (θ, D, B) to $M_{IP} =$

$\{p(\theta|B), p(D|\theta, B)\}$ is also typically

not unique. Level 2 uncertainty is called (naturally enough) model uncertainty;

it has been systematically studied in

detail since the 1990s. I will offer ^{some} advice on how to cope with model uncertainty.

In this course (Draper (1995))

Advice: A simple approach to assessing the magnitude of model uncertainty is sensitivity analysis: Vary the aspect of the modeling about which you're uncertain across a plausible range & see how much difference it makes to ^{the} results that you care the most about.

CS 1: Let α = prevalence

β = sensitivity γ = specificity

	truth		
	HIV +	HIV -	
blood test (+)	$\alpha\beta$	$(1-\alpha)(1-\gamma)$	$\alpha\beta + (1-\alpha)(1-\gamma)$
(-)	$\alpha(1-\beta)$	$(1-\alpha)\gamma$	
	α	$1-\alpha$	1

symbolically the false positive rate

is $FPR = \frac{(1-\alpha)(1-\gamma)}{\alpha\beta + (1-\alpha)(1-\gamma)}$

and the false negative

rate is $FNR = \frac{\alpha(1-\beta)}{\alpha(1-\beta) + \gamma(1-\alpha)}$

Now we

can (e.g.) hold β at 0.999 and γ at $\textcircled{39}$ 0.994 and vary $\alpha = P(\theta=1 | \beta)$ from (say) 0.005 to 0.02 (a factor of 2 lower & higher than the previous value of 0.01).

$(\beta = 0.999, \gamma = 0.994)$		
α	FPR	FNR
0.005	0.544	0.00000506
0.01	0.373	0.0000102
0.02	0.227	0.0000205

cutting the prevalence in half increases FPR by $\left| \frac{.544 - .373}{.373} \right|$.
 $100\% = 46\%$; doubling

the prevalence decreases FPR by $\left| \frac{.227 - .373}{.373} \right|$.
 $100\% = 39\%$. (15 Jun 19)

(i.e., the false positive rate is quite sensitive to prevalence (prior information))

Cutting α in half almost exactly cuts the FNR in half, and doubling α almost exactly doubles FNR, so the false negative rate is also quite sensitive to prevalence (although its value remains extremely low).