

Basic Bayes

Bayes rule (or theorem, law) and its application is a huge subject. Here we only scratch the surface.

The subject is *conditional* probability, where our interest is in whether something, an event (E), is true *given* that we know something else (F) is true or has happened.

$$p(E|F)$$

For example, you may take a rapid test for cancer C and receive a positive result $+$. The question you will ask is what is the probability that I have cancer *given* that positive result.

$$p(C|+)$$

A very common response is to give the probability that a person with cancer will have a positive result $p(+|C)$. But this is not what we asked, and although it is relevant, it cannot give us the whole of what we want to know.

table

The usual characterization of a test is by its error rates. There are two types of error: first is a false negative, FN , which for this example is the occurrence of a negative result when cancer is present. The second is a false positive, FP , a positive result when the subject does not actually have cancer.

The other possible results are sometimes called TP (true positive) and TN (true negative). Let's make a table:

	C	no C
+	TP	FP
-	FN	TN

Two descriptors are used for these error rates. Unfortunately, they are easily confused. *Sensitivity* is the ability to correctly identify persons *with* disease. High sensitivity means a low rate of FN . A very sensitive test is good at finding cases.

Specificity is the ability to correctly identify persons *without* disease. High specificity means a low rate of FP . A very specific test does not give many false positives. Both qualities are obviously desirable.

For example, a very good test might have a FN rate of 0.01 and a FP rate of 0.05. The other values are 1 minus those.

	C	no C
+	0.99	0.05
-	0.01	0.95
total	1.00	1.00

While this is a nice visual representation of the error rates, it is not enough. It is not true, for example, that all the people with a $+$ test result can be assigned to the C and no C groups in the proportion 99 : 5.

This would only be true if just as many people had cancer as were cancer-free. And of course, there is also the problem that these "rates", which look like probabilities, don't add up to 1.

The other piece of data we need is the prevalence of cancer in the population that is tested. The easiest way to see the math work out is to use a large round number for the total tested, say, 10,000.

	C	no C	
+	0.99	0.05	—
-	0.01	0.95	—
			10000

Suppose the prevalence is 2%. What we have, then, is a population consisting of 200 persons with cancer and the rest without. We leave the internal cells blank for the moment, because we are going to shift from proportions to total numbers.

	C	no C	
+	—	—	—
-	—	—	—
total	200	9800	10000

Multiply the total for each of the two groups by the proportions that we had above (the error rates). I obtain

	C	no C	
+	198	490	—
-	2	9310	—
total	200	9800	10000

We can, if we wish, fill in the marginal totals by addition.

	C	no C	total
+	198	490	688
-	2	9310	9312
total	200	9800	10000

There are two conclusions. The first is not surprising and even reassuring. If your result is negative, there is almost no chance that you have cancer. The test is very sensitive.

However, the second is that, of 688 people with a positive test result, only 198 or about $2 \text{ in } 7 = 29\%$ of them actually have cancer.

That is

$$p(+|C) = 0.99$$

but

$$p(C|+) = 0.29$$

The reason is that even though the test is both sensitive and specific, the number of cases actually present is low. This often feels counter-intuitive.

Bayes

Bayes theorem is a rule to tell us the probability that someone has cancer and also has a positive test result. It requires knowledge of FN , FP and the prevalence.

A general result from the study of probability is that

$$\begin{aligned} p(E \text{ AND } F) &= p(E|F) p(F) \\ &= p(F|E) p(E) \end{aligned}$$

(Actually, the symbol usually used for AND is \cap).

The probability of two events together both happening is equal to the probability of one times the *conditional* probability of the second, given the first. It works in both directions.

The equation can be rearranged to give

$$p(E|F) = \frac{p(F|E) p(E)}{p(F)}$$

using the symbols from our example

$$p(C \text{ AND } +) = p(C|+) p(+) = p(+|C) p(C)$$

thus

$$p(C|+) = \frac{p(+|C) p(C)}{p(+)}$$

The left-hand side is what we really want to know.

The first term in the numerator on the right-hand side is $p(+|C) = 1 - FN = 0.99$.

Next, the probability of cancer $p(C)$ is the prevalence, which we said is 2% or 0.02. This is often taken as the prevalence in the general population, although it should really be the prevalence in all people who have been administered a test. Since there is likely to be some reason for the referral to a test, this may not be the prevalence in the general population.

The third term (the denominator) is the fraction of all tests which are positive, which is 0.69. This is the part of the calculation that is a bit more complicated, it is $198 + 490 = 688$ divided by N .

After multiplying everything out, I get 29%, as expected by comparison our earlier result.

There is a slightly easier way to do the calculation, which also makes clear how the Bayes approach reflects the way we think about science.

odds form of Bayes

First, we must explain the relationship between probability and odds.

Suppose we consider the proverbial *urn* used in all probability books, which is a huge opaque jar. Our urn contains 10 red (r) balls and 90 white (w) ones.

We can call the event of drawing a red ball R , and its probability is

$p(R) = 0.10$. That is

$$p(R) = \frac{r}{r+w} \quad p(W) = \frac{w}{r+w} = 1 - p(R)$$

Odds are the ratio between two possible choices or events. They are normally given in whole numbers so here we would say that the odds *against* the less likely possibility are 9 to 1.

$$o(W) = \frac{90}{10} = \frac{w}{r} = \frac{p(W)}{p(R)}$$

Odds are very commonly used in horse racing, where for most horses the chance of winning is low and adds against of 5 to 1 are usual and 50 to 1 is rare but still seen, e.g. in the Kentucky Derby where there are many horses in the race.

If the probability of some event (my horse winning the 3rd race) is 20% or 0.20, then the probability of losing is 0.80. Even though there is only one actual Derby, we can take a frequentist perspective for the moment and say that we expect for every 80 losses we should have 20 wins so the odds against winning are 4 to 1.

The odds form of Bayes rule is

$$o(G|E) = L \cdot o(G)$$

We start with some odds of an event, G , writing $o(G)$ on the right. These odds are updated, multiplying by a factor we will describe, and the result, on the left, is the odds of G , *given* E . The updating factor is called the *likelihood ratio*.

$$L = \frac{p(E|G)}{p(E|G^c)}$$

If we want to write the whole thing out it is

$$o(G|E) = \frac{p(E|G)}{p(E|G^c)} \cdot o(G)$$

example

In the cancer test example, based on what we knew before the results of the test, the individual has some odds $o(C)$.

$$o(C|+) = L \cdot o(C)$$

The odds are updated based on the likelihood ratio according to the test statistics.

The notation is a bit more complicated in this case because normally, odds are expressed as a whole number for the more common result, namely

$$o(C^c) = 49 \text{ to } 1$$

The symbol C^c means the case where C is not true (cancer is not present). The c stands for *complement*. The meaning is that the odds $o(C^c)$ of not having cancer before we receive any data, are 49 to 1.

As a probability $p(C^c) = 1 - p(C)$. Given as odds

$$o(C^c) = \frac{1 - p(C)}{p(C)} = \frac{p(C^c)}{p(C)}$$

This part is just from the prevalence.

These odds are updated based on a different ratio of probabilities. To write this for the cancer example is a little tricky because C^c is the

common outcome. Nevertheless, with our specific symbols:

$$o(C^c|+) = \frac{p(+|C^c)}{p(+|C)} \cdot o(C^c)$$

But we can write this more simply. Recall that $p(+|C^c)$, is the false positive rate FP , the probability that a person without cancer will test positive, which was 5% or 0.05. So the numerator above is just FP .

Similarly, $p(+|C)$, is the probability that a person with cancer will test positive, that is, the true positive rate or TP . As we said is $TP = 1$ minus FN or 0.99.

The likelihood ratio is then just FP/TP so we can rewrite what we had above as

$$o(C^c|+) = \frac{FP}{TP} \cdot o(C^c)$$

Numerically, we have that $L = 5/99$. or about $1/20$.

In the update operation, the prior odds of 49 to 1 are multiplied by $1/20$ to obtain $2\frac{1}{2}$ to 1 or 5 to 2. These are the new odds that the test subject does not have cancer, based on a positive screening result. It is still more probable that we are cancer-free than not.

The ratio of probabilities may take a minute to think through, but the logic of the equation is just like what we suppose happens in the scientific enterprise. We start with some prior belief about the hypothesis, and then update that belief on the basis of new information.

If the business with the complement C^c is confusing, we can also allow odds that are less than 1, and just invert everything, since $o(C^c) = 1/o(C)$.

Hence it is also the case that

$$o(C|+) = \frac{TP}{FP} \cdot o(C)$$

The prior odds $o(C)$ before the test were 1 to 49, we multiply by 20 to obtain posterior odds of 1 to 2.5. That is 2.5 chances of no cancer for every 1 to have cancer or a probability of $1/3.5 = 29\%$.

It is even more natural to do these calculations in terms of logarithms, where updating a prior is simple addition. In that case we can talk about the *log odds*.