

The Terminology and Mathematical Analysis of Diagnostic Tests Carl Wagner 2012

When a pharmaceutical company develops a diagnostic test for a given disease, it wants the test to be as accurate as possible, i.e., to minimize the occurrence of *false positives* (individuals who are healthy, but test positive) and *false negatives* (individuals who have the disease, but test negative). So, once the test is developed, they identify a set \mathcal{D} of individuals who, based on their symptoms, almost certainly have the disease, give them the test, and identify the subset \mathcal{P} who test positive. The proportion $P(\mathcal{P}|\mathcal{D})$ is known as the *sensitivity* of the test, and is reported to physicians and public health officials. The company also identifies a set \mathcal{H} of individuals who, based on exhibiting none of the known symptoms of the disease, are almost certainly free of that disease, give them the test, and identify the subset \mathcal{N} who test negative. The proportion $P(\mathcal{N}|\mathcal{H})$ is known as the *specificity* of the disease, and is also reported to interested parties.

THE BIG QUESTION: You are given this diagnostic test, and the result is positive. No diagnostic test is perfectly reliable. So what is the probability that you actually have the disease? Or, suppose you test negative. What is the probability that you are actually healthy?

Your first response might be to ask: What do you mean by “probability” here? All the probabilities that we have encountered in this course so far have been observed relative frequencies or conditional relative frequencies in a given set S . Of course we’ve mentioned that the word “probability” is also used to describe estimated relative frequencies or conditional relative frequencies. But what does the probability that **you** have the disease, given that **you** test positive, have to do with observed or estimated relative frequencies? This is actually a very deep and subtle question that is still vigorously discussed by statisticians and philosophers and psychologists interested in how one should reason about uncertainty. Here, in oversimplified form, is how some of these scholars would answer this question:

I. Consider

S = a set of individuals among whom you are a “typical” member. This is your *reference class*.

D = the subset of S consisting of all those individuals with the given disease

$+$ = the subset of S consisting of all those individuals who would test “positive” for the disease if given the test

$-$ = the subset of S consisting of all those individuals who would test “negative” for the disease if given the test.

II. We now want to find the following relative frequencies and conditional relative frequencies for S:

$P(D | +) = PV^+$ = the *predictive value of positive* = the proportion of those who have the disease among those who test positive.

$P(\text{not } D | -) = PV^-$ = the *predictive value of negative* = the proportion of those who are healthy among those who test negative.

These important numbers can be calculated if we know

$P(D)$ = the *prevalence* of the disease in your reference class S.

$P(+|D)$ = the sensitivity of the test when applied to S.

$P(-|D)$ = the specificity of the test when applied to S.

III. A step-by-step outline of how to do this.

1. Get an estimate of $P(D)$ from a public health organization or other medical source

2. Set $P(+|D)$ = the sensitivity of the test reported by the pharmaceutical company.

3. Set $P(-| \text{not } D)$ = the specificity of the test reported by the pharmaceutical company.

Now there are actually formulas for $P(D|+)$ and $P(\text{not } D| -)$:

$$(1) \quad P(D|+) = PV^+ = \frac{\text{prevalence} \times \text{sensitivity}}{\text{prevalence} \times \text{sensitivity} + (1 - \text{prevalence}) \times (1 - \text{specificity})}$$

$$(2) \quad P(\text{not } D| -) = PV^- = \frac{(1 - \text{prevalence}) \times \text{specificity}}{(1 - \text{prevalence}) \times \text{specificity} + \text{prevalence} \times (1 - \text{sensitivity})}$$

But you do not need to memorize these formulas in order to calculate PV^+ and PV^- . As we'll see in the lecture, we can use a make-believe frequency table in order to calculate these numbers in a way that just requires basic common sense.

IV. Here's where the big conceptual leap takes place. If you test positive for the disease, adopt $PV^+ = P(D|+)$ as your **personal probability of having the disease**. If you test negative for the disease, adopt PV^- as your **personal probability of not having the disease**.

A Note on the Odds Formulation of Probability

We have been thinking of probabilities as relative frequencies, either observed (empirical probabilities) or estimated (probability models). In any case, we have $P(A) = \#A / \#S$, where S is the set of all possible individuals or possibilities under consideration. Notice that $S = A \cup A^c$. We can think of each member of A as being “favorable” to A and each member of A^c as being “unfavorable” to A . Whereas $P(A)$ is the ratio of the number of favorable cases ($\#A$) to the number of possible cases ($\#S$), what is called the *odds in favor of A*, and denoted $O(A)$ is the ratio of favorable cases to unfavorable cases, i.e.,

$$(3) \quad O(A) := \frac{\#A}{\#A^c}.$$

Many times, the actual ratio above is not computed. Instead, one says or writes : The odds in favor of A are $\#A$ to $\#A^c$, or the odds in favor of A are $\#A : \#A^c$. One can also say or write in this situation: The *odds against A* are $\#A^c$ to $\#A$. When we think of odds as given by the single number defined by (3), that number is uniquely determined (and it can be any nonnegative number; unlike probabilities, odds, are not restricted to numbers between 0 and 1). But we can state a given set of odds in many different ways. If the odds in favor of A are x to y , then for any positive number m , we can also say that the odds in favor of A are mx to my .

Example: Suppose the odds in favor of A are 2 to 1. We could also say that the odds are 4 to 2, or 200 to 100, etc. , etc.

How to convert odds to probabilities:

If the odds in favor of A are x to y , then $P(A) = \frac{x}{x+y}$.

Example: If the odds in favor of A are 2 to 1, then $P(A) = 2/3$. If the odds in favor of A are 4 to 7, then $P(A) = 4/11$.

How to convert probabilities to odds:

If $P(A) = p$, then the odds in favor of A are p to $1 - p$ (or mp to $m(1 - p)$).

Example: If $P(A) = 3/4$, then the odds in favor of A can be expressed as $3/4$ to $1/4$. But usually, we would multiply each of these numbers by some number m (here, $m = 4$) converts them to whole numbers. So here it would be more common to say that the odds in favor of A are 3 to 1.

The vast majority of journalists don't know the difference between odds and probability.

When a journalist writes “The odds of A happening are 1 in 100,” they really mean $P(A) = 1/100$. If using the word “probability” makes them break out in hives, then they should say “the

chance of A happening is 1 in 100.” When you talk about odds, you always say x **to** y. When talking about probabilities or chances, you always say x **in** y.

Answer to the reply, “You nitpicking mathematicians are making a big deal out of nothing. Suppose I, the journalist, say that the odds of A happening are 1 in 100. You say that I really mean that the probability of A is 1 in 100, and that I should say that the odds in favor of A are 1 to 99. What’s the big deal? All you did was change my 100 to a 99.”

When probabilities are small, the corresponding odds are very close to the probabilities. But this is not always the case. If the probability of A is $\frac{3}{4}$, or put another way, the chance of A is 3 in 4, then it is certainly not correct to say that the odds of A happening are 3 in 4. In fact, the odds in favor of A are 3 to 1. If the probability of A is $\frac{1}{2}$, then the odds in favor of A are 1 to 1 (so-called ‘even odds’.)

Here are the odds versions of formulas (1) and (2) above.

$$(4) \quad \frac{P(D|+)}{P(\text{not}D|+)} = \frac{\textit{prevalence}}{1 - \textit{prevalence}} \times \frac{\textit{sensitivity}}{1 - \textit{specificity}}.$$

$$(5) \quad \frac{P(\text{not}D|-)}{P(D|-)} = \frac{1 - \textit{prevalence}}{\textit{prevalence}} \times \frac{\textit{specificity}}{1 - \textit{sensitivity}}.$$

Both these formulas have the general form

$$(6) \quad \text{“posterior odds”} = \text{“prior odds”} \times \text{“likelihood ratio”}.$$