

Conditional Probability, Hypothesis Testing, and the Monty Hall Problem

Ernie Croot

September 17, 2008

On more than one occasion I have heard the comment “Probability does not exist in the real world”, and most recently I heard this in the context of plane crashes: “If you ask me what the probability that the next plane I get on will crash, I would say probably 1 in 10,000,000; but if you told me that the pilot was drunk, I might say 1 in 1,000,000. ” The problem here is that there is no “the probability” – by “the probability” we mean “the probability conditioned on the given information, making certain natural assumptions”. Thus, as the information and our assumptions change, so will our probability measure. That is the difficulty in applying probability theory to real-world problems: We often have to infer from the data a natural probability measure to use; however, in the purely theoretical world, that measure is usually just given to us.

In this short note I will work through two examples in depth to show how conditional probability and Bayes’s theorem can be used to solve them: the testing problem, and the Monty Hall problem.

1 Hypothesis Testing

If someone tells you that a test for cancer (or alcohol, or drugs, or lies etc.) is “98 percent accurate”, it would be wise to ask them what they mean, as the following example will demonstrate:

Suppose that a person goes in for a routine medical test, and one of the tests says that he has cancer. The doctor tells him that the test is not foolproof, but is nonetheless “98 percent accurate”. The patient asks the

doctor what he means, and the doctor tells him that if someone with cancer is given the test, it will correctly say “cancer” 98% of the time, and if someone does not have cancer, the test responds “no cancer” 98% of the time.

Assuming that only 2% of the population has cancer, if a person selected at random from the population tests positive for cancer, how reliable is the test? The answer turns out to be only 50%.

Here is a commonsense explanation: Suppose you have 5000 people, 2% or 100 of whom have cancer, and so 4900 of whom do not. If you give these people the test, 2% of those 4900 will test positive for “cancer” – these are what are called “false positives”. At the same time, 98 of the 100 people who really have cancer will test positive. In total you will get

$$(0.02)4900 + 98 = 196$$

test results that read “positive”; and, only 98 of these will actually have cancer. Thus, the test is only 50% reliable. In the language of conditional probability and Bayes’s theorem,

$$\begin{aligned} P(\text{cancer}|\text{tests+}) &= \frac{P(\text{tests+}|\text{cancer})P(\text{cancer})}{P(\text{tests+}|\text{cancer})P(\text{cancer}) + P(\text{tests+}|\text{no cancer})P(\text{no cancer})} \\ &= \frac{(0.98)(0.02)}{(0.98)(0.02) + (0.02)(0.98)} \\ &= \frac{1}{2}. \end{aligned}$$

1.1 Eliminating Suspects

In tests of evidence or even medical tests, there is the notion of “eliminating suspects” as not being the same as “pointing the finger of blame”. Just to give you an idea of how good a test can be to “eliminate suspects”, while at the same time being very poor at “pointing the finger of blame”, let us continue with our cancer test example: Suppose you are given that the test responds “negative for cancer”. What is the probability that it is correct – i.e. that the person really does not have cancer?

Applying Bayes’s theorem as above, one will quickly discover (try it yourself) that

$$P(\text{no cancer}|\text{tests-}) \approx 99.96\%.$$

So, the test is *extremely* good at screening out people, but very poor at deciding whether a person actually has it!

1.2 Pay attention to how the test is presented

As we have seen, if a test is presented to you as being “98% accurate”, that may not be at all true, due to a high false positive rate. What one should do if someone advertises a test as “98% accurate” is one ask them to be precise, and one should pay close attention to their answer! Also, it would be a good idea to ask them how high the false positive rate is.

If the test giver provides you with the following sort of answer, know that it is bogus: The giver of the cancer test tries to demonstrate to you that it is “highly effective” by taking 10 samples of people with cancer, and 10 samples without, and feeding the samples to the testing device. Upon getting correct results for the 20 samples, the tester tells you, “see!... the test is accurate.”

Why is the test giver’s demonstration completely bogus?

2 The Monty Hall Problem

Here I will state a variant of the problem, where the math is easier to analyze: There are three doors. Behind two are goats, and behind one is a prize. You begin by picking door 1 as the door you want the host to open. The host looks at you and says “Wait a minute. Let me show you what is behind this door,” and you see a goat behind one of doors 2 or 3. He then asks if you want to stay with door 1, or switch to the remaining unopened door.

It turns out that your chances of winning if you switch are $2/3$, making certain natural assumptions, such as that the prize was chosen at random to be behind doors 1,2,3, each with equal probability; and, that if the prize was behind door 1, then the host picked door 2 or 3 with equal probability, for the one to reveal.

A quick way to get to this answer of $2/3$ is to realize that if after initially selecting door 1 you shut your ears and eyes, and pay no attention to the host (Monty Hall), your chances of winning with door 1 must be $1/3$. So, the chances that the remaining door was the winner must be $1 - 1/3 = 2/3$. In the next subsections we will give alternative ways to solve the problem.

2.1 Exhausting the Possibilities

One way to solve the problem is to write down the sample space S , and reason for the probability of all the elementary elements; and, then, interpret the event “you win if you switch” as a subset of S .

A natural definition for S is

$$\{(1, 2), (1, 3), (2, 3), (3, 2)\},$$

where the first coordinates are the doors that the prize is behind, and the second coordinate is the door Monty decides to reveal. Since it is equally likely that the prize is behind doors 1, 2, or 3, we know that

$$P(\{(1, 2), (1, 3)\}) = P(\{(2, 3)\}) = P(\{(3, 2)\}) = \frac{1}{3}.$$

Also, we assume that

$$P(\{(1, 2)\}) = P(\{(1, 3)\}),$$

as Monty is not picky about whether to reveal door 2 or door 3 when the prize is behind door 1. With these assumptions, it follows that

$$\begin{aligned} P(\{(1, 2)\}) &= \frac{1}{6} \\ P(\{(1, 3)\}) &= \frac{1}{6} \\ P(\{(2, 3)\}) &= \frac{1}{3} \\ P(\{(3, 2)\}) &= \frac{1}{3}. \end{aligned}$$

Now, the event “You win if you switch to door 2 or 3 that Monty did not reveal” is

$$E = \{(2, 3), (3, 2)\},$$

which has probability

$$P(E) = P(\{(2, 3)\}) + P(\{(3, 2)\}) = \frac{2}{3}.$$

A common mistake here is to assume that the probabilities of all the events of S are equal, which would give $1/2$ for $P(E)$. This problem is a good example of a case where your probability measure on S is not uniform (like it is for fair die or fair coin experiments)!

2.2 Conditional Probability Approach

Here is a somewhat better way to do it using the following fact (which itself is a part of the general Bayes's Theorem): Suppose that E_1, E_2, \dots, E_k are k disjoint events whose union is all of S , the sample space. Then, given an event E we have

$$P(E) = P(E|E_1)P(E_1) + \dots + P(E|E_k)P(E_k).$$

The way we apply this theorem is as follows: As in the previous subsection, let E be the event that you win if you switch to the door 2 or 3 that Monty didn't show you; and, let E_1 be the event that the prize is behind door 1, E_2 is the event the prize is behind door 2, and E_3 is the event the prize is behind door 3. Clearly, E_1, E_2, E_3 are disjoint and their union is S . We now have

$$P(E) = P(E|E_1)P(E_1) + P(E|E_2)P(E_2) + P(E|E_3)P(E_3).$$

Clearly, $P(E|E_1) = 0$, $P(E_1) = P(E_2) = P(E_3) = 1/3$, and $P(E|E_2) = P(E|E_3) = 1$. So,

$$P(E) = 0 + \frac{1}{3} + \frac{1}{3} = \frac{2}{3}.$$

2.3 Two Variations

We consider the following two variations of the problem which look almost identical, but are subtly different (actually, it isn't that subtle if one thinks about it).

Problem 1. Given that the host reveals door 3, what is the probability of winning if you switch?

Problem 2. Given that door 3 has a goat, what is the probability of winning if you switch?

The conditional information in problem 1 is the set $A = \{(1, 3), (2, 3)\}$, and the conditional information in problem 2 is the set $B = \{(1, 3), (1, 2), (2, 3)\}$. So, in problem 2 you are given less information (that is, a larger set of possibilities).

The solution to problem 1 is thus

$$P(E|A) = \frac{P(E \cap A)}{P(A)} = \frac{P(\{(2, 3)\})}{P(A)} = \frac{1/3}{1/6 + 1/3} = \frac{2}{3}.$$

The solution to problem 2 is

$$P(E|B) = \frac{P(E \cap B)}{P(B)} = \frac{1/3}{1/6 + 1/6 + 1/3} = \frac{1}{2}.$$