

# Why $p$ -Values Make Sense, but Maybe Not in the Way You Think

Jonas Groschwitz

October 2024

**Summary:**  $p$ -Values are a valuable statistical measure. But they don't tell us directly whether a hypothesis is true or false. Rather, they can tell us how much to *shift* our *already existing* (*=prior*) *beliefs*. In this context, the arbitrary cutoff for "significant"  $p$ -values at e.g. 0.05 makes sense as a cutoff for "this experiment provides a meaningful shift to our beliefs", as opposed to a sort of "cutoff for the truth".

## 1 Basics of Significance Testing

The goal of significance testing is to mathematically determine, from a set of data, how likely a certain claim is to be true. Or more simply put, to answer the question, "should I believe this claim?". Somewhat counterintuitively, but mathematically very usefully, significance testing does not look at the likelihood of the claim directly, but looks at the likelihood of a *null hypothesis*, usually denoted as  $H_0$ . A null hypothesis describes essentially the opposite of the claim, the absence of an effect.

Take for example the classic [lady tasting tea](#) scenario described by Ronald Fisher. In this scenario, a woman claims to be able to taste whether the milk or the tea was added first to a cup of tea. The null hypothesis to this claim is that the woman is guessing at random.

To reject the null hypothesis, one can conduct an experiment – for example, giving the woman some cups of tea and writing down how often she is right, and how often she is wrong. Say, for example, the woman tastes 8 cups and is right 7 times. Let  $O$  then describe the set of experiment outcomes that are *at least as extreme* as the observed outcome. The idea of significance testing is then to look at the probability that an outcome in  $O$  occurs, assuming that the null hypothesis  $H_0$  is true; in mathematical notation this is the conditional probability

$$P(O|H_0)$$

(read " $P$  of  $O$  given  $H_0$ "). By convention, the null hypothesis is rejected if this probability is lower than some  $p$ -value, usually 0.05 (i.e. 5%), that is, if

$$P(O|H_0) < 0.05$$

To summarize, if it is sufficiently unlikely that the null hypothesis could produce the observed outcome by random chance, then the experiment provides *significant* evidence against the null hypothesis and in favor of the claim.

However, this view on statistical significance always leaves me with two questions. First, we are primarily interested in the likelihood of the null hypothesis, not in the likelihood of the experiment outcomes. That is, we really care about  $P(H_0|O)$ , not  $P(O|H_0)$ . How do the two relate? And second, what's up with this seemingly arbitrary choice of the  $p$ -value as 0.05? A cutoff point of 5% seems like an odd choice if we really care about the truth. I found quite satisfying answers to these questions when looking through the lense of Bayes' theorem.

## 2 The Bayesian View on Significance Testing

Bayes' Theorem, a central theorem in probability theory, connects two conditional probabilities. It states that for two events  $A$  and  $B$ ,

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}.$$

This provides us with an answer to the first question: how to relate  $P(H_0|O)$  and  $P(O|H_0)$ . Namely:

$$P(H_0|O) = \frac{P(O|H_0)P(H_0)}{P(O)}.$$

This is great news: we can now compute  $P(H_0|O)$ , which is what we really care about, from  $P(O|H_0)$ , which is what the statistical significance tests give us! The problem is only that there are two other pesky terms in the equation,  $P(H_0)$  and  $P(O)$ .  $P(H_0)$  is called the *prior* or *prior belief*. We'll look at  $P(O)$  in the next section.

The prior belief is how much we believe that the null hypothesis is true *before* we conduct the experiment. It is also the main reason why scientific analysis tends to report  $P(O|H_0)$  rather than the  $P(H_0|O)$  (also called the *posterior* – our belief in the null hypothesis *after* the experiment). The problem here is that the prior belief is subjective (you can imagine that people will differ widely in how much they believe the lady that she can distinguish the tea, in the absence of any evidence). And one cannot compute the posterior without knowing the prior.

In other words, it is not possible to compute the posterior fully objectively. The only thing we can really objectively get is  $P(O|H_0)$ . But the Bayesian perspective is still very useful, as we will see in the next section.

### 3 Updating One's Beliefs with $p$ -Values

In this section, we will first work out how exactly to update our beliefs using the  $p$ -value. That is, how to compute the posterior from the  $p$ -value and the prior. To shape our intuition, we'll then look at some concrete numbers of  $p$ -values and priors, and the posteriors they yield. Finally, we'll get back to our earlier two questions about significance testing, and see how this perspective resolves them quite nicely.

We left off above with the Bayesian equation

$$P(H_0|O) = \frac{P(O|H_0)P(H_0)}{P(O)}.$$

To describe the posterior purely in terms of  $p$ -value  $P(O|H_0)$  and prior  $P(H_0)$ , let us spell out  $P(O)$ , the likelihood of the evidence (without assuming the null-hypothesis). By the law of total probability, we have

$$P(O) = P(O|H_0)P(H_0) + P(O|\neg H_0)P(\neg H_0),$$

where  $\neg H_0$  means "the null hypothesis does not hold". The first summand are our old acquaintances  $p$ -value and prior. In the second summand, we can write  $P(\neg H_0)$  as  $1 - P(H_0)$ .

The remaining term,  $P(O|\neg H_0)$ , is the likelihood of the observation (or a more extreme result) under the condition that the null hypothesis does *not* hold. This number can often be assumed to be so close to one that we can just round it,  $P(O|\neg H_0) \approx 1$ . For example, if the lady can indeed taste the difference between putting the milk or the tea first into the cup, then the likelihood of her getting at least 7 out of 8 cups right is very high. Note that this assumption does not always hold, and what to do when it doesn't is a whole topic on its own. But here, we'll just focus on the simple cases where this approximation  $P(O|\neg H_0) \approx 1$  is appropriate.

Under this assumption, we then get

$$P(O) = P(O|H_0)P(H_0) + 1 - P(H_0),$$

and thus our final formula for the posterior, expressed purely in terms of the  $p$ -value and the prior:

$$P(H_0|O) = \frac{P(O|H_0)P(H_0)}{P(O|H_0)P(H_0) + 1 - P(H_0)}.$$

What a mouthful! To better understand what this means, let's plug in some numbers.

TODO: plots

To summarize, let us come back to our original questions. First, if we care about the posterior  $P(H_0|O)$ , why do we compute the  $p$ -value  $P(O|H_0)$  instead? And how do they relate? The answer is that only the  $p$ -value can be computed objectively, but our subjective beliefs can be updated with the methodology we just discussed. In other words, the  $p$ -value cannot tell us directly whether we should believe in the null-hypothesis (or disbelieve it, with related belief in the claim), but it can tell us how

to shift our existing prior belief. This allows us to make our beliefs more and more based on evidence, bringing us as close to objective consensus as we can get.

Second, why the  $p$ -value cutoff at 0.05? From a Bayesian perspective, this arbitrary cutoff can be interpreted as stating that experiments of at least this  $p$ -value provide a shift in belief meaningful enough ("significant") that they are worth taking note of (and publishing a paper about them etc.). While the value 0.05 is arbitrary, to me at least it actually makes intuitive sense as an "update threshold" – much more so than a "truth threshold".

While even in this view, issues such as the [multiple comparisons problem](#) and [publication bias](#) are a problem, I hope I was able to give you an appreciation for why  $p$ -values are useful, and maybe a bit more precise intuition for how to interpret them.

(Note: is  $p$ -value the cutoff point for  $P(O|H_0)$ , or is it  $P(O|H_0)$  itself? – It is the latter, and I'll have to adapt the text to that. Actually makes it easier to write!)