

## 246 Excursion 4: Objectivity and Auditing

---

But we cannot just deny one-sided tests, nor does Cumming. In fact, he encourages their use: “it’s unfortunate they are usually ignored” (2012, p. 113). (He also says he is happy for people to decide afterwards whether to report it as a one- or two-sided interval (ibid., p. 112), only doubling  $\alpha$ , which I do not mind.) Still needed is a justification for bringing in the upper limit when applying a one-sided estimator, and severity supplies it. You should always be interested in at least *two benchmarks*: discrepancies well warranted and those terribly warranted. In test T+, our handy benchmark for the terrible is to set the lower limit to  $\bar{x}$ . The severity for  $(\mu > \bar{x})$  is 0.5. Two side notes:

First I grant it would be wrong to charge Cumming with treating all parameter values within the confidence interval *on par*, because he does suggest distinguishing them by their likelihoods (by how probable each renders the outcome). Take just the single 0.975 lower CI bound with  $n = 100$  and  $\bar{x} = 0.2$ . A  $\mu$  value closer to the observed 0.2 has higher likelihood (in the technical sense) than ones close to the 0.975 lower limit 0. For example,  $\mu = 0.15$  is more likely than  $\mu = 0.05$ . However, this moves away from CI reasoning (toward likelihood comparisons). The claim  $\mu > 0.05$  has a *higher* confidence level (0.93) than does  $\mu > 0.15$  ( $0.7$ )<sup>3</sup> even though the point hypothesis  $\mu = 0.05$  is less likely than  $\mu = 0.15$  (the latter is closer to  $\bar{x} = 0.2$  than is the former). Each point in the lower CI corresponds to a different lower bound, each associated with a different confidence level, and corresponding severity assessment. That’s how to distinguish them.

Second there’s an equivocation, or at least a potential equivocation, in Cumming’s assertion “that for [2.5%] of replications the [lower limit] will exceed the true value” (Cumming 2012, p. 112 replacing 5% with 2.5%). This is not a true claim if “lower limit” is replaced by a *particular* lower limit:  $\hat{\mu}_{0.025}(\bar{x})$ , it holds only for the *generic* lower limit  $\hat{\mu}_{0.025}(\bar{X})$ . That is, we can’t say  $\mu$  exceeds zero 2.5% of the time, which would be to assign a probability of 0.975 to  $\mu > 0$ . Yet this misinterpretation of CIs is legion, as we’ll see in a historical battle about fiducial intervals (Section 5.8).

### 4.4 Do P-Values Exaggerate the Evidence?

“Significance levels overstate the evidence against the null hypothesis,” is a line you may often hear. Your first question is:

What do you mean by overstating the evidence against a hypothesis?

Several (honest) answers are possible. Here is one possibility:

<sup>3</sup> Subtract 1.5 SE and 0.5 SE from  $\bar{x} = 0.2$ , respectively.

---

**Tour II: Rejection Fallacies: Who's Exaggerating What? 247**

---

What I mean is that when I put a lump of prior weight  $\pi_0$  of 1/2 on a point null  $H_0$  (or a very small interval around it), the  $P$ -value is smaller than my Bayesian posterior probability on  $H_0$ .

More generally, the “ $P$ -values exaggerate” criticism typically boils down to showing that if inference is appraised via one of the probabilisms – Bayesian posteriors, Bayes factors, or likelihood ratios – the evidence against the null (or against the null and in favor of some alternative) isn't as big as  $1 - P$ .

You might react by observing that: (a)  $P$ -values are not intended as posteriors in  $H_0$  (or Bayes ratios, likelihood ratios) but rather are used to determine if there's an indication of discrepancy from, or inconsistency with,  $H_0$ . This might only mean it's worth getting more data to probe for a real effect. It's not a degree of belief or comparative strength of support to walk away with. (b) Thus there's no reason to suppose a  $P$ -value should match numbers computed in very different accounts, that differ among themselves, and are measuring entirely different things. Stephen Senn gives an analogy with “height and stones”:

... [S]ome Bayesians in criticizing  $P$ -values seem to think that it is appropriate to use a threshold for significance of 0.95 of the probability of the alternative hypothesis being true. This makes no more sense than, in moving from a minimum height standard (say) for recruiting police officers to a minimum weight standard, declaring that since it was previously 6 foot it must now be 6 stone. (Senn 2001b, p. 202)

To top off your rejoinder, you might ask: (c) Why assume that “the” or even “a” correct measure of evidence (relevant for scrutinizing the  $P$ -value) is one of the probabilist ones?

All such retorts are valid, and we'll want to explore how they play out here. Yet, I want to push beyond them. Let's be open to the possibility that evidential measures from very different accounts can be used to scrutinize each other.

**Getting Beyond “I'm Rubber and You're Glue”.** The danger in critiquing statistical method  $X$  from the standpoint of the goals and measures of a distinct school  $Y$ , is that of falling into begging the question. If the  $P$ -value is exaggerating evidence against a null, meaning it seems too small from the perspective of school  $Y$ , then  $Y$ 's numbers are too big, or just irrelevant, from the perspective of school  $X$ . Whatever you say about me bounces off and sticks to you. This is a genuine worry, but it's not fatal. The goal of this journey is to identify minimal theses about “bad evidence, no test (BENT)” that enable some degree of scrutiny of any statistical inference account – at least on the meta-level. Why assume all schools of statistical inference embrace the minimum severity principle? I don't, and they don't. But by identifying when methods violate

## 248 Excursion 4: Objectivity and Auditing

---

severity, we can pull back the veil on at least one source of disagreement behind the battles.

Thus, in tackling this latest canard, let's resist depicting the critics as committing a gross blunder of confusing a  $P$ -value with a posterior probability in a null. We resist, as well, merely denying we care about their measure of support. I say we should look at exactly what the critics are on about. When we do, we will have gleaned some short-cuts for grasping a plethora of critical debates. We may even wind up with new respect for what a  $P$ -value, the least popular girl in the class, really *does*.

To visit the core arguments, we travel to 1987 to papers by J. Berger and Sellke, and Casella and R. Berger. These, in turn, are based on a handful of older ones (Cox 1977, E, L, & S 1963, Pratt 1965), and current discussions invariably revert back to them. Our struggles through quicksand of Excursion 3, Tour II, are about to pay large dividends.

**J. Berger and Sellke, and Casella and R. Berger.** Berger and Sellke (1987a) make out the conflict between  $P$ -values and Bayesian posteriors by considering the two-sided test of the Normal mean,  $H_0: \mu = 0$  vs.  $H_1: \mu \neq 0$ . "Suppose that  $\mathbf{X} = (X_1, \dots, X_n)$ , where the  $X_i$  are IID  $N(\mu, \sigma^2)$ ,  $\sigma^2$  known" (p. 112). Then the test statistic  $d(\mathbf{X}) = \sqrt{n}|\bar{X} - \mu_0|/\sigma$ , and the  $P$ -value will be twice the  $P$ -value of the corresponding one-sided test.

Starting with a lump of prior, generally 0.5, on the point hypothesis  $H_0$ , they find the posterior probability in  $H_0$  is larger than the  $P$ -value for a variety of different priors on the alternative. However, the result depends entirely on how the remaining 0.5 is allocated or smeared over the alternative (a move dubbed spike and smear). Using what they call a Jeffreys-type prior, the 0.5 is spread out over the alternative parameter values as if the parameter is itself distributed  $N(\mu_0, \sigma)$ . Now Harold Jeffreys recommends the lump prior only to capture cases where a special value of a parameter is deemed plausible, for instance, the GTR deflection effect  $\lambda = 1.75''$ , after about 1960. The rationale is to avoid a 0 prior on  $H_0$  and enable it to receive a reasonable posterior probability.

By subtitling their paper "The irreconcilability of  $P$ -values and evidence," Berger and Sellke imply that if  $P$ -values disagree with posterior assessments, they can't be measures of evidence at all. Casella and R. Berger (1987) retort that "reconciling" is at hand, if you move away from the lump prior. So let's see how this unfolds. I assume throughout, as do the critics, that the  $P$ -values are "audited," so that neither selection effects nor violated model assumptions are in question at this stage. I see no other way to engage their arguments.

Table 4.1  $\Pr(H_0|\mathbf{x})$  for Jeffreys-type prior

P one-sided	$z_\alpha$	$n$ (sample size)				
		10	20	50	100	1000
0.05	1.645	0.47	0.56	0.65	0.72	0.89
0.025	1.960	0.37	0.42	0.52	0.60	0.82
0.005	2.576	0.14	0.16	0.22	0.27	0.53
0.0005	3.291	0.024	0.026	0.034	0.045	0.124

(From Table 1, J. Berger and T. Sellke (1987) p. 113 using the one-sided  $P$ -value)

Table 4.1 gives the values of  $\Pr(H_0|\mathbf{x})$ . We see that we would declare no evidence against the null, and even evidence for it (to the degree indicated by the posterior) whenever  $d(\mathbf{x})$  fails to reach a 2.5 or 3 standard error difference. With  $n = 50$ , "one can classically 'reject  $H_0$  at significance level  $p = 0.05$ ,' although  $\Pr(H_0|\mathbf{x}) = 0.52$  (which would actually indicate that the evidence favors  $H_0$ )" (J. Berger and Sellke 1987, p. 113).

If  $n = 1000$ , a result statistically significant at the 0.05 level results in the posterior probability to  $\mu = 0$  going up from 0.5 (the lump prior) to 0.82! From their Bayesian perspective, this suggests  $P$ -values are exaggerating evidence against  $H_0$ . Error statistical testers, on the other hand, balk at the fact that using the recommended priors allows statistically significant results to be interpreted as no evidence against  $H_0$  – or even evidence for it. They point out that 0 is excluded from the two-sided confidence interval at level 0.95. Although a posterior probability doesn't have an error probability attached, a tester can evaluate the error probability credentials of these inferences. Here we'd be concerned with a Type II error: failing to find evidence against the null, and providing a fairly high posterior for it, when it's false (Souvenir I).

Let's use a less extreme example where we have some numbers handy: our water-plant accident. We had  $\sigma = 10$ ,  $n = 100$  leading to the nice  $(\sigma/\sqrt{n})$  value of 1. Here it would be two-sided, to match their example:  $H_0: \mu = 150$  vs.  $H_1: \mu \neq 150$ . Look at the second entry of the 100 column, the posterior when  $z_\alpha = 1.96$ . With the Jeffreys prior, perhaps championed by the water coolant company, J. Berger and Sellke assign a posterior of 0.6 to  $H_0: \mu = 150$  degrees when a mean temperature of 152 (151.96) degrees is observed – reporting decent evidence the cooling mechanism is working just fine. How often would this occur even if the actual underlying mean temperature is, say, 151 degrees? With a two-sided test, cutting off 2 standard errors on either side, we'd reject whenever either  $\bar{X} \geq 152$  or  $\bar{X} \leq 148$ . The probability of the second is negligible under  $\mu = 151$ , so the probability we want is

## 250 Excursion 4: Objectivity and Auditing

---

$\Pr(\bar{X} < 152; \mu = 151) = 0.84$  ( $Z = (152 - 151) = 1$ ). The probability of declaring evidence for 150 degrees (with posterior of 0.6 to  $H_0$ ) even if the true increase is actually 151 degrees is around 0.84; 84% of the time they erroneously fail to ring the alarm, and would boost their probability of  $\mu = 150$  from 0.5 to 0.6. Thus, from our minimal severity principle, the statistically significant result can't even be taken as evidence for compliance with 151 degrees, let alone as evidence for the null of 150 (Table 3.1).

Is this a problem for them? It depends what you think of that prior. The N-P test, of course, does not use a prior, although, as noted earlier, one needn't rule out a frequentist prior on mean water temperature after an accident (Section 3.2). For now our goal is making out the criticism.

### Jeffreys–Lindley “Paradox” or Bayes/Fisher Disagreement

But how, intuitively, does it happen that a statistically significant result corresponds to a Bayes boost for  $H_0$ ? Go back to J. Berger and Sellke's example of Normal testing of  $H_0: \mu = 0$  vs.  $H_1: \mu \neq 0$ . Some sample mean  $\bar{x}$  will be close enough to 0 to increase the posterior for  $H_0$ . By choosing a sufficiently large  $n$ , even a statistically significant result can correspond to large posteriors on  $H_0$ . This is the Jeffreys–Lindley “paradox,” which some more aptly call the Bayes/Fisher disagreement. Lindley's famous result dealt with just this example, two-sided Normal testing with known variance. With a lump given to the point null, and the rest appropriately spread over the alternative, an  $n$  can be found such that an  $\alpha$  significant result corresponds to  $\Pr(H_0|x) = (1 - \alpha)!$  We can see by extending Table 4.1 to arbitrarily large  $n$ , we can get a posterior for the null of 0.95, when the (two-sided)  $P$ -value is 0.05. Many say you should decrease the required  $P$ -value for significance as  $n$  increases; and Cox and Hinkley (1974, p. 397) provide formulas to achieve this and avoid the mismatch. There's nothing in N-P or Fisherian theory to oppose this. I won't do that here, as I want to make out the criticism. We need only ensure that the interpretation takes account of the (obvious) fact that, with a fixed  $P$ -value and increasing  $n$ , the test is more and more sensitive to smaller and smaller discrepancies. Using a smaller plate at the French restaurant may make the portion appear bigger, but, Jackie Mason notwithstanding, knowing the size of the plate, I can see there's not much there.

Why assign the lump of  $\frac{1}{2}$  as prior to the point null? “The choice of  $\pi_0 = 1/2$  has obvious intuitive appeal in scientific investigations as being ‘objective’” say J. Berger and Sellke (1987, p. 115). But is it? One starts by making  $H_0$  and  $H_1$  equally probable, then the 0.5 accorded to  $H_1$  is spread out over all the values in  $H_1$ : “The net result is that all values of  $[\mu]$  are far from being equally likely”

(Senn 2015a). Any small group of  $\mu$  values in  $H_1$  gets a tiny prior. David Cox describes how it happens:

... if a sample say at about the 5% level of significance is achieved, then either  $H_0$  is true or some alternative in a band of order  $1/\sqrt{n}$ ; the latter possibility has, as  $n \rightarrow \infty$ , prior probability of order  $1/\sqrt{n}$  and hence at a fixed level of significance the posterior probabilities shift in favour of  $H_0$  as  $n$  increases. (Cox 1977, p. 59)

What justifies the lump prior of 0.5?

### A Dialogue at the Water Plant Accident

EPA REP: The mean temperature of the water was found statistically significantly higher than 150 degrees at the 0.025 level.

SPIKED PRIOR REP: This even strengthens my belief the water temperature's no different from 150. If I update the prior of 0.5 that I give to the null hypothesis, my posterior for  $H_0$  is still 0.6; it's not 0.025 or 0.05, that's for sure.

EPA REP: Why do you assign such a high prior probability to  $H_0$ ?

SPIKED PRIOR REP: If I gave  $H_0$  a value lower than 0.5, then, if there's evidence to reject  $H_0$ , at most I would be claiming an improbable hypothesis has become more improbable.

[W]ho, after all, would be convinced by the statement 'I conducted a Bayesian test of  $H_0$ , assigning prior probability 0.1 to  $H_0$ , and my conclusion is that  $H_0$  has posterior probability 0.05 and should be rejected?' (J. Berger and Sellke 1987, p. 115).

This quote from J. Berger and Sellke is peculiar. They go on to add: "We emphasize this obvious point because some react to the Bayesian–classical conflict by attempting to argue that [prior]  $\pi_0$  should be made small in the Bayesian analysis so as to force agreement" (ibid.). We should not force agreement. But it's scarcely an obvious justification for a lump of prior on the null  $H_0$  – one which results in a low capability to detect discrepancies – that it ensures, if they *do* reject  $H_0$ , there will be a meaningful drop in its probability. Let's listen to the pushback from Casella and R. Berger (1987a), the Berger being Roger now (I use initials to distinguish them).

**The Cult of the Holy Spike and Slab.** Casella and R. Berger (1987a) charge that the problem is not  $P$ -values but the high prior, and that "concentrating mass on the point null hypothesis is biasing the prior in favor of  $H_0$  as much as possible" (p. 111) whether in one- or two-sided tests. According to them:

The testing of a point null hypothesis is one of the most misused statistical procedures. In particular, in the location parameter problem, the point null

## 252 Excursion 4: Objectivity and Auditing

---

hypothesis is more the mathematical convenience than the statistical method of choice. (ibid., p. 106)

Most of the time “there is a direction of interest in many experiments, and saddling an experimenter with a two-sided test would not be appropriate” (ibid.). The “cult of the holy spike” is an expression I owe to Sander Greenland (personal communication).

By contrast, we can reconcile  $P$ -values and posteriors in one-sided tests if we use more diffuse priors. (e.g., Cox and Hinkley 1974, Jeffreys 1939/1961, Pratt 1965). In fact, Casella and Berger show that for sensible priors in that case, the  $P$ -value is at least as big as the minimum value of the posterior probability on the null, again contradicting claims that  $P$ -values exaggerate the evidence.<sup>4</sup>

J. Berger and Sellke (1987) adhere to the spikey priors, but following E, L, & S (1963), they’re keen to show that  $P$ -values exaggerate evidence even in cases less extreme than the Jeffreys posteriors in Table 4.1. Consider the likelihood ratio of the null hypothesis over the hypothesis most generous to the alternative, they say. This is the point alternative with maximum likelihood,  $H_{\max}$  – arrived at by setting  $\mu = \bar{x}$ . Through their tunnel, it’s disturbing that even using this likelihood ratio, the posterior for  $H_0$  is still larger than 0.05 – when they give a 0.5 spike to both  $H_0$  and  $H_{\max}$ . Some recent authors see this as the key to explain today’s lack of replication of significant results. Through the testing tunnel, things look different (Section 4.5).

**Why Blame Us Because You Can’t Agree on Your Posterior?** Stephen Senn argues that the reason for the wide range of variation of the posterior is the fact that it depends radically on the choice of alternative to the null and its prior.<sup>5</sup> According to Senn, “. . . the reason that Bayesians can regard  $P$ -values as overstating the evidence against the null is simply a reflection of the fact that Bayesians can disagree *sharply* with each other” (Senn 2002, p. 2442). Senn

<sup>4</sup> Casella and R. Berger (1987b) argue, “We would be surprised if most researchers would place even a 10% prior probability of  $H_0$ . We hope that the casual reader of Berger and Delampady realizes that the big discrepancies between  $P$ -values  $P(H_0|x)$  . . . are due to a large extent to the large value of [the prior of 0.5 to  $H_0$ ] that was used.” The most common uses of a point null, asserting the difference between means is 0, or the coefficient of a regression coefficient is 0, merely describe a potentially interesting feature of the population, with no special prior believability. “Berger and Delampady admit . . .,  $P$ -values are reasonable measures of evidence when there is no a priori concentration of belief about  $H_0$ ” (ibid., p. 345). Thus, “the very argument that Berger and Delampady use to dismiss  $P$ -values can be turned around to argue *for*  $P$ -values” (ibid., p. 346).

<sup>5</sup> In defending spiked priors, Berger and Sellke move away from the importance of effect size. “Precise hypotheses . . . ideally relate to, say, some precise theory being tested. Of primary interest is whether the theory is right or wrong; the amount by which it is wrong may be of interest in developing alternative theories, but the initial question of interest is that modeled by the precise hypothesis test” (1987, p. 136).

illustrates how “two Bayesians having the same prior probability that a hypothesis is true and having seen the same data can come to radically different conclusions because they differ regarding the alternative hypothesis” (Senn 2001b, p. 195). One of them views the problem as a one-sided test and gets a posterior on the null that matches the  $P$ -value; a second chooses a Jeffreys-type prior in a two-sided test, and winds up with a posterior to the null of  $1 - p$ !

Here's a recap of Senn's example (ibid., p. 200): Two scientists A and B are testing a new drug to establish its treatment effect,  $\delta$ , where positive values of  $\delta$  are good. Scientist A has a vague prior whereas B, while sharing the same distribution about the probability of positive values of  $\delta$ , is less pessimistic than A regarding the effect of the drug. If it's not useful, B believes it will have no effect. They “share the same belief that the drug has a positive effect. Given that it has a positive effect, they share the same belief regarding its effect. . . . They differ only in belief as to how harmful it might be.” A clinical trial yields a difference of 1.65 standard units, a one-sided  $P$ -value of 0.05. The result is that A gives 1/20 posterior probability to  $H_0$ : the drug does *not* have a positive effect, while B gives a probability of 19/20 to  $H_0$ . B is using the two-sided test with a lump of prior on the null ( $H_0: \mu = 0$  vs.  $H_1: \mu \neq 0$ ), while A is using a one-sided test T+ ( $H_0: \mu \leq 0$  vs.  $H_1: \mu > 0$ ). The contrast, Senn observes, is that of Cox's distinction between “precise and dividing hypothesis” (Section 3.3). “[F]rom a common belief in the drug's efficacy they have moved in opposite directions” (ibid., pp. 200–201). Senn riffs on Jeffreys' well-known joke that we heard in Section 3.4:

It would require that a procedure is dismissed [by significance testers] because, when combined with information which it doesn't require and which may not exist, it disagrees with a [Bayesian] procedure that disagrees with itself. (ibid., p. 195)

In other words, if Bayesians disagree with each other even when they're measuring the same thing – posterior probabilities – why be surprised that disagreement is found between posteriors and  $P$ -values? The most common argument behind the “ $P$ -values exaggerate evidence” appears not to hold water. Yet it won't be zapped quite so easily, and will reappear in different forms.

#### **Exhibit (vii): Contrasting Bayes Factors and Jeffreys–Lindley Paradox.**

We've uncovered some interesting bones in our dig. Some lead to seductive arguments purporting to absolve the latitude in assigning priors in Bayesian tests. Take Wagenmakers and Grünwald (2006, p. 642): “Bayesian hypothesis tests are often criticized because of their dependence on prior distributions

## 254 Excursion 4: Objectivity and Auditing

... [yet] no matter what prior is used, the Bayesian test provides substantially less evidence against  $H_0$  than”  $P$ -values, in the examples we’ve considered. Be careful in translating this. We’ve seen that what counts as “less” evidence runs from seriously underestimating to overestimating the discrepancy we are entitled to infer with severity. Begin with three types of priors appealed to in some prominent criticisms revolving around the Fisher – Jeffreys disagreement.

1. *Jeffreys-type prior with the “spike and slab” in a two-sided test.* Here, with large enough  $n$ , a statistically significant result becomes evidence for the null; the posterior to  $H_0$  exceeds the lump prior.
2. *Likelihood ratio most generous to the alternative.* Here, there’s a spike to a point null,  $H_0: \theta = \theta_0$  to be compared to the point alternative that’s maximally likely  $\theta_{\max}$ . Often, both  $H_0$  and  $H_{\max}$  are given 0.5 priors.
3. *Matching.* Instead of a spike prior on the null, it uses a smooth diffuse prior, as in the “dividing” case. Here, the  $P$ -value “is an approximation to the posterior probability that  $\theta < 0$ ” (Pratt 1965, p. 182).

In sync with our attention to high-energy particle physics (HEP) in Section 3.6, consider an example that Aris Spanos (2013b) explores in relation to the Jeffreys–Lindley paradox. The example is briefly noted in Stone (1997).

A large number ( $n = 527,135$ ) of independent collisions that can be of either type A or type B are used to test if the proportion of type A collisions is exactly 0.2, as opposed to any other value. It’s modeled as  $n$  Bernoulli trials testing  $H_0: \theta = 0.2$  vs.  $H_1: \theta \neq 0.2$ . The observed proportion of type A collisions is scarcely greater than the point null of 0.2:

$$\bar{x} = k/n = 0.20165233, \text{ where } n = 527,135, k = 106,298.$$

**The significance level against  $H_0$  is small** (so there’s evidence against  $H_0$ )

The test statistic  $d(\mathbf{X}) = [\sqrt{n}(\bar{X} - 0.2)/\sigma] = 3$ ,  $\sigma = \sqrt{[\theta(1 - \theta)]}$ , which under the null is  $\sqrt{[0.2(0.8)]} = 0.4$ . The significance level associated with  $d(\mathbf{x}_0)$  in this two-sided test is

$$\Pr(|d(\mathbf{X})| > |d(\mathbf{x}_0)|; H_0) = 0.0027.$$

So the result  $\bar{x}$  is highly significant, even though it’s scarcely different from the point null.

**The Bayes factor in favor of  $H_0$  is high**

$H_0$  is given the spiked prior of 0.5, and the remaining 0.5 is spread equally among the values in  $H_1$ . I follow Spanos' computations:<sup>6</sup>

$$\Pr(k|H_0) = \binom{n}{k} 0.2^k (0.8)^{n-k},$$

$$\Pr(k|H_1) = \int_0^1 \binom{n}{k} \theta^k (1-\theta)^{n-k} d\theta = 1/(n+1),$$

where  $n = 527,135$  and  $k = 106,298$ .

$$\begin{aligned} \text{The Bayes factor } B_{01} &= \Pr(k|H_0) / \Pr(k|H_1) = 0.000015394 / 0.000001897 \\ &= 8.115. \end{aligned}$$

While the likelihood of  $H_0$  in the numerator is tiny, the likelihood of  $H_1$  is even tinier. Since  $B_{01}$  in favor of  $H_0$  is 8, which is greater than 1, the posterior for  $H_0$  goes up, even though the outcome is statistically significantly greater than the null.

There's no surprise once you consider the Bayesian question here: compare the likelihood of a result scarcely different from 0.2 being produced by a universe where  $\theta = 0.2$  – where this has been given a spiked prior of 0.5 under  $H_0$  – with the likelihood of that result being produced by any  $\theta$  in a small band of  $\theta$  values, which have been given a very low prior under  $H_1$ . Clearly,  $\theta = 0.2$  is more likely, and we have an example of the Jeffreys–Fisher disagreement.

Who should be afraid of this disagreement (to echo the title of Spanos' paper)? Many tribes, including some Bayesians, think it only goes to cast doubt on this particular Bayes factor. Compare it with proposal 2 in Exhibit (vii): the *Likelihood ratio most generous to the alternative*:  $\text{Lik}(0.2)/\text{Lik}(\theta_{\max})$ . We know the maximally likely value for  $\theta$ ,  $\theta_{\max} = \bar{x}$ :

$$\bar{x} = k/n = 0.20165233 = \theta_{\max},$$

$$\Pr(k|H_{\max}) = \binom{n}{k} 0.20165233^k (1-0.20165233)^{n-k} = 0.0013694656,$$

$$\text{Lik}(0.2) = 0.000015394, \text{ and } \text{Lik}(\theta_{\max}) = 0.0013694656.$$

Now  $B_{01}$  is 0.01 and  $B_{10}$ ,  $\text{Lik}(\theta_{\max})/\text{Lik}(0.2) = 89$ .

Why should a result 89 times more likely under alternative  $\theta_{\max}$  than under  $\theta = 0.2$  be taken as strong evidence for  $\theta = 0.2$ ? It shouldn't, according to some, including Lindley's own student, default Bayesian José Bernardo (2010).

<sup>6</sup> The spiked prior drops out, so the result is the same as a uniform prior on the null and alternative.

## 256 Excursion 4: Objectivity and Auditing

---

Presumably, the Likelihoodist concurs. There are family feuds within and between the diverse tribes of probabilisms.<sup>7</sup>

**Greenland and Poole** Given how often spiked priors arise in foundational arguments, it's worth noting that even Bayesians Edwards, Lindman, and Savage (1963, p. 235), despite raising the “*P*-values exaggerate” argument, aver that for Bayesian statisticians, “no procedure for testing a sharp null hypothesis is likely to be appropriate unless the null hypothesis deserves special initial credence.” Epidemiologists Sander Greenland and Charles Poole, who claim not to identify with any one statistical tribe, but who often lead critics of significance tests, say:

Our stand against spikes directly contradicts a good portion of the Bayesian literature, where null spikes are used too freely to represent the belief that a parameter ‘differs negligibly’ from the null. In many settings ... even a tightly concentrated probability near the null has no basis in genuine evidence. Many scientists and statisticians exhibit quite a bit of irrational prejudice in favor of the null ... (2013, p. 77).

They angle to reconcile *P*-values and posteriors, and to this end they invoke the matching result in # 3, Exhibit (vii). An uninformative prior, assigning equal probability to all values of the parameter, allows the *P*-value to approximate the posterior probability that  $\theta < 0$  in one-sided testing ( $\theta \leq 0$  vs.  $\theta > 0$ ). In two-sided testing, the posterior probability that  $\theta$  is on the opposite side of 0 than the observed is  $P/2$ . They proffer this as a way “to live with” *P*-values. Commenting on them, Andrew Gelman (2013, p. 72) raises this objection:

[C]onsider what would happen if we routinely interpreted one-sided *P* values as posterior probabilities. In that case, an experimental result that is 1 standard error from zero – that is, exactly what one might expect from chance alone – would imply an 83% posterior probability that the true effect in the population has the same direction as the observed pattern in the data at hand. It does not make sense to me to claim 83% certainty – 5 to 1 odds [to  $H_1$ ] ...

(The *P*-value is 0.16.) Rather than relying on non-informative priors, Gelman prefers to use prior information that leans towards the null. This avoids as high a posterior to  $H_1$  as when using the matching result.

Greenland and Poole respond that Gelman is overlooking the hazard of “strong priors that are not well founded. ... Whatever our prior opinion and

<sup>7</sup> Bernardo shocked his mentor in announcing that the Lindley paradox is really an indictment of the Bayesian computations: “Whether you call this a paradox or a disagreement, the fact that the Bayes factor for the null may be arbitrarily large for sufficiently large *n*, however relatively unlikely the data may be under  $H_0$  is, ... deeply disturbing” (Bernardo 2010, p. 59).

its foundation, we still need reference analyses with weakly informative priors to alert us to how much our prior probabilities are driving our posterior probabilities” (2013, p. 76). They rightly point out that, in some circles, giving weight to the null can be the outgrowth of some ill-grounded metaphysics about “simplicity.” Or it may be seen as an assumption akin to a presumption of innocence in law. So the question turns on the appropriate prior on the null.

Look what has happened! The problem was simply to express “I’m not impressed” with a result reaching a  $P$ -value of 0.16: Differences even larger than 1 standard error are not so very infrequent – they occur 16% of the time – even if there’s zero effect. So I’m not convinced of the reality of the effect, based on this result.  $P$ -values did their job, reflecting as they do the severity requirement.  $H_1$  has passed a lousy test. That’s that. No prior probability assignment to  $H_0$  is needed. Problem solved.

But there’s a predilection for changing the problem (if you’re a probabilist). Greenland and Poole feel they’re helping us to live with  $P$ -values without misinterpretation. By choosing the prior so that the  $P$ -value matches the posterior on  $H_0$ , they supply us “with correct interpretations” (ibid., p. 77) where “correct interpretations” are those where the misinterpretation (of a  $P$ -value as a posterior in the null) is not a misinterpretation. To a severe tester, this results in completely changing the problem from an assessment of how well tested the reality of the effect is, with the given data, to what odds I would give in betting, or the like. We land in the same uncharted waters as with other attempts to fix  $P$ -values, when we could have stayed on the cruise ship, interpreting  $P$ -values as intended.

### Souvenir Q: Have We Drifted From Testing Country? (Notes From an Intermission)

Before continuing, let’s pull back for a moment, and take a coffee break at a place called Spike and Smear. Souvenir Q records our notes. We’ve been exploring the research program that appears to show, quite convincingly, that significance levels exaggerate the evidence against a null hypothesis, based on evidential assessments endorsed by various Bayesian and Likelihoodist accounts. We suspended the impulse to deny it can make sense to use a rival inference school to critique significance tests. We sought to explore if there’s something to the cases they bring as ammunition to this conflict. The Bayesians say the disagreement between their numbers and  $P$ -values is relevant for impugning  $P$ -values, so we try to go along with them.

Reflect just on the first argument, pertaining to the case of two-sided Normal testing  $H_0: \mu = 0$  vs.  $H_0: \mu \neq 0$ , which was the most impressive, particularly with  $n \geq 50$ . It showed that a statistically significant difference from a test hypothesis

## 258 Excursion 4: Objectivity and Auditing

---

at familiar levels, 0.05 or 0.025, can correspond to a result that a Bayesian takes as evidence for  $H_0$ . The prior for this case is the spike and smear, where the smear will be of the sort leading to J. Berger and Sellke's results, or similar. The test procedure is to move from a statistically significant result at the 0.025 level, say, and infer the posterior for  $H_0$ .

Now our minimal requirement for data  $x$  to provide evidence for a claim  $H$  is that

(S-1)  $H$  accords with (agrees with)  $x$ , and

(S-2) there's a reasonable, preferably a high, probability that the procedure would have produced disagreement with  $H$ , if in fact  $H$  were false.

So let's apply these severity requirements to the data taken as evidence for  $H_0$  here.

Consider (S-1). Is a result that is 1.96 or 2 standard errors away from 0 in good accord with 0? Well, 0 is excluded from the corresponding 95% confidence interval. That does not seem to be in accord with 0 at all. Still, they have provided measures whereby  $x$  does accord with  $H_0$ , the likelihood ratio or posterior probability on  $H_0$ . So, in keeping with the most useful and most generous way to use severity, let's grant (S-1) holds.

What about (S-2)? Has anything been done to probe the falsity of  $H_0$ ? Let's allow that  $H_0$  is not a precise point, but some very small set of values around 0. This is their example, and we're trying to give it as much credibility as possible. Did the falsity of  $H_0$  have a good chance of showing itself? The falsity of  $H_0$  here is  $H_1: \mu \neq 0$ . What's troubling is that we found the probability of failing to pick up on population discrepancies as much as 1 standard error in excess of 0 is rather high (0.84) with  $n = 100$ . Larger sample sizes yield even less capability. Nor are they merely announcing "no discrepancy from 0" in this case. They're finding evidence for 0!

So how did the Bayesian get the bump in posterior probability on the null? It was based on a spiked prior of 0.5 to  $H_0$ . All the other points get minuscule priors having to share the remaining 0.5 probability. What was the warrant for the 0.5 prior to  $H_0$ ? J. Berger and Sellke are quite upfront about it: if they allowed the prior spike to be low, then a rejection of the null would merely be showing an improbable hypothesis got more improbable. "[W]ho, after all, would be convinced," recall their asking: if "my conclusion is that  $H_0$  has posterior probability 0.05 and should be rejected" since it previously had probability, say 0.1 (1987, p. 115). A slight lowering of probability won't cut it. Moving from a low prior to a slightly higher one also lacks punch.

---

**Tour II: Rejection Fallacies: Who's Exaggerating What?**259

---

This explains their high prior (at least 0.5) on  $H_0$ , but is it evidence for it? Clearly not, nor does it purport to be. We needn't deny there are cases where a theoretical parameter value has passed severely (we saw this in the case of GTR in Excursion 3). But that's not what's happening here. Here they intend for the 0.5 prior to show, *in general*, that statistically significant results problematically exaggerate evidence.<sup>8</sup>

A tester would be worried when the rationale for a spike is to avoid looking foolish when rejecting with a small drop; she'd be worried too by a report: "I don't take observing a mean temperature of 152 in your 100 water samples as indicating it's hotter than 150, because I give a whopping spike to our coolants being in compliance." That is why Casella and R. Berger describe J. Berger and Sellke's spike and smear as maximally biased toward the null (1987a, p. 111). Don't forget the powerful role played by the choice of how to smear the 0.5 over the alternative! Bayesians might reassure us that the high Bayes factor for a point null doesn't depend on the priors given to  $H_0$  and  $H_1$ , when what they mean is that it depends only on the priors given to discrepancies under  $H_1$ . It was the diffuse prior to the effect size that gave rise to the Jeffreys–Lindley Paradox. It affords huge latitude in what gets supported.

We thought we were traveling in testing territory; now it seems we've drifted off to a different place. It shouldn't be easy to take data as evidence for a claim when that claim is false; but here it is easy (the claim here being  $H_0$ ). How can this be one of a handful of main ways to criticize significance tests as exaggerating evidence? Bring in a navigator from a Popperian testing tribe before we all feel ourselves at sea:

Mere supporting instances are as a rule too cheap to be worth having . . . any support capable of carrying weight can only rest upon ingenious tests, undertaken with the aim of refuting our hypothesis, if it can be refuted. (Popper 1983, p. 130)

The high spike and smear tactic can't be taken as a basis from which to launch a critique of significance tests because it fails rather glaringly a minimum requirement for evidence, let alone a test. We met Bayesians who don't approve of these tests either, and I've heard it said that Bayesian testing is still a work in progress (Bernardo). Yet a related strategy is at the heart of some recommended statistical reforms.

<sup>8</sup> In the special case, where there's appreciable evidence for a special parameter, Senn argues that Jeffreys only required  $H_1$ 's posterior probability to be greater than 0.5. One has, so to speak, used up the prior belief by using the spiked prior (Senn 2015a).