

# Statistical Evidence in Experimental Psychology: An Empirical Comparison Using 855 $t$ Tests

Perspectives on Psychological Science  
6(3) 291–298

© The Author(s) 2011

Reprints and permission:

sagepub.com/journalsPermissions.nav

DOI: 10.1177/1745691611406923

http://pps.sagepub.com



Ruud Wetzels<sup>1</sup>, Dora Matzke<sup>1</sup>, Michael D. Lee<sup>2</sup>, Jeffrey N. Rouder<sup>3</sup>,  
Geoffrey J. Iverson<sup>2</sup>, and Eric-Jan Wagenmakers<sup>1</sup>

<sup>1</sup>Department of Psychology, University of Amsterdam, Amsterdam, The Netherlands; <sup>2</sup>Department of Cognitive Sciences, University of California, Irvine; and <sup>3</sup>Department of Psychological Sciences, University of Missouri-Columbia

## Abstract

Statistical inference in psychology has traditionally relied heavily on  $p$ -value significance testing. This approach to drawing conclusions from data, however, has been widely criticized, and two types of remedies have been advocated. The first proposal is to supplement  $p$  values with complementary measures of evidence, such as effect sizes. The second is to replace inference with Bayesian measures of evidence, such as the Bayes factor. The authors provide a practical comparison of  $p$  values, effect sizes, and default Bayes factors as measures of statistical evidence, using 855 recently published  $t$  tests in psychology. The comparison yields two main results. First, although  $p$  values and default Bayes factors almost always agree about what hypothesis is better supported by the data, the measures often disagree about the strength of this support; for 70% of the data sets for which the  $p$  value falls between .01 and .05, the default Bayes factor indicates that the evidence is only anecdotal. Second, effect sizes can provide additional evidence to  $p$  values and default Bayes factors. The authors conclude that the Bayesian approach is comparatively prudent, preventing researchers from overestimating the evidence in favor of an effect.

## Keywords

hypothesis testing,  $t$  test,  $p$  value, effect size, Bayes factor

Experimental psychologists use statistical procedures to convince themselves and their peers that the effect of interest is real, reliable, replicable, and hence worthy of academic attention. A representative example comes from Mussweiler (2006), who studied whether particular actions can activate a corresponding stereotype. To test this hypothesis empirically, Mussweiler unobtrusively induced half the participants, the experimental group, to move in a portly manner that is stereotypical for the overweight. The other half, the control group, made no such movements. Next, all participants were given an ambiguous description of a target person and then used a 9-point scale (ranging from 1 = *not at all* to 9 = *very*) to rate this person on dimensions that correspond to the overweight stereotype (e.g., “unhealthy,” “sluggish,” and “insecure”). To assess whether performing the stereotypic motion affected the rating of the ambiguous target person, Mussweiler computed a  $t$  statistic,  $t(18) = 2.1$ , and found that this value corresponded to a low  $p$  value ( $p < .05$ ).<sup>1</sup> Following conventional protocol, Mussweiler concluded that the low  $p$  value should be taken to provide “initial support for the hypothesis that

engaging in stereotypic movements activates the corresponding stereotype” (Mussweiler, 2006, p. 28).

The use of  $t$  tests and corresponding  $p$  values in this way constitutes a common and widely accepted practice in the psychological literature. It is, however, not the only possible or reasonable approach to measuring evidence and making statistical and scientific inferences. Indeed, the use of  $t$  tests and  $p$  values has been widely criticized (e.g., Cohen, 1994; Cumming, 2008; Dixon, 2003; Howard, Maxwell, & Flemming, 2000; Lee & Wagenmakers, 2005; Loftus, 1996; Nickerson, 2000; Wagenmakers, 2007). There are at least two different criticisms, coming from different perspectives and resulting in different remedies. First, many have argued that null hypothesis tests should be supplemented with other

## Corresponding Author:

Ruud Wetzels, Department of Psychology, University of Amsterdam, Roetersstraat 15, 1018 WB Amsterdam, The Netherlands  
E-mail: wetzels.ruud@gmail.com

statistical measures, such as confidence intervals and effect sizes. Within psychology, this approach to remediation has sometimes been institutionalized, being required by journal editors or recommended by the American Psychological Association (e.g., American Psychological Association, 2010; Cohen, 1988; Erdfelder, 2010; Wilkinson & the Task Force on Statistical Inference, 1999).

A second, more fundamental criticism that comes from Bayesian statistics is that there are basic conceptual and practical problems with  $p$  values. Although Bayesian criticism of psychological statistical practice dates back to at least Edwards, Lindman, and Savage (1963), it has become especially prominent and increasingly influential in the last decade (e.g., Dienes, 2008; Gallistel, 2009; Kruschke, 2010a, 2010c; Lee, 2008; Myung, Forster, & Browne, 2000; Rouder, Speckman, Sun, Morey, & Iverson, 2009). One standard Bayesian measure for quantifying the amount of evidence from the data in support of an experimental effect is the Bayes factor (Gönen, Johnson, Lu, & Westfall, 2005; Rouder et al., 2009; Wetzels, Raaijmakers, Jakab, & Wagenmakers, 2009). The measure takes the form of an odds ratio: It is the probability of the data under one hypothesis relative to that under another (Dienes, 2011; Kass & Raftery, 1995; Lee & Wagenmakers, 2005).

With this background, it seems that psychological statistical practice currently stands at a three-way fork in the road. Staying on the current path means continuing to rely on  $p$  values. A modest change is to place greater focus on the additional inferential information provided by effect sizes and confidence intervals. A radical change is struck by moving to Bayesian approaches, such as Bayes factors. The path that psychological science chooses seems likely to matter. It is not just that there are philosophical differences between the three choices. It is also clear that the three measures of evidence can be mutually inconsistent (e.g., Berger & Sellke, 1987; Rouder et al., 2009; Wagenmakers, 2007; Wagenmakers & Grünwald, 2006; Wagenmakers, Lodewyckx, Kuriyal, & Grasman, 2010).

In this article, we assess the practical consequences of choosing among inference by  $p$  values, by effect sizes, and by Bayes factors. By practical consequences, we mean the extent to which conclusions of extant studies change according to the inference measure that is used. To assess these practical consequences, we reanalyzed 855  $t$  tests reported in articles from the 2007 issues of *Psychonomic Bulletin & Review (PBR)* and *Journal of Experimental Psychology: Learning, Memory, and Cognition (JEP:LMC)*. For each  $t$  test, we compute the  $p$  value, the effect size, and the Bayes factor and study the extent to which they provide information that is redundant, complementary, or inconsistent. On the basis of these analyses, we suggest the best direction for measuring statistical evidence from psychological experiments.

### Three Measures of Evidence

In this section, we describe how to calculate and interpret the  $p$  value, the effect size, and the Bayes factor. For concreteness, we use Mussweiler's (2006) study on the effect of action on

stereotypes. The mean score of the control group,  $M_c$ , was 5.8 on a weight-stereotype scale ( $s_c = 0.69$ ,  $n_c = 10$ ), and the mean score of the experimental group,  $M_e$ , was 6.4 ( $s_e = 0.66$ ,  $n_e = 10$ ).

### The $p$ value

The interpretation of  $p$  values is not straightforward, and their use in hypothesis testing is heavily debated (Cohen, 1994; Cortina & Dunlap, 1997; Cumming, 2008; Dixon, 2003; Frick, 1996; Gigerenzer, 1993, 1998; Hagen, 1997; Killeen, 2005, 2006; Kruschke, 2010a, 2010c; Lee & Wagenmakers, 2005; Loftus, 1996; Nickerson, 2000; Schmidt, 1996; Wagenmakers & Grünwald, 2006; Wainer, 1999). The  $p$  value is the probability of obtaining a test statistic (in this case, the  $t$  statistic) at least as extreme as the one that was observed in the experiment, given that the null hypothesis is true and the sample is generated according to a specific intended procedure, such as fixed sample size. Fisher (1935) interpreted these  $p$  values as evidence against the null hypothesis. The smaller the  $p$  value, the more evidence there was against the null hypothesis. Fisher viewed these values as self-explanatory measures of evidence that did not need further guidance. In practice, however, most researchers (and reviewers) adopt a .05 cutoff:  $p$  values less than .05 constitute evidence for an effect, and those greater than .05 do not. More fine-grained categories are possible, and Wasserman (2004, p. 157) proposes the gradations shown in the top of Table 1. Note that the top part of Table 1 lists various categories of evidence against the null hypothesis. A basic limitation of null hypothesis significance testing is that it does not allow a researcher to gather evidence in favor of the null (Dennis, Lee, & Kinnell, 2008; Gallistel, 2009; Rouder et al., 2009; Wetzels et al., 2009).

For the data from Mussweiler (2006), we compute a  $p$  value based on the  $t$  test. The  $t$  test is designed to test whether a difference between two means is significant. First, we calculate the  $t$  statistic:

$$t = \frac{M_e - M_c}{\sqrt{s_{pooled}^2 \left( \frac{1}{n_e} + \frac{1}{n_c} \right)}} = \frac{6.42 - 5.79}{\sqrt{0.46 \left( \frac{1}{10} + \frac{1}{10} \right)}} = 2.09$$

where  $M_e$  and  $M_c$  are the means of both groups,  $n_e$  and  $n_c$  are the sample sizes, and  $s_{pooled}^2$  estimates the common population variance:

$$s_{pooled}^2 = \frac{(n_e - 1)S_e^2 + (n_c - 1)S_c^2}{n_e + n_c - 2}$$

Next, the  $t$  statistic with  $n_e + n_c - 2 = 18$  degrees of freedom results in a  $p$  value slightly larger than .05 ( $\approx .051$ ). For our concrete example, Table 1 leads to the conclusion that the  $p$  value is on the cusp between "no evidence against  $H_0$ " and "positive evidence against  $H_0$ ."

**Table 1.** Evidence Categories for  $p$  Values (adapted from Wasserman, 2004, p. 157), for Effect Sizes (as proposed by Cohen, 1988), and for Bayes Factor  $BF_{A0}$  (Jeffreys, 1961)

Statistic	Interpretation
$p$ value	
<.001	Decisive evidence against $H_0$
.001–.01	Substantive evidence against $H_0$
.01–.05	Positive evidence against $H_0$
>.05	No evidence against $H_0$
Effect size	
<0.2	Small effect size
0.2–0.5	Small to medium effect size
0.5–0.8	Medium to large effect size
0.8	Large to very large effect size
Bayes factor	
>100	Decisive evidence for $H_A$
30–100	Very strong evidence for $H_A$
10–30	Strong evidence for $H_A$
3–10	Substantial evidence for $H_A$
1–3	Anecdotal evidence for $H_A$
1	No evidence
1/3–1	Anecdotal evidence for $H_0$
1/10–1/3	Substantial evidence for $H_0$
1/30–1/10	Strong evidence for $H_0$
1/100–1/30	Very strong evidence for $H_0$
<1/100	Decisive evidence for $H_0$

Note: For the Bayes factor categories, we replaced the label “worth no more than a bare mention” with “anecdotal.” Also, in contrast to  $p$  values, the Bayes factor can quantify evidence in favor of the null hypothesis.

### The effect size

Effect sizes quantify the magnitude of an effect and serve as a measure of how much the results deviate from the null hypothesis (Cohen, 1988; Richard, Bond, & Stokes-Zoota, 2003; Rosenthal, 1990; Rosenthal & Rubin, 1982; Thompson, 2002). For the data from Mussweiler (2006), the effect size,  $d$ , is calculated as follows:

$$d = \frac{M_e - M_c}{S_{pooled}} = \frac{6.42 - 5.79}{0.68} = 0.93$$

Note that in contrast to the  $p$  value, the effect size is independent of sample size; increasing the sample size does not increase effect size but instead allows it to be estimated more accurately.

Effect sizes are often interpreted in terms of the categories introduced by Cohen (1988), as listed in the middle of Table 1, ranging from “small” to “very large.” For our concrete example,  $d = 0.93$ , and we conclude that this effect is large to very large. Interestingly, the  $p$  value was on the cusp between the categories “no evidence against  $H_0$ ” and “positive evidence against  $H_0$ ,” whereas the effect size indicates the effect to be strong.

### The Bayes factor

In Bayesian statistics, uncertainty (or degree of belief) is quantified by probability distributions over parameters. This makes the Bayesian approach fundamentally different from the

classical “frequentist” approach, which relies on sampling distributions of data (Berger & Delampady, 1987; Berger & Wolpert, 1988; Jaynes, 2003; Lindley, 1972).

Within the Bayesian framework, one may quantify the evidence for one hypothesis relative to another. The Bayes factor is the most commonly used (although certainly not the only possible) Bayesian measure for doing so (Jeffreys, 1961; Kass & Raftery, 1995). The Bayes factor is the probability of the data under one hypothesis relative to the other. When a hypothesis is a simple point, such as the null, then the probability of the data under this hypothesis is simply the likelihood evaluated at that point. When a hypothesis consists of a range of points, such as all positive effect sizes, then the probability of the data under this hypothesis is the weighted average of the likelihood across that range. This averaging automatically controls for the complexity of different models, as has been emphasized in Bayesian literature in psychology (e.g., Pitt, Myung, & Zhang, 2002; Rouder et al., 2009).

We take as the null that a parameter  $\alpha$  is restricted to 0 (i.e.,  $H_0: \alpha = 0$ ), and we take as the alternative that  $\alpha$  is not zero (i.e.,  $H_A: \alpha \neq 0$ ). In this case, the Bayes factor given data  $D$  is simply the ratio where the integral in the denominator takes the average evidence over all values of  $\alpha$ , weighted by the prior probability of those values  $p(\alpha | H_A)$  under the alternative hypothesis.

An alternative—but formally equivalent—conceptualization of the Bayes factor is

$$BF_{A0} = \frac{p(D|H_A)}{p(D|H_0)} = \frac{\int p(D|H_A, \alpha)p(\alpha|H_A)d\alpha}{p(D|H_0)},$$

as a measure of the change from prior model odds to posterior model odds, brought about by the observed data. This change is often interpreted as the *weight of evidence* (Good, 1983, 1985). Before seeing the data  $D$ , the two hypotheses  $H_0$  and  $H_A$  are assigned prior probabilities  $p(H_0)$  and  $p(H_A)$ . The ratio of the two prior probabilities defines the *prior odds*. When the data  $D$  are observed, the prior odds are updated to *posterior odds*, which is defined as the ratio of the posterior probabilities,  $p(H_0 | D)$  and  $p(H_A | D)$ :

$$\frac{p(H_A|D)}{p(H_0|D)} = \frac{p(D|H_A)}{p(D|H_0)} \times \frac{p(H_A)}{p(H_0)}. \quad (1)$$

Equation 1 shows that the change from prior odds to posterior odds is quantified by  $p(D|H_A)/p(D|H_0)$ : the Bayes factor,  $BF_{A0}$ .

Under either conceptualization, the Bayes factor has an appealing and direct interpretation as an odds ratio. For example,  $BF_{A0} = 2$  implies that the data are twice as likely to have occurred under  $H_A$  than under  $H_0$ . Jeffreys (1961) proposed a set of verbal labels to categorize the Bayes factor according to its evidential impact. This set of labels, presented at the bottom of Table 1, facilitates scientific communication but should only be considered an approximate descriptive articulation of different standards of evidence (Kass & Raftery, 1995).

In general, calculating Bayes factors is more difficult than calculating  $p$  values and effect sizes. However, psychologists can now turn to easy-to-use Web pages to calculate the Bayes

factor for many common experimental situations or use software such as WinBUGS (Lunn, Thomas, Best, & Spiegelhalter, 2000; Wetzels, Lee, & Wagenmakers, 2010; Wetzels et al., 2009).<sup>2</sup> In this article, we use the Bayes factor calculation described in Rouder et al. (2009). Rouder et al.'s development is suitable for one-sample and two-sample designs, and the only necessary input is the  $t$  value and sample size.

The Bayes factor that we report in this article is the result of a default Bayesian  $t$  test (for details, see Rouder et al., 2009). The test is default because it applies regardless of the phenomenon under study: For every experiment, one uses the same prior distribution on effect size for the alternative hypothesis, the Cauchy (0,1) distribution. This prior distribution has statistical advantages that make it an appropriate default choice (for example, it has excellent theoretical properties in the limit,  $N \rightarrow \infty$  and  $t \rightarrow \infty$ ; for details, see Liang, Paulo, Molina, Clyde, & Berger, 2008).

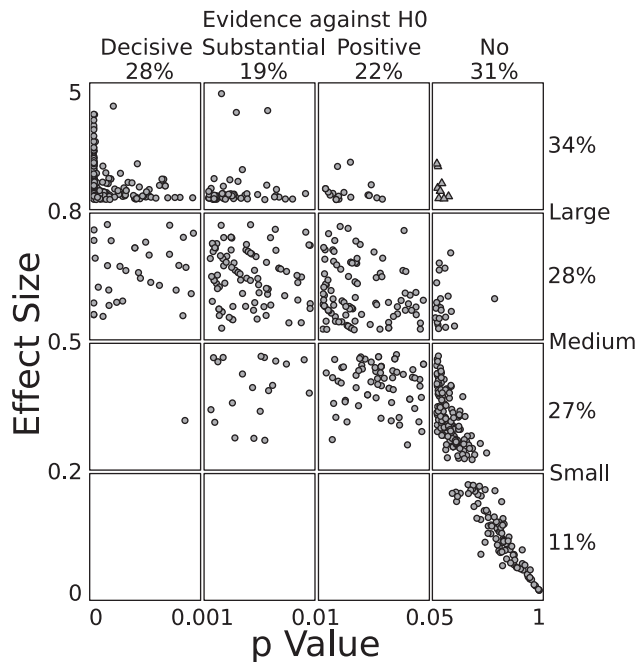
The default test is easy to use and avoids informed specification of prior distributions that other researchers may contest. Conversely, one may argue that the informed specification of priors is the appropriate way to take problem-specific prior knowledge into account. Bayesian statisticians are divided over the relative merits of default versus informed specifications of prior distributions (Press, Chib, Clyde, Woodworth, & Zaslavsky, 2003). In our opinion, the default test provides an excellent starting point of analysis, one that may later be supplemented with a detailed problem-specific analysis (see Dienes, 2008, 2011, this issue; Kruschke, 2010a, 2010b, 2011, this issue, for additional discussion of informed priors).

In our concrete example, the resulting Bayes factor for  $t = 2.09$  and a sample size of 20 observations is  $BF_{A0} = 1.56$ . Accordingly, the data are 1.56 times more likely to have occurred under the alternative hypothesis than under the null hypothesis. This Bayes factor falls into the category “anecdotal.” In other words, this Bayes factor indicates that although the alternative hypothesis is slightly favored, we do not have sufficiently strong evidence from the data to reject or accept either hypothesis.

## Comparing $p$ Values, Effect Sizes, and Bayes Factors

For our concrete example, the three measures of evidence are not in agreement. The  $p$  value was on the cusp between the categories “no evidence against  $H_0$ ” and “positive evidence against  $H_0$ ,” the effect size indicates a large to very large effect size, and the Bayes factor indicates that the data support the null hypothesis almost as much as they support the alternative hypothesis. If this example is not an isolated one, and the measures differ in many psychological applications, then it is important to understand the nature of those differences.

To address this question, we studied all of the empirical results evaluated by a  $t$  test in the 2007 volumes of *PBR* and *JEP:LMC*. This sample was composed of 855  $t$  tests from 252 articles. These articles covered 2,394 journal pages and addressed many topics that are important in modern experimental psychology. Our sample suggests, on average, that an article published in *PBR* and



**Fig. 1.** The relationship between effect size and  $p$  values. Points denote comparisons (855 in total). Points denoted by circles indicate relative consistency between the effect size and  $p$  value, whereas those denoted by triangles indicate gross inconsistencies. The scale of the axes is based on the decision categories, as given in Table 1.

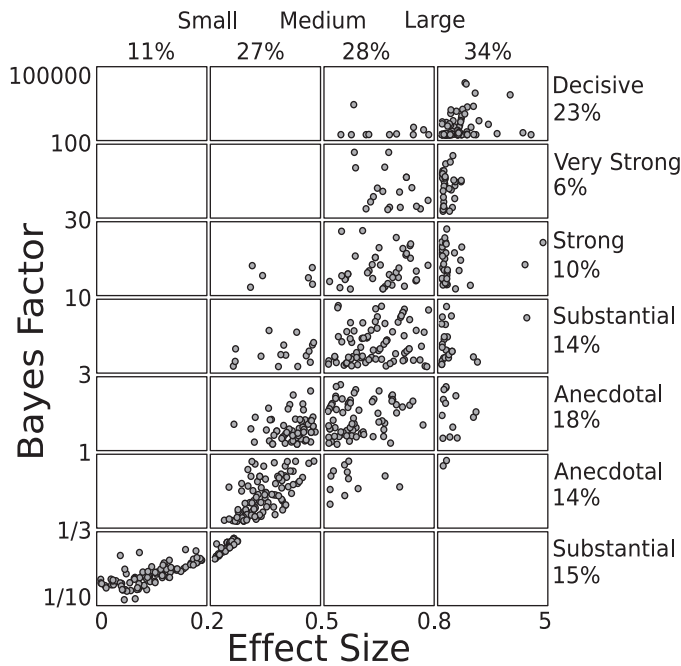
*JEP:LMC* contains about 3.4  $t$  tests, which amounts to one  $t$  test for every 2.8 pages. For simplicity, we did not include  $t$  tests that resulted from multiple comparisons in analysis of variance designs (for a Bayesian perspective on multiple comparisons, see Scott and Berger, 2006). Even though our  $t$  tests are sampled from the field of experimental and cognitive psychology, we expect our findings to generalize to many other subfields of psychology, as long as the studies in these subfields use the same level of statistical significance, approximately the same number of participants, and approximately the same number of trials per participant (Howard et al., 2000).

In the next sections, we describe the empirical relation between the three measures of evidence, starting with the relation between effect sizes and  $p$  values.

## Comparing effect sizes and $p$ values

The relationship between the obtained  $p$  values and effect sizes is shown as a scatter plot in Figure 1. Each point corresponds to one of the 855 comparisons. Different panels are introduced to distinguish the different evidence categories, as given in Table 1.

Figure 1 suggests that  $p$  values and effect sizes capture roughly the same information in the data. Large effect sizes tend to correspond to low  $p$  values, and small effect sizes tend to correspond to large  $p$  values. The two measures, however, are far from identical. For instance, a  $p$  value of .01 can correspond to effect sizes ranging from about 0.2 to 1, and an effect size near 0.5 can correspond to  $p$  values ranging from about



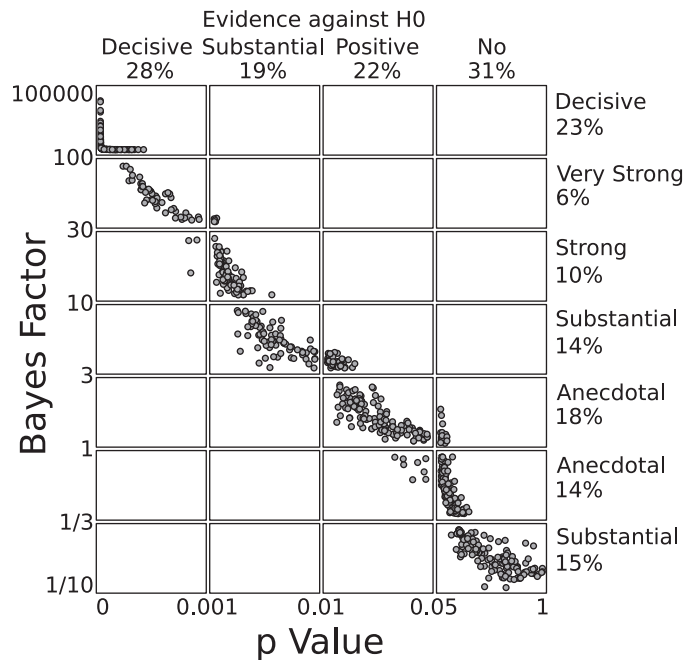
**Fig. 2.** The relationship between Bayes factor and effect size. Points denote comparisons (855 in total). The scale of the axes is based on the decision categories, as given in Table 1.

.001 to .05. The triangular points in the top-right panel of Figure 1 highlight gross inconsistencies. These eight studies have a large effect size, above 0.8, but their *p* values do not indicate evidence against the null hypothesis. A closer examination revealed that these studies had *p* values very close to .05 and were comprised of small sample sizes.

**Comparing effect sizes and Bayes factors**

The relationship between the obtained Bayes factors and effect sizes is shown in Figure 2. Much as with the comparison of *p* values with effect sizes, it seems clear that the default Bayes factor and effect size generally agree, though not exactly. No striking inconsistencies are apparent: No study with an effect size greater than 0.8 coincides with a Bayes factor below 1/3, nor does a study with very low effect size below 0.2 coincide with a Bayes factor above 3. The two measures, however, are not identical. They differ in the assessment of strength of evidence. Effect sizes above 0.8 range all the way from anecdotal to decisive evidence in terms of the Bayes factor. Also note that small to medium effect sizes (i.e., those between 0.2 and 0.5) can correspond to Bayes factor evidence in favor of either the alternative or the null hypothesis.

This last observation supports the premise that Bayes factors may quantify support for the null hypothesis. Figure 2 shows that about one-third of all studies produced evidence in favor of the null hypothesis. In about half of these studies favoring the null, the evidence is substantial. Because of the file-drawer problem (i.e., only significant effects tend to get published), this is an underestimate of the true number of null findings and their Bayes factor support.



**Fig. 3.** The relationship between Bayes factor and *p* value. Points denote comparisons (855 in total). The scale of the axes is based on the decision categories, as given in Table 1.

**Comparing *p* values and Bayes factors**

The relationship between the obtained Bayes factors and *p* values is shown in Figure 3, again using interpretative panels. It is clear that default Bayes factors and *p* values largely covary with each other. Low Bayes factors correspond to high *p* values, and high Bayes factors correspond to low *p* values, a relationship that is much more exact than for our previous two comparisons. The main difference between default Bayes factors and *p* values is one of calibration; *p* values accord more evidence against the null than do Bayes factors. Consider the *p* values between .01 and .05, values that correspond to “positive evidence” and that usually pass the bar for publishing in academia. According to the default Bayes factor, 70% of these experimental effects convey evidence in favor of the alternative hypothesis that is only “anecdotal.” This difference in the assessment of the strength of evidence is dramatic and consequential.

**Conclusion**

We compared *p* values, effect sizes, and default Bayes factors as measures of statistical evidence in empirical psychological research. Our comparison was based on a total of 855 different *t* statistics from all published articles in two major empirical journals in 2007. In virtually all studies, the three different measures of evidence are broadly consistent: Small *p* values correspond to large effect sizes and large Bayes factors in favor of the alternative hypothesis. Despite the fact that the measures of evidence reach the same conclusion about what hypothesis is best supported by the data, however, the measures differ with respect to the strength of that support. In particular, we noted

that  $p$  values between .01 and .05 often correspond to what, in Bayesian terms, is only anecdotal evidence in favor of the alternative hypothesis. The practical ramifications of this are considerable.

### **Practical ramifications**

Our results showed that when the  $p$  value falls in the interval from .01 to .05, there is a 70% chance that the default Bayes factor indicates the evidence for the alternative hypothesis to be only anecdotal or “worth no more than a bare mention”; this means that the data are no more than three times more likely under the alternative hypothesis than they are under the null hypothesis. Hence, for the studies under consideration here, it seems that a  $p$ -value criterion more conservative than .05 is appropriate. Alternatively, researchers could avoid computing a  $p$  value altogether and instead compute the Bayes factor. Both methods help prevent researchers from overestimating the strength of their findings and help keep the field from incorporating ambiguous findings as if these were real and reliable (Ioannidis, 2005).

As a practical illustration, consider a series of recent experiments on precognition (Bem, 2011). In nine experiments with over 1,000 participants, Bem intended to show that precognition exists, that is, that people can foresee the future. And indeed, eight out of nine experiments yielded a significant result. However, most  $p$  values fell in the ambiguous range of .01 to .05, and across all nine experiments, a Bayes factor analysis indicates about as much evidence for the alternative hypothesis as against it (Kruschke, 2011; Wagenmakers, Wetzels, Borsboom, & van der Maas, 2011). We believe that this situation typifies part of what could be improved in psychological research today. It is simply too easy to obtain a  $p$  value below .05 and to subsequently publish the result.

When researchers publish ambiguous results as if they were real and reliable, this damages the field as a whole: Time, effort, and money will be invested to replicate the phenomenon, and when replication fails, the burden of proof is almost always on the part of the researcher who, after all, failed to replicate a phenomenon that was demonstrated to be present (with a  $p$  value between .01 and .05).

Thus, our empirical comparison shows that the academic criterion of .05 is too liberal. Note that this problem would not be solved by opting for a stricter significance level, such as .01. It is well known that the  $p$  value decreases as the sample size,  $n$ , increases. Hence, if psychologists switch to a significance level of .01 but inevitably increase their sample sizes to compensate for the stricter statistical threshold, then the phenomenon of anecdotal evidence will start to plague  $p$  values even when these  $p$  values are lower than .01. Therefore, we make a case for Bayesian statistics in the next section.

### **A case for Bayesian statistics**

We have compared the conclusions from the different measures of evidence. It is easy to make a case for Bayesian statistical

inference in general, based on arguments already well documented in statistics and psychology (e.g., Dienes, 2008; Jaynes, 2003; Kruschke, 2010a, 2010c; Lee & Wagenmakers, 2005; Lindley, 1972; Wagenmakers, 2007). We briefly mention three arguments here.

First, unlike null hypothesis testing, Bayesian inference does not violate basic principles of rational statistical decision making, such as the stopping rule principle or the likelihood principle (Berger & Delampady, 1987; Berger & Wolpert, 1988; Dienes, 2011). This means that the results of Bayesian inference do not depend on the intention with which the data were collected. As stated by Edwards et al. (1963, p. 193), “the rules governing when data collection stops are irrelevant to data interpretation. It is entirely appropriate to collect data until a point has been proven or disproven, or until the data collector runs out of time, money, or patience.”

Second, Bayesian inference takes model complexity into account in a rational way. Specifically, the Bayes factor has the attraction of not assigning a special status to the null hypothesis and so makes it theoretically possible to measure evidence in favor of the null (e.g., Dennis et al., 2008; Gallistel, 2009; Kass & Raftery, 1995; Rouder et al., 2009).

Third, we believe that Bayesian inference provides the kind of answers that researchers care about. In our experience, researchers are usually not that interested in the probability of encountering data at least as extreme as those that were observed, given that the null hypothesis is true and the sample was generated according to a specific intended procedure. Instead, most researchers want to know what they have learned from the data about the relative plausibility of the hypotheses under consideration. This is exactly what is quantified by the Bayes factor.

These advantages notwithstanding, the Bayes factor is not a measure of the mere size of an effect. Hence, the measure of effect size confers additional information, particularly when small numbers of participants or trials are involved. So, especially for these sorts of studies, there is an argument for reporting both a Bayes factor and an effect size. We note that, from a Bayesian perspective, the effect size can naturally be conceived as (a summary statistic of) the posterior distribution of a parameter representing the effect, under an uninformative prior distribution. In this sense, a standard Bayesian combination of parameter estimation and model selection could encompass all of the useful measures of evidence we observed (for an example of how Bayes factor estimation can be incorporated in a Bayesian estimation framework, see, for instance, Kruschke, 2011).

Our final thought is that reasons for adopting a Bayesian approach now are amplified by the promise of using an extended Bayesian approach in the future. In particular, we think the hierarchical Bayesian approach, which is standard in statistics (e.g., Gelman & Hill, 2007) and is becoming more common in psychology (e.g. Kruschke, 2010b, 2010c; Lee, in press; Rouder & Lu, 2005), could fundamentally change how psychologists identify effects. Hierarchical Bayesian analysis can be a valuable tool both for meta-analyses and for the

analysis of a single study. In the meta-analytical context, multiple studies can be integrated, so that what is inferred about the existence of effects and their magnitude is informed, in a coherent and quantitative way, by a domain of experiments. In the context of a single experiment, a hierarchical analysis can be used to take variability across participants or items into account.

In sum, our empirical comparison of 855  $t$  tests shows that three often-used measures of evidence— $p$  values, effect sizes, and Bayes factors—almost always agree about what hypothesis is better supported by the data. However, the measures often disagree about the strength of this support: for those data sets with  $p$  values in between .01 and .05, about 70% are associated with a Bayes factor that indicates the evidence to be only anecdotal or “worth no more than a bare mention” (Jeffreys, 1961). This analysis suggests that many results that have been published in the literature are not established as strongly as one would like.

### Declaration of Conflicting Interests

The authors declared that they had no conflicts of interest with respect to their authorship or the publication of this article.

### Funding

This research was supported by a Vidi grant from the Netherlands Organization for Scientific Research.

### Notes

1. The findings suggest that Mussweiler (2006) conducted a one-sided  $t$  test. In the remainder of this article, we conduct two-sided  $t$  tests.
2. A Web page for computing a Bayes factor online is <http://pcl.missouri.edu/bayesfactor>, and a Web page to download a tutorial and a flexible R/WinBUGS function to calculate the Bayes factor can be found at <http://www.ruudwetzels.com>.

### References

- American Psychological Association. (2010). *Publication manual of the American Psychological Association* (6th ed.). Washington, DC: Author.
- Bem, D.J. (2011). Feeling the future: Experimental evidence for anomalous retroactive influences on cognition and affect. *Journal of Personality and Social Psychology, 100*, 407–425.
- Berger, J.O., & Delampady, M. (1987). Testing precise hypotheses. *Statistical Science, 2*, 317–352.
- Berger, J.O., & Sellke, T. (1987). Testing a point null hypothesis: The irreconcilability of  $p$  values and evidence. *Journal of the American Statistical Association, 82*, 112–139.
- Berger, J.O., & Wolpert, R.L. (1988). *The likelihood principle* (2nd ed.). Hayward, CA: Institute of Mathematical Statistics.
- Cohen, J. (1988). *Statistical power analysis for the behavioral sciences*. Hillsdale, NJ: Erlbaum.
- Cohen, J. (1994). The earth is round ( $p < .05$ ). *American Psychologist, 49*, 997–1003.
- Cortina, J.M., & Dunlap, W.P. (1997). On the logic and purpose of significance testing. *Psychological Methods, 2*, 161–172.
- Cumming, G. (2008). Replication and  $p$  intervals:  $p$  values predict the future only vaguely, but confidence intervals do much better. *Perspectives on Psychological Science, 3*, 286–300.
- Dennis, S., Lee, M., & Kinnell, A. (2008). Bayesian analysis of recognition memory: The case of the list-length effect. *Journal of Memory and Language, 59*, 361–376.
- Dienes, Z. (2008). *Understanding psychology as a science: An introduction to scientific and statistical inference*. New York: Palgrave Macmillan.
- Dienes, Z. (2011). Bayesian versus orthodox statistics: Which side are you on? *Perspectives on Psychological Science, 6*, 274–290.
- Dixon, P. (2003). The  $p$ -value fallacy and how to avoid it. *Canadian Journal of Experimental Psychology, 57*, 189–202.
- Edwards, W., Lindman, H., & Savage, L.J. (1963). Bayesian statistical inference for psychological research. *Psychological Review, 70*, 193–242.
- Erdfelder, E. (2010). A note on statistical analysis. *Experimental Psychology, 57*, 1–4.
- Fisher, R.A. (1935). *The design of experiments*. Edinburgh: Oliver and Boyd.
- Frick, R.W. (1996). The appropriate use of null hypothesis testing. *Psychological Methods, 1*, 379–390.
- Gallistel, C. (2009). The importance of proving the null. *Psychological Review, 116*, 439–453.
- Gelman, A., & Hill, J. (2007). *Data analysis using regression and multilevel/hierarchical models*. Cambridge, England: Cambridge University Press.
- Gigerenzer, G. (1993). The Superego, the ego, and the id in statistical reasoning. In G. Keren & C. Lewis (Eds.), *A handbook for data analysis in the behavioral sciences: Methodological issues* (pp. 311–339). Hillsdale, NJ: Erlbaum.
- Gigerenzer, G. (1998). We need statistical thinking, not statistical rituals. *Behavioral and Brain Sciences, 21*, 199–200.
- Gönen, M., Johnson, W.O., Lu, Y., & Westfall, P.H. (2005). The Bayesian two-sample  $t$  test. *American Statistician, 59*, 252–257.
- Good, I.J. (1983). *Good thinking: The foundations of probability and its applications*. Minneapolis: University of Minnesota Press.
- Good, I.J. (1985). Weight of evidence: A brief survey. In J.M. Bernardo, M.H. DeGroot, D.V. Lindley, & A.F.M. Smith (Eds.), *Bayesian statistics 2* (pp. 249–269). New York: Elsevier.
- Hagen, R.L. (1997). In praise of the null hypothesis statistical test. *American Psychologist, 52*, 15–24.
- Howard, G., Maxwell, S., & Fleming, K. (2000). The proof of the pudding: An illustration of the relative strengths of null hypothesis, meta-analysis, and Bayesian analysis. *Psychological Methods, 5*, 315–332.
- Ioannidis, J.P.A. (2005). Why most published research findings are false. *PLoS Medicine, 2*, 696–701.
- Jaynes, E.T. (2003). *Probability theory: The logic of science*. Cambridge, UK: Cambridge University Press.
- Jeffreys, H. (1961). *Theory of probability*. Oxford, UK: Oxford University Press.
- Kass, R.E., & Raftery, A.E. (1995). Bayes factors. *Journal of the American Statistical Association, 90*, 377–395.
- Killeen, P.R. (2005). An alternative to null-hypothesis significance tests. *Psychological Science, 16*, 345–353.

- Killeen, P.R. (2006). Beyond statistical inference: A decision theory for science. *Psychonomic Bulletin & Review*, *13*, 549–562.
- Kruschke, J.K. (2010a). Bayesian data analysis. *Wiley Interdisciplinary Reviews: Cognitive Science*, *1*, 658–676.
- Kruschke, J.K. (2010b). *Doing Bayesian data analysis: A tutorial introduction with R and BUGS*. Burlington, MA: Academic Press.
- Kruschke, J.K. (2010c). What to believe: Bayesian methods for data analysis. *Trends in Cognitive Sciences*, *14*, 293–300.
- Kruschke, J.K. (2011). Bayesian assessment of null values via parameter estimation and model comparison. *Perspectives on Psychological Science*, *6*, 299–312.
- Lee, M.D. (2008). Three case studies in the Bayesian analysis of cognitive models. *Psychonomic Bulletin & Review*, *15*, 1–15.
- Lee, M.D. (in press). How cognitive modeling can benefit from hierarchical Bayesian models. *Journal of Mathematical Psychology*.
- Lee, M.D., & Wagenmakers, E.-J. (2005). Bayesian statistical inference in psychology: Comment on Trafimow (2003). *Psychological Review*, *112*, 662–668.
- Liang, F., Paulo, R., Molina, G., Clyde, M., & Berger, J. (2008). Mixtures of  $g$  priors for Bayesian variable selection. *Journal of the American Statistical Association*, *103*, 410.
- Lindley, D.V. (1972). *Bayesian statistics; a review*. Philadelphia: Society for Industrial and Applied Mathematics.
- Loftus, G.R. (1996). Psychology will be a much better science when we change the way we analyze data. *Current Directions in Psychological Science*, *5*, 161–171.
- Lunn, D.J., Thomas, A., Best, N., & Spiegelhalter, D. (2000). WinBUGS—a Bayesian modelling framework: Concepts, structure, and extensibility. *Statistics and Computing*, *10*, 325–337.
- Mussweiler, T. (2006). Doing is for thinking! *Psychological Science*, *17*, 17–21.
- Myung, I.J., Forster, M.R., & Browne, M.W. (2000). A special issue on model selection. *Journal of Mathematical Psychology*, *44*.
- Nickerson, R.S. (2000). Null hypothesis statistical testing: A review of an old and continuing controversy. *Psychological Methods*, *5*, 241–301.
- Pitt, M.A., Myung, I.J., & Zhang, S. (2002). Toward a method of selecting among computational models of cognition. *Psychological Review*, *109*, 472–491.
- Press, S., Chib, S., Clyde, M., Woodworth, G., & Zaslavsky, A. (2003). *Subjective and objective Bayesian statistics: Principles, models, and applications*. Hoboken, NJ: Wiley-Interscience.
- Richard, F.D., Bond, C.F.J., & Stokes-Zoota, J.J. (2003). One hundred years of social psychology quantitatively described. *Review of General Psychology*, *7*, 331–363.
- Rosenthal, R. (1990). How are we doing in soft psychology? *American Psychologist*, *45*, 775–777.
- Rosenthal, R., & Rubin, D. (1982). A simple, general purpose display of magnitude of experimental effect. *Journal of Educational Psychology*, *74*, 166–169.
- Rouder, J.N., & Lu, J. (2005). An introduction to Bayesian hierarchical models with an application in the theory of signal detection. *Psychonomic Bulletin & Review*, *12*, 573–604.
- Rouder, J.N., Speckman, P.L., Sun, D., Morey, R.D., & Iverson, G. (2009). Bayesian  $t$  tests for accepting and rejecting the null hypothesis. *Psychonomic Bulletin & Review*, *16*, 225–237.
- Schmidt, F.L. (1996). Statistical significance testing and cumulative knowledge in psychology: Implications for training of researchers. *Psychological Methods*, *1*, 115–129.
- Scott, J., & Berger, J. (2006). An exploration of aspects of Bayesian multiple testing. *Journal of Statistical Planning and Inference*, *136*, 2144–2162.
- Thompson, B. (2002). What future quantitative social science research could look like: Confidence intervals for effect sizes. *Educational Researcher*, *31*, 25–32.
- Wagenmakers, E.-J. (2007). A practical solution to the pervasive problems of  $p$  values. *Psychonomic Bulletin & Review*, *14*, 779–804.
- Wagenmakers, E.-J., & Grünwald, P. (2006). A Bayesian perspective on hypothesis testing. *Psychological Science*, *17*, 641–642.
- Wagenmakers, E.-J., Lodewyckx, T., Kuriyal, H., & Grasman, R. (2010). Bayesian hypothesis testing for psychologists: A tutorial on the Savage-Dickey method. *Cognitive Psychology*, *60*, 158–189.
- Wagenmakers, E.-J., Wetzels, R., Borsboom, D., & van der Maas, H.L.J. (2011). Why psychologists must change the way they analyze their data: The case of  $\psi$ : Comment on Bem (2011). *Journal of Personality and Social Psychology*, *100*, 426–432.
- Wainer, H. (1999). One cheer for null hypothesis significance testing. *Psychological Methods*, *4*, 212–213.
- Wasserman, L. (2004). *All of statistics: A concise course in statistical inference*. New York: Springer.
- Wetzels, R., Lee, M., & Wagenmakers, E.-J. (2010). Bayesian inference using WDev: A tutorial for social scientists. *Behavior Research Methods*, *42*, 884–897.
- Wetzels, R., Raaijmakers, J., Jakab, E., & Wagenmakers, E.-J. (2009). How to quantify support for and against the null hypothesis: A flexible WinBUGS implementation of a default Bayesian  $t$  test. *Psychonomic Bulletin & Review*, *16*, 752–760.
- Wilkinson, L., & the Task Force on Statistical Inference. (1999). Statistical methods in psychology journals: Guidelines and explanations. *American Psychologist*, *54*, 594–604.