

# How to Quantify the Evidence for the Absence of an Association Between Loneliness and Bathing Habits

Eric-Jan Wagenmakers, Josine Verhagen, and Alexander Ly  
University of Amsterdam

Correspondence concerning this article should be addressed to:

Eric-Jan Wagenmakers  
University of Amsterdam, Department of Psychology  
Weesperplein 4  
1018 XA Amsterdam, The Netherlands  
E-mail may be sent to EJ.Wagenmakers@gmail.com.

## Abstract

We present a suite of Bayes factor hypothesis tests that allow researchers to grade the decisiveness of the evidence that the data provide for the presence versus the absence of a correlation between two variables. We apply our methods to the recent work of Donnellan, Lucas, and Cesario (in press) who conducted nine replication studies with over 3,000 participants and failed to replicate the phenomenon that lonely people compensate for a lack of social warmth by taking warmer baths or showers. We show how the Bayes factor hypothesis test can quantify evidence in favor of the null hypothesis, and how the prior specification for the correlation coefficient can be used to define a broad range of tests that address complementary questions. Specifically, we show how the prior specification can be adjusted to create a two-sided test, a one-sided test, a sensitivity analysis, and a replication test.

**Keywords:** Hypothesis test; Statistical evidence; Bayes factor.

After a Herculean effort involving a series of nine replication experiments, Donnellan et al. (in press) ultimately failed to reject the null hypothesis that people do not use warm showers and baths to compensate for a lack of social warmth, contradicting an earlier claim by Bargh and Shalev (2012). Unfortunately, the standard  $p$  value methodology does not allow one to quantify evidence in favor of the null hypothesis (Gallistel, 2009; Rouder, Speckman, Sun,

---

We thank Brent Donnellan and Joe Cesario for sending us the data from their nine replication studies. This work was supported by an ERC grant from the European Research Council. Correspondence concerning this article may be addressed to Eric-Jan Wagenmakers, University of Amsterdam, Department of Psychology, Weesperplein 4, 1018 XA Amsterdam, the Netherlands. Email address: EJ.Wagenmakers@gmail.com.

Morey, & Iverson, 2009; Wagenmakers, 2007). This is a major limitation, particularly for replication studies in which there is an important distinction between the statement “ $p > .05$ , the data are uninformative” versus the statement “ $p > .05$ , the data are informative and support the null hypothesis”.

It should be noted that the experiments from Donnellan et al. (in press) featured a total of 3073 participants; for such high-power experiments one expects the outcome to be diagnostic, and hence it may be tempting to conclude that the non-significant  $p$  values reported by Donnellan et al. (in press) do indicate support in favor of the null hypothesis. However, this argument from power is insufficient, for two reasons. First, power is a pre-experimental expectation over all possible outcomes, only one of which is relevant after the data are observed. In other words, even when conducting high-power experiments, researchers can be unlucky and obtain uninformative outcomes. Second, even if the data could be argued to provide support in favor of the null hypothesis, the quantitative impact of this support remains unclear: are the observed data twice as likely under the null hypothesis  $\mathcal{H}_0$  than under the alternative hypothesis  $\mathcal{H}_1$ , or 20 times, or perhaps 200 times?

Here we provide a series of Bayesian hypothesis tests to grade the decisiveness of the evidence that the data from Donnellan et al. (in press) provide in favor of the null hypothesis that people do not use warm showers and baths to compensate for a lack of social warmth. Throughout this article we display a suite of Bayesian hypothesis tests: a default two-sided test for correlations (Jeffreys, 1961), a default one-sided test for correlations (Boekel et al., 2014), a sensitivity analysis (Ly, Verhagen, & Wagenmakers, 2014), and a replication test for correlations (extending the work by Verhagen & Wagenmakers, 2014).

Our results show that although most  $p$  values from Donnellan et al. (in press) are non-significant, the evidence in favor of  $\mathcal{H}_0$ —as quantified by the default two-sided Bayesian hypothesis test—differs widely across the nine replication attempts: for the least informative attempt, the observed data are only 2 times more likely under  $\mathcal{H}_0$  than under  $\mathcal{H}_1$ ; for the most informative attempt, the observed data are 17 times more likely under  $\mathcal{H}_0$  than under  $\mathcal{H}_1$ . Overall, the combined data from studies 1-4 and studies 5-9 are 16 and 29 times more likely under  $\mathcal{H}_0$  than under  $\mathcal{H}_1$ , respectively.

The methods outlined here are general and can be used for other research that seeks to test correlations as well. The relevant R code is provided in an online appendix.

### The Donnellan Data

In their Studies 1a and 1b, Bargh and Shalev (2012) found that Loneliness—as measured by the UCLA Loneliness Scale—correlated positively with the “Physical Warmth Index”, a composite variable based on self-reported average frequency, duration, and temperature of showers and baths ( $N = 51$ ,  $r = .57$ ,  $p < .0001$ ;  $N = 41$ ,  $r = .37$ ,  $p < .017$ ). Based in part on these results, Bargh and Shalev (2012, p. 155) hypothesized that people “self-regulate their feelings of social warmth (connectedness to others) with applications of physical warmth (through taking warm baths or showers)”.

In this article we reanalyze the data from the nine replication experiments conducted by Donnellan et al. (in press). As explained by Donnellan et al. (in press), Studies 1-4 were near-exact replications (e.g., using a UCLA Loneliness Scale slightly different from the one used by Bargh & Shalev, 2012) and Studies 5-9 were exact replications. In all nine studies,

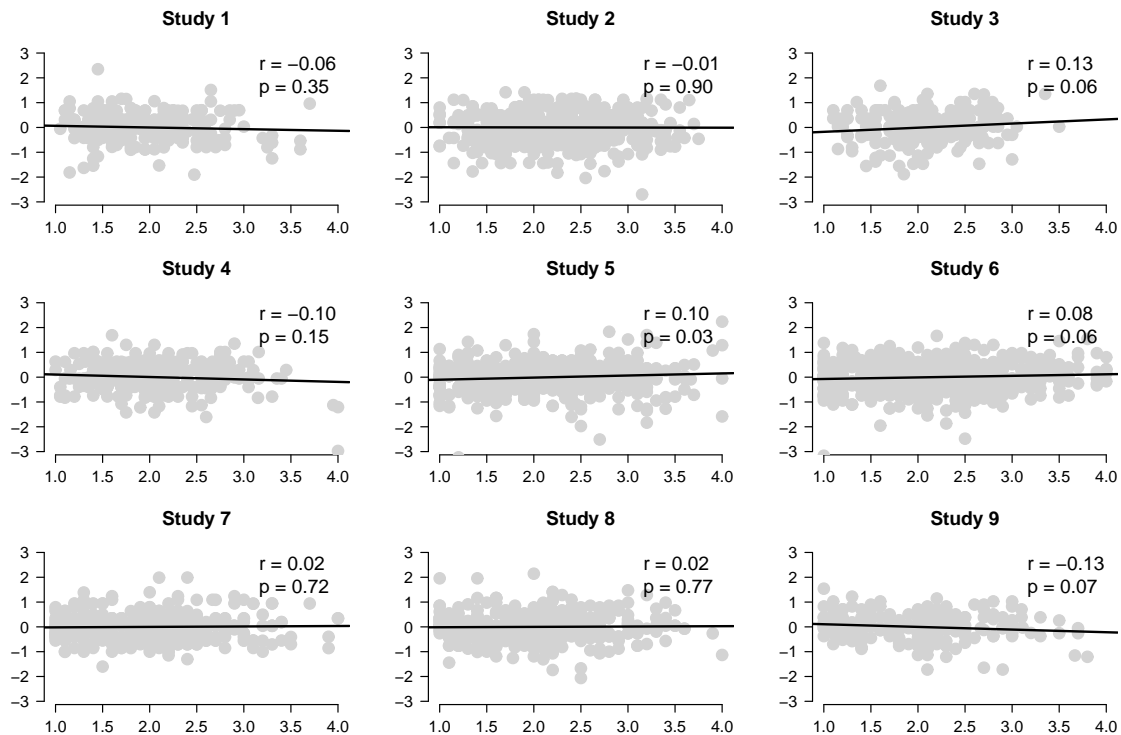


Figure 1. Data for the nine replication experiments from Donnellan et al. (in press). Scores for the Loneliness scale are on the  $x$ -axis and scores for the Physical Warmth Index are on the  $y$ -axis. Each panel also shows the sample Pearson correlation coefficient  $r$  and the two-sided  $p$  value.

the focus of our analysis is the statistical association between Loneliness and the Physical Warmth Index used by Bargh and Shalev (2012).

The first step in analyzing correlations is to plot the data and confirm that the assumption of a linear relation is appropriate (Anscombe, 1973). For instance, a zero correlation between Loneliness and the Physical Warmth Index is misleading when the empirical relation is U-shaped. Figure 1 shows the raw data and confirms the validity of a standard correlational analysis. Across the nine experiments, the sample Pearson correlation values range from  $-0.13$  to  $+0.13$ , and the associated two-sided  $p$  values range from  $.03$  to  $.77$ .

### Posterior Distributions

To quantify the evidence that the data provide for the presence and absence of a correlation  $\rho$  between Loneliness and the Physical Warmth Index we need to contrast two statistical models: the null hypothesis  $\mathcal{H}_0 : \rho = 0$  and the alternative hypothesis  $\mathcal{H}_1 : \rho \neq 0$ . In Bayesian inference, the complete specification of a statistical model requires that its parameters be assigned prior distributions (Dienes, 2008; Lee & Wagenmakers, 2013; Lindley, 2004). For the Pearson correlation, the data are assumed to come from a bivariate normal, and this means that the model has five parameters: parameters  $\mu_1$  and  $\sigma_1^2$  are the

mean and variance of the first variable,  $\mu_2$  and  $\sigma_2^2$  are the mean and variance of the second variable, and  $\rho$  is the correlation (see appendix for details).

We start the specification of  $\mathcal{H}_1$  by assigning uninformative, widely spread-out prior distributions to parameters  $\mu_1$ ,  $\mu_2$ ,  $\sigma_1^2$ , and  $\sigma_2^2$  (Jeffreys, 1961; Lee & Wagenmakers, 2013). This leaves the specification of the prior distribution for the parameter of interest, the correlation  $\rho$ . At first we follow Jeffreys (1961) and assign  $\rho$  a prior that is uniform from  $-1$  to  $1$ ; this prior reflects the belief that each value for  $\rho$  is equally likely before seeing the data. Hence, the alternative hypothesis is specified as  $\mathcal{H}_1 : \rho \sim U(-1, 1)$ .

Assume for the moment that we know with certainty that  $\mathcal{H}_0$  is false and  $\mathcal{H}_1$  is true. In that case our prior knowledge about  $\rho$  is completely captured by its prior distribution  $\rho \sim U(-1, 1)$ . When data  $D$  arrive, this prior distribution  $p(\rho)$  is updated to a posterior distribution  $p(\rho | D)$ . The posterior distribution describes all that we know about  $\rho$  after seeing the data (and ignoring the fact that  $\mathcal{H}_0$  may be true). To provide an initial intuitive impression about what the Donnellan data tell us about the correlation between Loneliness and the Physical Warmth Index, Figure 2 shows prior and posterior distributions separately for each of the nine experiments.<sup>1</sup>

As is evident from each panel in Figure 2, the data are informative in the sense that there is a substantial difference between the prior distribution and the posterior distribution. For Studies 2, 7, and 8, the posterior distribution is approximately centered on  $\rho = 0$ ; for Studies 1, 4, and 9, most of the posterior distribution is concentrated on negative values of  $\rho$ ; and for Studies 3, 5, 6, most of the posterior distribution is concentrated on positive values of  $\rho$ . Although useful, a visual impression of the posterior distribution alone cannot serve to quantify the evidence that the data provide for the hypothesis that the correlation is present or absent, a topic we turn to next.

### Default Bayes Factors

The Bayesian model selection and hypothesis testing machinery works as follows (Jeffreys, 1961). Assume for simplicity that there are only two models under consideration,  $\mathcal{H}_0 : \rho = 0$  and  $\mathcal{H}_1 : \rho \sim U(-1, 1)$ . We start by assigning complementary prior probabilities to both hypotheses, that is  $p(\mathcal{H}_0)$  and  $p(\mathcal{H}_1) = 1 - p(\mathcal{H}_0)$ . Dividing these probabilities yields the prior model odds. For instance, a proponent of the relation between loneliness and bathing habits may believe that  $p(\mathcal{H}_0) = .05$ ; hence, the proponent's prior model odds equal  $p(\mathcal{H}_0)/p(\mathcal{H}_1) = .05/.95 = 1/19$ . Hence, this proponent believes that the presence of a correlation between loneliness and bathing habits is a priori 19 times more plausible than its absence.

Of course, the specification of prior model odds is subjective. In this case, a skeptic may well have prior odds equal to  $p(\mathcal{H}_0)/p(\mathcal{H}_1) = .99/.01 = 99$ , meaning that this skeptic believes that the absence of a correlation between loneliness and bathing habits is a priori 99 times more plausible than its presence. In sum, the prior model odds can be used to measure an individual's initial enthusiasm or skepticism regarding the hypotheses at hand.

Bayesian hypothesis testing, however, does not depend on prior odds; instead, it concerns itself with the change in prior odds brought about by the data. When the data

<sup>1</sup>A complete Bayesian analysis can update the posterior for  $\rho$  across experiments. Here we wish to provide an indication of the informativeness of each experiment separately.

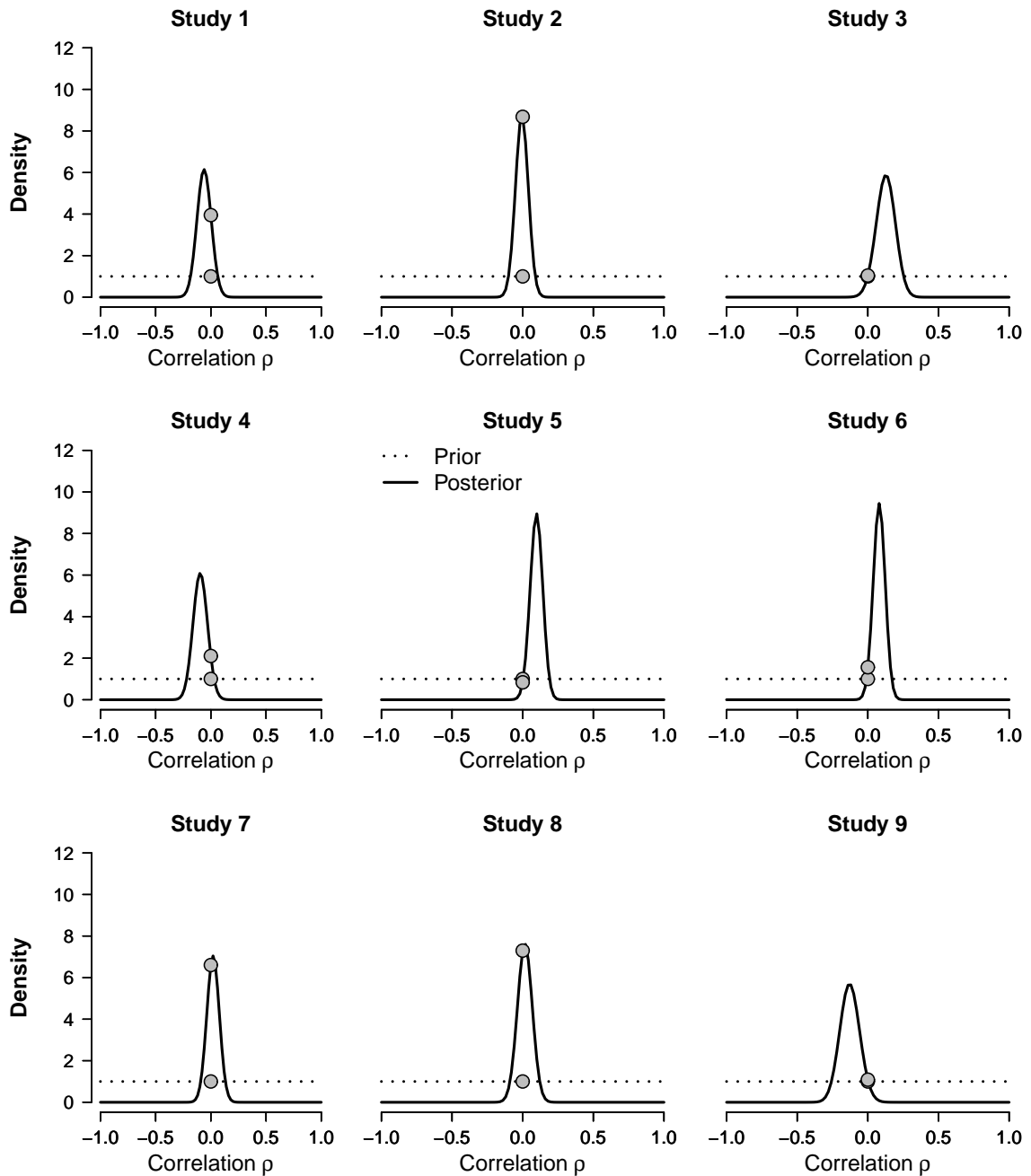


Figure 2. Prior and posterior distributions for the correlation  $\rho$  between Loneliness and the Physical Warmth Index across the nine replication experiments from Donnellan et al. (in press). The statistical model is defined as  $\mathcal{H}_1 : \rho \sim U(-1, 1)$ . The filled dots indicate the height of the prior and posterior distributions at  $\rho = 0$ ; the ratio of these heights equals the evidence that the data provide for  $\mathcal{H}_1$  versus  $\mathcal{H}_0$  (Wagenmakers et al., 2010).

$D$  arrive, the prior model odds are updated to posterior model odds. Mathematically, the updating process proceeds as follows:

$$\underbrace{\frac{p(\mathcal{H}_0 | D)}{p(\mathcal{H}_1 | D)}}_{\text{Posterior odds}} = \underbrace{\frac{p(\mathcal{H}_0)}{p(\mathcal{H}_1)}}_{\text{Prior odds}} \times \underbrace{\frac{p(D | \mathcal{H}_0)}{p(D | \mathcal{H}_1)}}_{\text{Bayes factor}}. \quad (1)$$

The Bayesian hypothesis test centers on the Bayes factor  $\text{BF}_{01}$ : the extent to which the data change one’s belief about the plausibility of the competing models (Jeffreys, 1961; Kass & Raftery, 1995; Lee & Wagenmakers, 2013). Thus, although proponent and skeptic may differ on their prior model odds (and, consequently, on their posterior model odds), as long as they agree on the model specification  $\mathcal{H}_1 : \rho \sim U(-1, 1)$  they will agree precisely on the extent to which the data have changed their initial opinion. For instance, when  $\text{BF}_{01} = 8.5$  the observed data are 8.5 times more likely under  $\mathcal{H}_0$  than under  $\mathcal{H}_1$ , and when  $\text{BF}_{01} = 0.2$  the observed data are five times more likely under  $\mathcal{H}_1$  than under  $\mathcal{H}_0$ .

Thus, in order to grade the decisiveness of the evidence in the nine studies by Donnellan et al. (in press) we need to compute the Bayes factor  $\text{BF}_{01} = p(D | \mathcal{H}_0)/p(D | \mathcal{H}_1)$ . When  $\mathcal{H}_1 : \rho \sim U(-1, 1)$ , this Bayes factor can be obtained easily (Jeffreys, 1961; see also the appendix). The  $\text{BF}_{01}$  column of Table 1 shows the result. As expected from considering the posterior distributions shown in Figure 2, the evidence in favor of  $\mathcal{H}_0$  is particularly high for Studies 2 (i.e.,  $\text{BF}_{01} = 17.34$ ), 7 (i.e.,  $\text{BF}_{01} = 13.19$ ), and 8 (i.e.,  $\text{BF}_{01} = 14.58$ ); each of these studies alone requires that we adjust our beliefs about the presence of a correlation between Loneliness and the Physical Warmth Index by more than an order of magnitude.

To visualize the Bayes factor results, Figure 2 uses filled dots to indicate the height of the prior distribution versus the height of the posterior distribution at  $\rho = 0$ , assuming  $\mathcal{H}_1$  holds. An identity known as the Savage-Dickey density ratio test (e.g., Dickey & Lientz, 1970; Wagenmakers et al., 2010) states that the ratio between these heights equals  $\text{BF}_{01}$ . For instance, consider the Study 2 panel in Figure 2. For that study, the data increased the plausibility of the point  $\rho = 0$  by a factor of 17.34, meaning that at  $\rho = 0$  the posterior distribution is 17.34 times higher than the prior distribution. This height ratio –obtained by considering only the prior and posterior distributions under  $\mathcal{H}_1$ – is identical to the Bayes factor  $\text{BF}_{01}$  between  $\mathcal{H}_0$  and  $\mathcal{H}_1$ .

In addition, the evidence for  $\mathcal{H}_0$  is rather weak for those studies in which the effect is in the predicted direction and most of the posterior mass is concentrated on positive values of  $\rho$ . Specifically, the results from Studies 3 (i.e.,  $\text{BF}_{01} = 2.10$ ), 5 (i.e.,  $\text{BF}_{01} = 1.68$ ), and 6 (i.e.,  $\text{BF}_{01} = 3.13$ ) do not necessitate a substantial adjustment of our beliefs about the presence of a correlation between Loneliness and the Physical Warmth Index, as can be confirmed by the relative closeness of the dots on the distributions in the corresponding panels of Figure 2. Note, however, that even for these relatively uninformative studies the evidence favors  $\mathcal{H}_0$ , whereas the respective classical  $p$  values equal  $p = .06$  (i.e., “marginally significant”),  $p = .03$  (i.e., “significant, reject  $\mathcal{H}_0$ ”), and  $p = .06$  (i.e., “marginally significant”).<sup>2</sup>

Finally, consider the evidence for the studies in which the effect is in the opposite direction and most of the posterior mass is concentrated on negative values of  $\rho$ . The

<sup>2</sup>Results such as these illustrate the strong statement by Edwards (1965, p. 400): “Classical significance tests are violently biased against the null hypothesis.”

results for Studies 1 (i.e.,  $\text{BF}_{01} = 7.90$ ), 4 (i.e.,  $\text{BF}_{01} = 4.22$ ), and 9 (i.e.,  $\text{BF}_{01} = 2.19$ ) yield somewhat more evidence for  $\mathcal{H}_0$  than did Studies 3, 5, and 6, but the overall impression is less compelling than one might expect. The main reason for this is that our current Bayes factor is two-sided such that positive correlations constitute just as much evidence against  $\mathcal{H}_0$  as negative correlations. For this particular scenario, however, there are strong expectations about the direction of the effect, and this warrants the application of a one-sided test.

	$N$	$r$	$p$	$\text{BF}_{01}$	$\text{BF}_{0d}$	$\text{BF}_{0r}(.57)$	$\text{BF}_{0r}(.37)$
Study 1	235	-0.06	0.35	7.90	22.52	<b>16718.92</b>	39.29
Study 2	480	-0.01	0.90	17.34	19.22	17635.98	<b>47.39</b>
Study 3	210	0.13	0.06	2.10	1.09	50.35	<b>1.16</b>
Study 4	228	-0.10	0.15	4.22	28.48	21804.45	<b>35.07</b>
Study 5	494	0.10	0.03	1.68	0.85	135.10	<b>1.32</b>
Study 6	553	0.08	0.06	3.13	1.61	398.61	<b>2.99</b>
Study 7	311	0.02	0.72	13.19	10.31	<b>4872.91</b>	23.71
Study 8	365	0.02	0.77	14.58	11.82	<b>6978.35</b>	28.70
Study 9	107	-0.13	0.07	2.19	30.73	<b>21670.70</b>	28.33
Study 1-4	1152	-0.03	0.31	16.17	52.18	49659.76	70.00
Study 5-9	1919	0.01	0.56	29.53	20.52	31050.36	70.35

Table 1: Results from different Bayes factor hypothesis tests for each of the nine experiments from Donnellan et al. (in press), as well as for the data collapsed over Studies 1-4 and Studies 5-9. Note:  $N$  is the total number of participants,  $r$  is the sample Pearson correlation coefficient between Loneliness and the Physical Warmth Index,  $p$  is the two-sided  $p$  value,  $\text{BF}_{01}$  is the two-sided default Bayes factor in favor of  $\mathcal{H}_0$ ,  $\text{BF}_{0d}$  is the one-sided default Bayes factor in favor of  $\mathcal{H}_0$ ,  $\text{BF}_{0r}(.57)$  is the replication Bayes factor in favor of  $\mathcal{H}_0$  based on Study 1a from Bargh and Shalev (2012) (featuring undergraduate participants, as in Studies 1, 7, 8, and 9), and  $\text{BF}_{0r}(.37)$  is the replication Bayes factor in favor of  $\mathcal{H}_0$  based on Study 1b from Bargh and Shalev (2012) (featuring participants from community samples, as in Studies 2-6).

### *The One-sided Test*

For the two-sided test discussed in the previous section, the alternative hypothesis was specified as  $\mathcal{H}_1 : \rho \sim U(-1, 1)$ . This model specification expresses the belief that every value of the correlation  $\rho$  is equally likely a priori. However, the hypothesis proposed by Bargh and Shalev (2012) and tested by Donnellan et al. (in press) is clearly directional: the assertion is that lonely people take showers and baths that are warmer, not colder.

Within the Bayesian framework, it is conceptually straightforward to account for the direction of the hypothesis. Specifically, for the one-sided test the prior mass is assigned only to positive values of  $\rho$  such that  $\mathcal{H}_d : \rho \sim U(0, 1)$ , where the subscript  $d$  stands for “directional”. The computation of the associated one-sided Bayes factor  $\text{BF}_{0d}$  is provided in the appendix (see also Boekel et al., 2014; Morey & Wagenmakers, 2014). The  $\text{BF}_{0d}$  column of Table 1 shows the result.

A comparison between the two-sided Bayes factor  $\text{BF}_{01}$  and the one-sided Bayes factor  $\text{BF}_{0d}$  reveals three regularities (see Table 1). The first regularity is that for the three studies

where the posterior distribution from Figure 2 was approximately symmetrical around  $\rho = 0$ , the evidence in favor of  $\mathcal{H}_0$  is virtually unaffected; Study 2:  $\text{BF}_{01} = 17.34$  vs.  $\text{BF}_{0d} = 19.22$ ; Study 7:  $\text{BF}_{01} = 13.19$  vs.  $\text{BF}_{0d} = 10.31$ ; Study 8:  $\text{BF}_{01} = 14.58$  vs.  $\text{BF}_{0d} = 11.82$ . In fact, when the posterior distribution is perfectly symmetrical around zero, the two Bayes factors are identical (Wagenmakers et al., 2010).

The second regularity is that for the studies where the effect was in the predicted direction, the evidence is now more favorable to  $\mathcal{H}_d$  than it was to  $\mathcal{H}_1$ ; Study 3:  $\text{BF}_{01} = 2.10$  vs.  $\text{BF}_{0d} = 1.09$ ; Study 5:  $\text{BF}_{01} = 1.68$  vs.  $\text{BF}_{0d} = 0.85$ ; Study 6:  $\text{BF}_{01} = 3.13$  vs.  $\text{BF}_{0d} = 1.61$ . Under the one-sided test, the data from these studies have become almost completely uninformative. The data from Study 5 even favor  $\mathcal{H}_1$ , although the strength of this support is so small that it does not deserve attention (i.e., the data are  $1/0.85 \approx 1.18$  times more likely under  $\mathcal{H}_d$  than under  $\mathcal{H}_0$ ). Thus, when the effect goes in the predicted direction the one-sided test makes the alternative hypothesis look better, but not by much. In fact, for a symmetrical prior a sign-restriction cannot increase the Bayes factor in favor of the alternative hypothesis more than two-fold (Klugkist, Laudy, & Hoijsink, 2005; Wagenmakers et al., 2010).

The third regularity is that for the studies where the effect was in the opposite direction, the evidence is much less favorable for  $\mathcal{H}_d$  than it was for  $\mathcal{H}_1$ ; Study 1:  $\text{BF}_{01} = 7.90$  vs.  $\text{BF}_{0d} = 22.52$ ; Study 4:  $\text{BF}_{01} = 4.22$  vs.  $\text{BF}_{0d} = 28.48$ ; Study 9:  $\text{BF}_{01} = 2.19$  vs.  $\text{BF}_{0d} = 30.73$ . This is then the major difference between specifying a two-sided alternative hypothesis  $\mathcal{H}_1$  and a one-sided alternative hypothesis  $\mathcal{H}_d$ : when the effect goes in the direction opposite to the one that was predicted, the evidence greatly favors  $\mathcal{H}_0$ . This happens because the evidence quantified by the Bayes factor is relative: when the observed effect is negative, this may be unlikely under  $\mathcal{H}_0$ , but it is even less likely under a model  $\mathcal{H}_d$  that stipulates the effect to be positive.

In sum, by changing the prior distribution on  $\rho$  we can implement a one-sided Bayes factor that quantifies the evidence that the data provide for a positive correlation between Loneliness and the Physical Warmth Index. This one-sided test is arguably a better reflection of the underlying directional hypothesis, which states that lonely people take warmer –but not colder– showers and baths. Application of the one-sided test showed that the out of the nine replication experiments by Donnellan et al. (in press), three were not very informative. The other six studies, however, provided highly diagnostic information, each separately requiring a shift in belief towards  $\mathcal{H}_0$  of more than an order of magnitude.

### *Sensitivity Analysis*

The comparison between the two-sided and the one-sided Bayes factor has highlighted how the prior distribution on  $\rho$  can be used to specify different alternative hypotheses; and when different hypotheses are put to the test, different results will (and should) emerge. A persistent concern, however, is that the presented Bayes factor may be delicately sensitive to the specification of the prior, and that by specifying the prior at will, researchers can obtain any desired result. This concern can be addressed in more than one way. The most general counterargument is that the prior is an integral part of the model specification process – yes, one can specify a highly implausible and idiosyncratic prior on  $\rho$  to obtain a nonsensical result, but the specification of the prior is subject to criticism just as the specification of a highly implausible and idiosyncratic model structure (e.g., an exponential



distribution for response times). In other words, silly models (whether through silly priors or silly structure) will lead to silly conclusions, but in many situations it is obvious when a model is silly and when it is not.

A related counterargument is that for many models, researchers can depend on default priors that are suitable for a reference-style analysis. This analysis can be refined if more knowledge is available, as was demonstrated above: we started with a two-sided default prior  $\mathcal{H}_1 : \rho \sim (-1, 1)$  and then refined the prior to  $\mathcal{H}_d : \rho \sim (0, 1)$ . An extreme form of refinement will be demonstrated in the next section. There, the prior distribution for the Bayes factor analysis of the Donnellan et al. (in press) studies is provided by the posterior distribution obtained from the Bargh and Shalev (2012) studies.

In this section we explore another counterargument, namely to take the critique and evaluate it explicitly by means of a sensitivity analysis (e.g., Wagenmakers, Wetzels, Borsboom, & van der Maas, 2011). In such an analysis, one calculates Bayes factors for a wide range of plausible prior distributions. If the conclusions depend on the prior specification in an important way, such that different plausible priors lead to qualitatively different Bayes factors, then it should be acknowledged that the data do not allow an unambiguous conclusion. However, it may also happen that the conclusions are qualitatively robust across a wide range of prior distributions (e.g., Wagenmakers et al., 2011). In our experience such robustness is the rule rather than the exception.

For consistency with the two-sided tests carried out by Bargh and Shalev (2012) and Donnellan et al. (in press), we return to the two-sided Bayes factor  $\text{BF}_{01}$  that compares  $\mathcal{H}_0 : \rho = 0$  to  $\mathcal{H}_1 : \rho \sim U(-1, 1)$ . One proposal for a sensitivity analysis could define a set of models by smoothly decreasing the range of the uniform distribution on  $\rho$ , such that  $\mathcal{H}_1 : \rho \sim U(-c, c)$ , with  $c \in (0, 1)$ . We prefer a similar but more elegant solution, where we first rescale the correlation to lie between 0 and 1, and then assign it a beta distribution. Hence,  $\rho' \sim \text{beta}(a, a)$ , and a measure of the spread of this distribution is  $\sigma = 1/a$ . We then transform the beta distribution back to the  $(-1, 1)$  scale and calculate the Bayes factors as a function of  $\sigma$ . When  $\sigma = 1$ , this corresponds to a uniform prior on the correlation coefficient, as per our default analysis. When  $\sigma = 0$ , which happens when  $a$  grows very large,  $\mathcal{H}_1$  becomes indistinguishable from  $\mathcal{H}_0$  and consequently the Bayes factor is 1. Values of  $\sigma$  in between 0 and 1 define a continuous range of different alternative hypotheses that represent different beliefs about the extent to which large values for the correlation are plausible.

Figure 3 shows the result of the sensitivity analysis for each of the nine experiments from Donnellan et al. (in press). The  $y$  axis shows the log of the Bayes factor  $\text{BF}_{01}$ , such that when  $\sigma = 0$ , all panels yield  $\log \text{BF}_{01} = \log(1)$  and  $\text{BF}_{01} = 1$ , as predicted. In all panels, for all reasonable values of  $\sigma$ , the evidence supports the null hypothesis. In addition, there is no value of  $\sigma$  for which the evidence supports the alternative hypothesis in compelling fashion. Furthermore, for a large range of  $\sigma$  the Bayes factor does not show large fluctuations. Overall, the sensitivity analysis shows that, although different priors instantiate different models and will therefore yield different Bayes factors, it is not the case that any results whatsoever can be obtained. Instead, the qualitative results are similar across a range of plausible values for  $\sigma$ : the data provide clear evidence in favor of  $\mathcal{H}_0$ , but some experiments provide stronger evidence than others.

The same sensitivity analysis can be carried out after collapsing the data in two

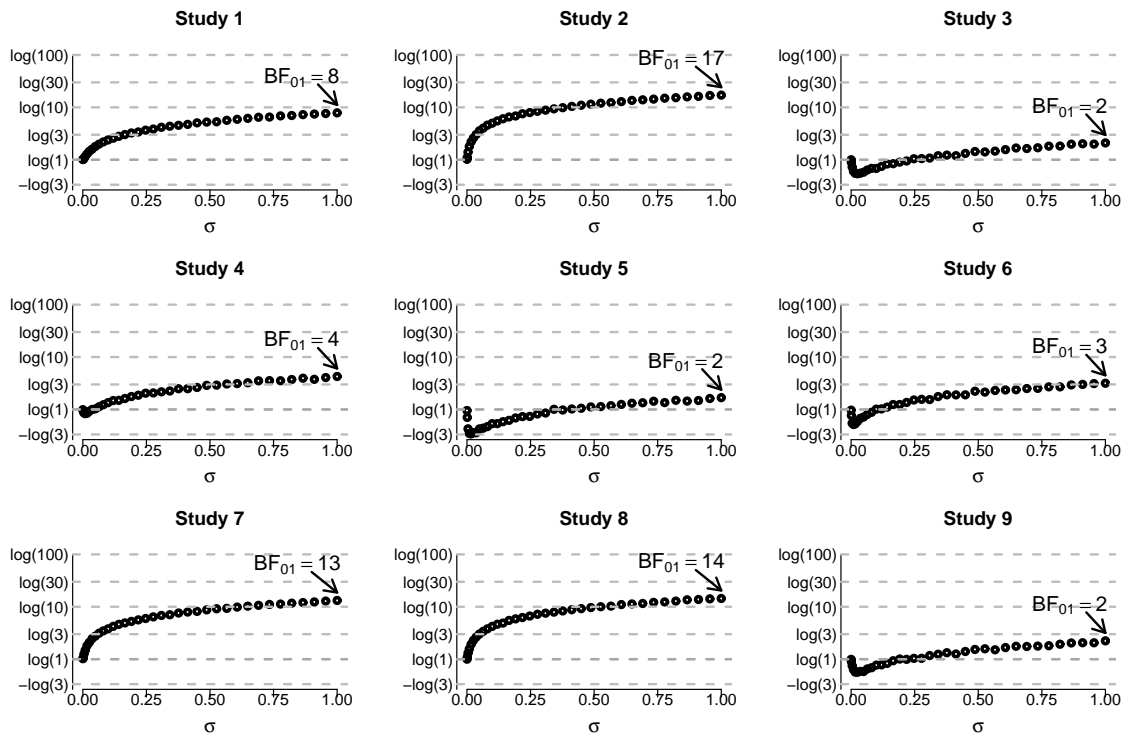


Figure 3. Sensitivity analysis for the Bayes factor  $BF_{01}$  across the nine replication experiments from Donnellan et al. (in press). The log of the Bayes factor  $BF_{01}$  is on the  $x$ -axis and the prior width  $\sigma$  is on the  $y$ -axis. When  $\sigma = 0$  the alternative hypothesis equals the null hypothesis; when  $\sigma = 1$  the alternative hypothesis is  $\rho \sim U(-1, 1)$ . The Bayes factor is qualitatively robust in the sense that the evidence favors the null hypothesis across a wide range of prior beliefs. See text for details.

classes: one based on Studies 1-4 and one based on Studies 5-9. The studies within these two classes were highly similar (Donnellan et al., in press). Figure 4 shows the result. All values for  $\sigma$  result in Bayes factors that indicate support in favor of  $\mathcal{H}_0$ . When  $\mathcal{H}_1$  is defined so as to predicts larger effects (i.e., through larger values of  $\sigma$ ), the evidence more strongly supports  $\mathcal{H}_0$ . Thus, the more the models become distinguishable, the more the Bayes factor prefers  $\mathcal{H}_0$ .

It is insightful to compare the Bayes factors for the collapsed data from Studies 1-4 (i.e.,  $BF_{01} = 16.17$ ) and Studies 5-9 (i.e.,  $BF_{01} = 29.53$ ) to those obtained by multiplying the Bayes factors from the individual experiments. For Studies 1-4, the multiplication yields  $7.90 \times 17.34 \times 2.10 \times 4.22 \approx 1214$ ; for Studies 5-9, the multiplication yields  $1.68 \times 3.13 \times 13.19 \times 14.58 \times 2.19 \approx 2215$ . The discrepancy with the collapsed-data Bayes factors is large, and this serves to demonstrate that when effect sizes in similar experiments are similar and not independent –as is reasonable to assume– Bayes factors should not be multiplied (e.g., as was done by Bem, Utts, & Johnson, 2011 in order to present evidence in favor of extra-sensory perception).

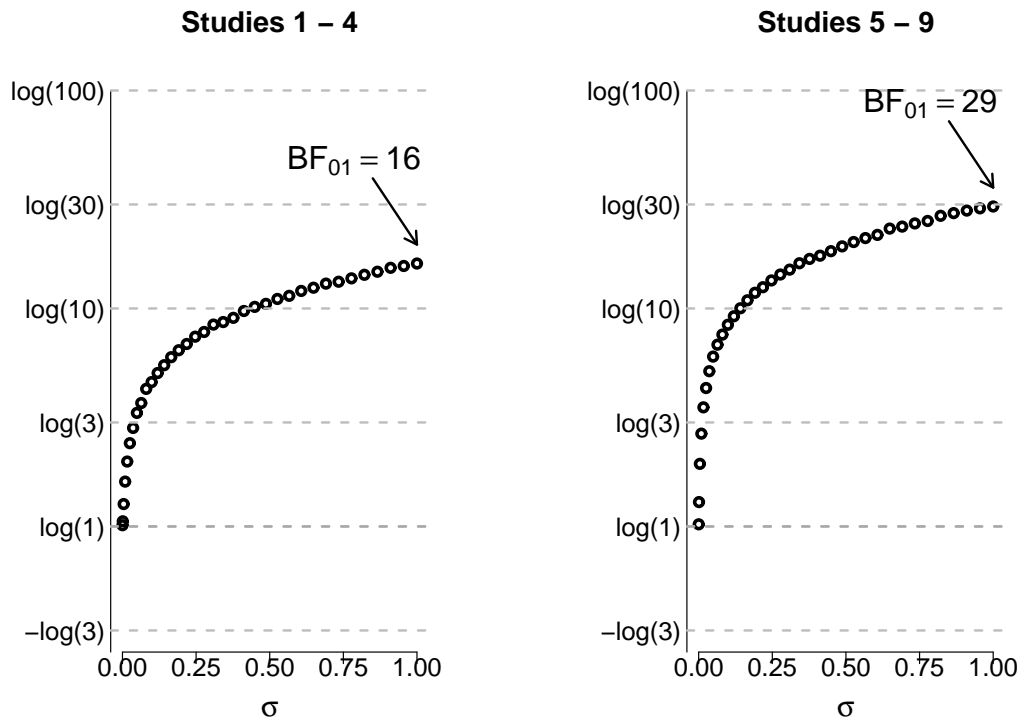


Figure 4. Sensitivity analysis for the Bayes factor  $BF_{01}$ , collapsing data across Studies 1-4 (left panel) and across Studies 5-9 (right panel) from Donnellan et al. (in press). The log of the Bayes factor  $BF_{01}$  is on the  $x$ -axis and the prior width  $\sigma$  is on the  $y$ -axis. When  $\sigma = 0$  the alternative hypothesis equals the null hypothesis; when  $\sigma = 1$  the alternative hypothesis is  $\rho \sim U(-1, 1)$ . See text for details.

### Replication Bayes Factors

As outlined above, for replication studies there exists another way to alleviate the concern over how to specify the alternative hypothesis (Verhagen & Wagenmakers, 2014). Specifically, one can use the data from the original study to obtain a posterior distribution, and then use that posterior distribution to specify the alternative hypothesis for the analysis of the replication studies. This “replication Bayes factor” therefore pits two models against one another. The first model,  $\mathcal{H}_0 : \rho = 0$ , represents the belief of a skeptic, and the second model,  $\mathcal{H}_r : \rho \sim$  “posterior distribution from original study”, represents the idealized belief of a proponent. As pointed out by Verhagen and Wagenmakers (2014, p. 1459), “(...) the default test addresses the question, “Given that we know relatively little about the expected effect size beforehand, is the effect present or absent in the replication attempt?”; our test addresses the question, “Is the effect similar to what was found before, or is it absent?”. The two tests therefore represent extremes on a continuum of sensitivity to past research; the default test completely ignores the outcomes of an earlier experiment, whereas the replication test takes these outcomes fully into account.”

The replication Bayes factor was developed by Verhagen and Wagenmakers (2014) for

the  $t$  test; here we extend that work to the Pearson correlation coefficient (for an application see Boekel et al., 2014; for mathematical details see the appendix). Table 1 shows the results for two replication Bayes factors; the first,  $\text{BF}_{0r}(.57)$ , is based on Study 1a from Bargh and Shalev (2012), featuring undergraduate participants and yielding  $r = .57$ ; the second,  $\text{BF}_{0r}(.37)$ , is based on Study 1b from Bargh and Shalev (2012), featuring a community sample of participants and yielding  $r = .37$ .

The  $\text{BF}_{0r}(.57)$  column of Table 1 shows that, across all studies, the data are much more likely under the skeptic's  $\mathcal{H}_0$  than under the proponent's  $\mathcal{H}_r$  based on Study 1a from Bargh and Shalev (2012). Even for the least compelling study, the data are 50.35 times more likely under  $\mathcal{H}_0$  than under  $\mathcal{H}_r$ . When the proponent's belief is based on Study 1b from Bargh and Shalev (2012), the results are less extreme: the results for Study 3 ( $\text{BF}_{0r}(.37) = 1.16$ ), Study 5 ( $\text{BF}_{0r}(.37) = 1.32$ ), and Study 6 ( $\text{BF}_{0r}(.37) = 2.99$ ) are relatively uninformative: the data are almost as likely under the skeptic's  $\mathcal{H}_0$  than under the proponent's  $\mathcal{H}_r$ . For the remaining studies, however, the results show compelling support for the skeptic's  $\mathcal{H}_0$ , with Bayes factors ranging from about 23 to about 47.

Figure 5 visualizes the results using the Savage-Dickey density ratio. In each panel, the dotted line indicates the idealized belief of a proponent, that is, the posterior distribution from the original study by Bargh and Shalev (2012).<sup>3</sup> Studies 1, 7, 8, and 9 featured undergraduate participants, and hence the dotted lines in the corresponding panels are based on Study 1a by Bargh and Shalev (2012) (i.e.,  $r = .57$ ); in contrast, Studies 2-6 featured community samples of participants, and hence the dotted lines are based on Study 1b by Bargh and Shalev (2012) (i.e.,  $r = .37$ ). In each panel, the solid line indicates the posterior distribution that was obtained after updating the beliefs based on the original study (i.e., the dotted distribution) with the data from the replication study.

As before, the Bayes factor  $\text{BF}_{0r}$  is given by the ratio of the height of the prior and posterior distribution at  $\rho = 0$ . For instance, the panel for Study 1 shows that the value  $\rho = 0$  is much more plausible after having seen the data from the replication study than before. In fact, the ratio of the prior and posterior height at  $\rho = 0$  equals 16,718.92, exactly equal to  $\text{BF}_{0r}(.57)$ .

Similarly, the panel for Study 3 shows that the data from the replication study have hardly altered the plausibility of the value  $\rho = 0$  at all; hence the dot that indicates the height of the prior at  $\rho = 0$  overlaps with the dot that indicates the height of the posterior at  $\rho = 0$ , and the Bayes factor equals  $\text{BF}_{0r}(.37) = 1.16$ .

### Concluding Comments

In this article we illustrated a suite of Bayesian hypothesis testing techniques that allow researchers to grade the decisiveness of the evidence that the data provide for the presence versus the absence of a correlation between two dependent variables. This approach is fundamentally different from Fisher's  $p$  value methodology, which does not take acknowledge the existence of an alternative hypothesis, and it is also fundamentally different from Neyman and Pearson's frequentist tool for making decisions. As stated eloquently by

---

<sup>3</sup>In order to obtain the posterior distribution from the original experiment we still require a prior. However, even for relatively small data sets the shape of the posterior distribution is not much affected by the choice of prior distribution, as expressed by the adage "the data overwhelm the prior".

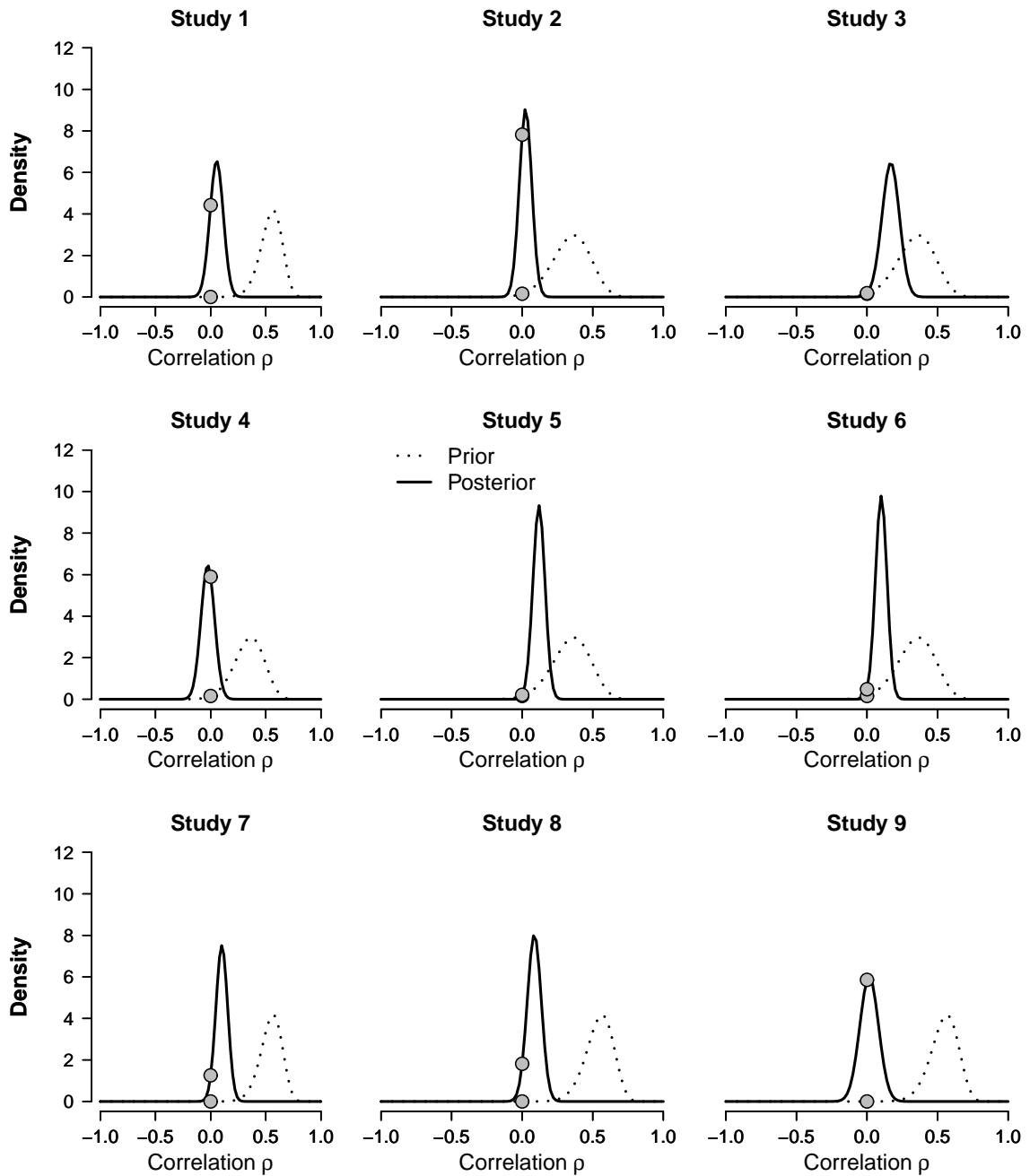


Figure 5. Prior and posterior distributions for the correlation  $\rho$  between Loneliness and the Physical Warmth Index across the nine replication experiments from Donnellan et al. (in press). The statistical model is defined as  $\mathcal{H}_r : \rho \sim$  “posterior distribution from original study”. The filled dots indicate the height of the prior and posterior distributions at  $\rho = 0$ ; the ratio of these heights equals the evidence that the data provide for the proponent’s  $\mathcal{H}_r$  versus the skeptic’s  $\mathcal{H}_0$  (Wagenmakers et al., 2010).

Rozeboom (1960, pp. 422-423): “The null-hypothesis significance test treats ‘acceptance’ or ‘rejection’ of a hypothesis as though these were decisions one makes. But a hypothesis is not something, like a piece of pie offered for dessert, which can be accepted or rejected by a voluntary physical action. Acceptance or rejection of a hypothesis is a cognitive process, a degree of believing or disbelieving which, if rational, is not a matter of choice but determined solely by how likely it is, given the evidence, that the hypothesis is true.”

What, then, are the practical advantages of the Bayes factor hypothesis test over its classical counterpart? Among the most salient are the following: (1) Bayes factors allow researchers to claim evidence in favor of the null hypothesis (Gallistel, 2009; Rouder et al., 2009; Wagenmakers, 2007), an advantage that is particularly prominent in replication research such as that conducted by Donnellan et al. (in press); (2) Bayes factors allow researchers to quantify the above claim, so that we may know whether the data are more likely under  $\mathcal{H}_0$  by a factor of 2, by a factor of 20, or by a factor of 200; (3) Bayes factors allow researchers to monitor the “evidential flow”<sup>4</sup> as the data come in and stop data collection whenever this is deemed desirable, without the need for corrections depending on the intent with which the data were collected (Rouder, 2014; Wagenmakers, Wetzels, Borsboom, van der Maas, & Kievit, 2012). This flexibility is a direct consequence of the Stopping Rule Principle (Berger & Wolpert, 1988), a principle that all Bayesian analyses respect.

One may be tempted to argue that sensible conclusions can be reached using classical statistics when, in addition to the  $p$  value, the concept of power is taken into account. However, as alluded to in the introduction, power is a pre-experimental concept that entails averaging across all possible data sets, only one of which ends up being observed. It is therefore entirely possible that an uninformative result is obtained even after conducting a high-power experiment. For instance, consider Studies 3, 5, and 6 from Donnellan et al. (in press); all our Bayes factor hypothesis tests indicated that these studies were virtually uninformative. Nevertheless, these studies featured 210, 494, and 553 participants, respectively. It is hard to argue that the uninformative nature of these data is due to a lack of power (see also Wagenmakers et al., in press).

Some psychologists and statisticians object to hypothesis testing on the grounds that a point null hypothesis (e.g.,  $\rho = 0$  exactly) is known to be false from the outset (e.g., Cohen, 1994; Meehl, 1978). We disagree with this claim on principle (e.g., Iverson, Wagenmakers, & Lee, 2010), but, more importantly, even if the claim were true it would not detract from the usefulness of hypothesis testing – instead, it could mean only that  $\mathcal{H}_0$  needs to be specified with more care. For instance, for a test of the Pearson correlation coefficient one may replace  $\mathcal{H}_0 : \rho = 0$  with  $\mathcal{H}'_0 : \rho \sim U(-.01, .01)$ . After specifying such an interval null hypothesis (Morey & Rouder, 2011), the same methods outlined in this article may then be applied, with virtually identical results. That is, “(...) the assignment of a lump of prior probability to the simple hypothesis is strictly a mathematical convenience and not at all fundamental.” (Cornfield, 1969, p. 637).

Finally, it should be acknowledged that, in many cases, the data pass the interocular traumatic test (i.e., when the data hit you right between the eyes; Edwards, Lindman, & Savage, 1963) and it does not matter whether one carries out a classical analysis, a

<sup>4</sup>To the best of our knowledge, this term was introduced in the blog of Eliezer Yudkowsky.

Bayesian analysis, or no analysis at all. This argument loses some of its force, however, when the data appear to support the null hypothesis and an intuitive assessment of evidential strength becomes non-trivial. At any rate, one purpose of statistics is to make our intuitive judgement precise and quantitative. We hope that the methods outlined in this article will help contribute to that purpose.

## References

- Anscombe, F. J. (1973). Graphs in statistical analysis. *The American Statistician*, *27*, 17–21.
- Bargh, J. A., & Shalev, I. (2012). The substitutability of physical and social warmth in daily life. *Emotion*, *12*, 154–162.
- Bem, D. J., Utts, J., & Johnson, W. O. (2011). Must psychologists change the way they analyze their data? *Journal of Personality and Social Psychology*, *101*, 716–719.
- Berger, J. O., & Wolpert, R. L. (1988). *The likelihood principle (2nd ed.)*. Hayward (CA): Institute of Mathematical Statistics.
- Boekel, W., Wagenmakers, E.-J., Belay, L., Verhagen, A. J., Brown, S. D., & Forstmann, B. (2014). A purely confirmatory replication study of structural brain-behavior correlations. *Manuscript submitted for publication*.
- Cohen, J. (1994). The earth is round ( $p < .05$ ). *American Psychologist*, *49*, 997–1003.
- Cornfield, J. (1969). The Bayesian outlook and its application. *Biometrics*, *25*, 617–657.
- Dickey, J. M., & Lientz, B. P. (1970). The weighted likelihood ratio, sharp hypotheses about chances, the order of a Markov chain. *The Annals of Mathematical Statistics*, *41*, 214–226.
- Dienes, Z. (2008). *Understanding psychology as a science: An introduction to scientific and statistical inference*. New York: Palgrave MacMillan.
- Donnellan, M. B., Lucas, R. E., & Cesario, J. (in press). On the association between loneliness and bathing habits: Nine replications of Bargh and Shalev (2012) Study 1. *Emotion*.
- Edwards, W. (1965). Tactical note on the relation between scientific and statistical hypotheses. *Psychological Bulletin*, *63*, 400–402.
- Edwards, W., Lindman, H., & Savage, L. J. (1963). Bayesian statistical inference for psychological research. *Psychological Review*, *70*, 193–242.
- Gallistel, C. R. (2009). The importance of proving the null. *Psychological Review*, *116*, 439–453.
- Iverson, G. J., Wagenmakers, E.-J., & Lee, M. D. (2010). A model averaging approach to replication: The case of  $p_{rep}$ . *Psychological Methods*, *15*, 172–181.
- Jeffreys, H. (1961). *Theory of probability* (3 ed.). Oxford, UK: Oxford University Press.
- Kass, R. E., & Raftery, A. E. (1995). Bayes factors. *Journal of the American Statistical Association*, *90*, 773–795.
- Klugkist, I., Laudy, O., & Hoijsink, H. (2005). Inequality constrained analysis of variance: A Bayesian approach. *Psychological Methods*, *10*, 477–493.
- Lee, M. D., & Wagenmakers, E.-J. (2013). *Bayesian modeling for cognitive science: A practical course*. Cambridge University Press.
- Lindley, D. V. (2004). That wretched prior. *Significance*, *1*, 85–87.
- Ly, A., Verhagen, A. J., & Wagenmakers, E.-J. (2014). Analysis and application of two Bayes factors proposed by Harold Jeffreys (1961). *Manuscript in preparation for the special issue on Bayes factors for the Journal of Mathematical Psychology*.
- Meehl, P. E. (1978). Theoretical risks and tabular asterisks: Sir Karl, Sir Ronald, and the slow progress of soft psychology. *Journal of Consulting and Clinical Psychology*, *46*, 806–834.
- Morey, R. D., & Rouder, J. N. (2011). Bayes factor approaches for testing interval null hypotheses. *Psychological Methods*, *16*, 406–419.



- Morey, R. D., & Wagenmakers, E.-J. (2014). Simple relation between Bayesian order-restricted and point-null hypothesis tests. *Statistics and Probability Letters*, *92*, 121–124.
- Oberhettinger, F. (1972). Hypergeometric functions. In M. Abramowitz & I. A. Stegun (Eds.), *Handbook of mathematical functions with formulas, graphs, and mathematical tables* (pp. 555–566). New York: Dover.
- Rouder, J. N. (2014). Optional stopping: No problem for Bayesians. *Psychonomic Bulletin & Review*, *21*, 301–308.
- 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.
- Rozeboom, W. W. (1960). The fallacy of the null-hypothesis significance test. *Psychological Bulletin*, *57*, 416–428.
- Verhagen, A. J., & Wagenmakers, E.-J. (2014). Bayesian tests to quantify the result of a replication attempt. *Journal of Experimental Psychology: General*, *143*, 1457–1475.
- Wagenmakers, E.-J. (2007). A practical solution to the pervasive problems of  $p$  values. *Psychonomic Bulletin & Review*, *14*, 779–804.
- 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., Verhagen, A. J., Ly, A., Bakker, M., Lee, M. D., Matzke, D., Rouder, J. N., & Morey, R. D. (in press). A power fallacy. *Behavior Research Methods*.
- 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. *Journal of Personality and Social Psychology*, *100*, 426–432.
- Wagenmakers, E.-J., Wetzels, R., Borsboom, D., van der Maas, H. L. J., & Kievit, R. A. (2012). An agenda for purely confirmatory research. *Perspectives on Psychological Science*, *7*, 627–633.

Appendix  
Statistical Details

*The Likelihood*

In Bargh and Shalev (2012), the Pearson's correlation coefficient is used to measure the linear association between Loneliness and the Physical Warmth Index, which we denote by  $X$  and  $Y$  respectively. The Pearson's population correlation coefficient is defined as follows:

$$\rho = \frac{\text{Cov}(X, Y)}{\sigma_x \sigma_y} \text{ or equivalently } \rho = E \left[ \left( \frac{X - \mu_x}{\sigma_x} \right) \left( \frac{Y - \mu_y}{\sigma_y} \right) \right], \quad (2)$$

where  $X$  is taken to be normally distributed with population mean  $\mu_x$  and population standard deviation  $\sigma_x$ . Similarly,  $Y$  is also assumed to be normally distributed such that  $Y \sim \mathcal{N}(\mu_y, \sigma_y^2)$ . Assume that  $\mathcal{H}_1$  holds and there exists a correlation between  $X$  and  $Y$ . In order to describe the pair  $X, Y$  simultaneously we then require five parameters, the normality parameters  $\mu_x, \mu_y, \sigma_x, \sigma_y$  and the linear association  $\rho$ , resulting in the following likelihood:

$$L(\mathcal{H}_1 | D) = \left( \frac{1}{2\pi\sigma_x\sigma_y\sqrt{1-\rho^2}} \right)^n \times \exp \left( - \frac{1}{2(1-\rho^2)} \sum_{i=1}^n \left[ \frac{(x_i - \mu_x)^2}{\sigma_x^2} + \frac{(y_i - \mu_y)^2}{\sigma_y^2} + \frac{2\rho(x_i - \mu_x)(y_i - \mu_y)}{\sigma_x\sigma_y} \right] \right), \quad (3)$$

where we have written  $\mathcal{H}_1 = (\mu_x, \mu_y, \sigma_x, \sigma_y, \rho)$  for the five parameters and  $D$  for the observed data with  $D = \begin{pmatrix} x_1 \\ y_1 \end{pmatrix}, \dots, \begin{pmatrix} x_n \\ y_n \end{pmatrix}$  and  $\begin{pmatrix} x_i \\ y_i \end{pmatrix}$  being the reported Loneliness and Physical Warmth Index of participant  $i$ .

When the null hypothesis of no linear association between  $X$  and  $Y$  holds true, this means that  $\rho$  is fixed at zero. Consequently, this yields a model with only four free parameters,  $\mathcal{H}_0 = (\mu_x, \mu_y, \sigma_x, \sigma_y)$ . More precisely, the likelihood for the null model given the observations  $D$  then depends on the four parameters as follows:

$$L(\mathcal{H}_0 | D) = \left( \frac{1}{2\pi\sigma_x\sigma_y} \right)^n \exp \left( - \frac{1}{2} \sum_{i=1}^n \left[ \frac{(x_i - \mu_x)^2}{\sigma_x^2} + \frac{(y_i - \mu_y)^2}{\sigma_y^2} \right] \right). \quad (4)$$

Note that Eq. (4) is simply Eq. (3) with  $\rho = 0$ .

*The Bayes Factor Test Proposed by Sir Harold Jeffreys*

To test whether the population correlation  $\rho$  is zero, we compare Eq. (4) to Eq. (3). Although the true values of the parameters are unknown, we can quantify our degree of belief about the true values by means of prior distributions. These prior distributions act as weighting functions for the likelihood. Sir Harold Jeffreys proposed that for  $\mathcal{H}_0$ , we weight the likelihood Eq. (4) with respect to the population means  $\mu_x$  and  $\mu_y$  proportional to 1, mathematically,  $p(\mu_x) \propto 1$  and  $p(\mu_y) \propto 1$ . Furthermore, the standard deviations are weighted proportional to their inverse, that is,  $p(\sigma_x) \propto \frac{1}{\sigma_x}$  and  $p(\sigma_y) \propto \frac{1}{\sigma_y}$ . These weighting functions then fully specify the marginal or average likelihood for  $\mathcal{H}_0$ . To obtain the

marginal likelihood for  $\mathcal{H}_1$  we use the same weighting functions for the common parameters and we weight the effects of  $\rho$  uniformly over  $(-1, 1)$ , that is,  $p(\rho) = \frac{1}{2}$ . Hence, the two marginal likelihoods are given by the following integrals (i.e., averages):

$$P(D | \mathcal{H}_0) \stackrel{\text{Eq. (4)}}{=} \int \int \int \int L(\mu_x, \mu_y, \sigma_x, \sigma_y | D) d\mu_x d\mu_y \frac{1}{\sigma_x} d\sigma_x \frac{1}{\sigma_y} d\sigma_y, \quad (5)$$

$$P(D | \mathcal{H}_1) \stackrel{\text{Eq. (3)}}{=} \int \int \int \int \int L(\mu_x, \mu_y, \sigma_x, \sigma_y, \rho | D) d\mu_x d\mu_y \frac{1}{\sigma_x} d\sigma_x \frac{1}{\sigma_y} d\sigma_y \frac{1}{2} d\rho. \quad (6)$$

The ratio of these two marginal likelihoods yields the Bayes factor  $\text{BF}_{01} = \frac{P(D | \mathcal{H}_0)}{P(D | \mathcal{H}_1)}$  that allows us to compare the two models as discussed in the main text.

The above equations suggest that we have to compute nine intensive integrals in order to obtain the Bayes factor. Fortunately, this is unnecessary; Jeffreys (1961) showed that only a single integral is required, one with respect to parameter of interest  $\rho$ :

$$\text{BF}_{01} = \frac{1}{\text{BF}_{10}}, \text{ where } \text{BF}_{10} = \frac{1}{2} \int_{-1}^1 \frac{(1 - \rho^2)^{\frac{n-1}{2}}}{(1 - \rho r)^{n - \frac{3}{2}}} d\rho, \quad (7)$$

where  $r$  refers to the sample correlation  $r$  that is defined as:

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}}. \quad (8)$$

#### *The One-Sided Extension*

It is straightforward to extend Jeffreys' result to a one-sided test of the null hypothesis  $\mathcal{H}_0$  that  $\rho = 0$  versus the directional restriction  $\rho > 0$ , which we denote by  $\mathcal{H}_d$ . The extension only requires us to change the uniform prior of  $\rho$  from on  $(-1, 1)$  to a uniform prior on  $(0, 1)$ , which yields:

$$\text{BF}_{0d} = \frac{1}{\text{BF}_{d0}}, \text{ where } \text{BF}_{d0} = \int_0^1 \frac{(1 - \rho^2)^{\frac{n-1}{2}}}{(1 - \rho r)^{n - \frac{3}{2}}} d\rho. \quad (9)$$

#### *The Replication Bayes Factor*

A replication Bayes factor (Verhagen & Wagenmakers, 2014) answers the question: "Is the effect from the replication attempt comparable to what was found before, or is it absent?". When a correlational study is replicated, the replication Bayes factor compares evidence in favor of the null hypothesis of no effect,  $\mathcal{H}_0 : \rho = 0$ , with the evidence in favor of the alternative hypothesis that the effect is equal to the effect found in the original study,  $\mathcal{H}_r : \rho \sim$  "posterior distribution from original study".

The replication Bayes factor is calculated in two steps. In the first step, the posterior distribution of the original study is obtained, assuming a uniform prior distribution on the correlation. The density of this posterior distribution was given by Jeffreys (1961, p. 175, equation 9), and simplifies to:

$$p(\rho | D_{orig}) = \frac{\frac{(1-\rho^2)^{\frac{1}{2}n}}{(1-\rho r)^{n-\frac{1}{2}}} \sqrt{\frac{\pi}{2}} \frac{\Gamma(n)}{\Gamma(n+\frac{1}{2})} {}_2F_1(\frac{1}{2}, \frac{1}{2}, n + \frac{1}{2}, \frac{1}{2} + \frac{1}{2}r\rho)}{\int \frac{(1-\rho^2)^{\frac{1}{2}n}}{(1-\rho r)^{n-\frac{1}{2}}} \sqrt{\frac{\pi}{2}} \frac{\Gamma(n)}{\Gamma(n+\frac{1}{2})} {}_2F_1(\frac{1}{2}, \frac{1}{2}, n + \frac{1}{2}, \frac{1}{2} + \frac{1}{2}r\rho) d\rho}, \quad (10)$$

where  ${}_2F_1$  is Gauss' hypergeometric function (Oberhettinger, 1972, section 15).

The second step consists of the computation of the Bayes factor by integration over this posterior distribution in order to obtain  $p(D | \mathcal{H}_1)$ :

$$\begin{aligned}
 \text{BF}_{10} &= \frac{p(D | \mathcal{H}_1)}{p(D | \mathcal{H}_0)} \\
 &= \frac{\int p(D | \delta, \mathcal{H}_1) p(\delta | \mathcal{H}_1) d\delta}{p(D | \mathcal{H}_0)} \\
 &= \frac{\int \frac{(1-\rho^2)^{\frac{n-1}{2}}}{(1-\rho r)^{n-\frac{3}{2}}} p(\rho | D_{orig}) d\delta}{p(\rho = 0 | D_{orig})} \\
 &= \int \frac{(1-\rho^2)^{\frac{n-1}{2}}}{(1-\rho r)^{n-\frac{3}{2}}} p(\rho | D_{orig}) d\delta
 \end{aligned}$$

which can be accomplished by one-dimensional integration. The R code to perform this analysis can be found at [http://www.josineverhagen.com/?page\\_id=76](http://www.josineverhagen.com/?page_id=76).