

# 1 Non-dichotomous inference using bootstrapped evidence

2

3 *D. Samuel Schwarzkopf*

4 Experimental Psychology, University College London, 26 Bedford Way, London WC1H 0AP, United Kingdom &

5 Institute of Cognitive Neuroscience, University College London, 17 Queen Square, London, WC1N 3AR, United Kingdom

6 Email: [s.schwarzkopf@ucl.ac.uk](mailto:s.schwarzkopf@ucl.ac.uk)

7

## 8 **Abstract**

9

10 The problems with classical frequentist statistics have recently received much attention, yet the  
11 enthusiasm of researchers to adopt alternatives like Bayesian inference remains modest. Here I  
12 present the *bootstrapped evidence test*, an objective resampling procedure that takes the  
13 precision with which both the experimental and null hypothesis can be estimated into account.  
14 Simulations and reanalysis of actual experimental data demonstrate that this test minimizes false  
15 positives while maintaining sensitivity. It is equally applicable to a wide range of situations and  
16 thus minimizes problems arising from analytical flexibility. Critically, it does not dichotomize the  
17 results based on an arbitrary significance level but instead quantifies how well the data support  
18 either the alternative or the null hypothesis. It is thus particularly useful in situations with  
19 considerable uncertainty about the expected effect size. Because it is non-parametric, it is also  
20 robust to severe violations of assumptions made by classical statistics.

21

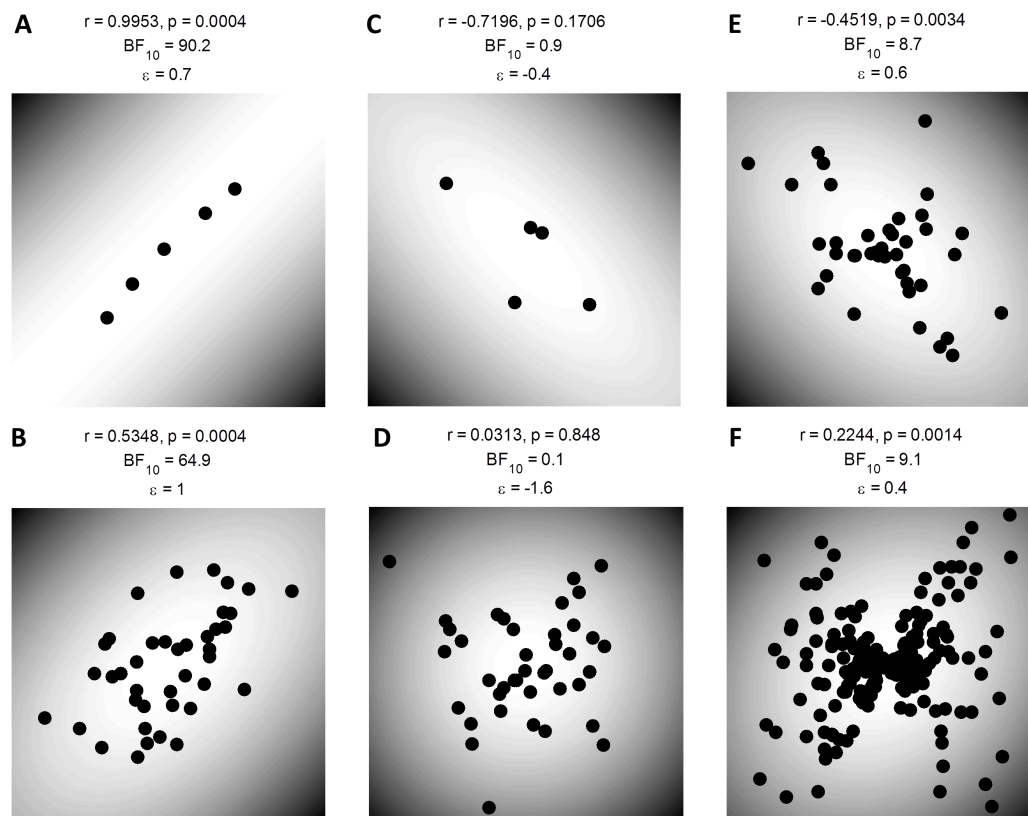
## 22 **Introduction**

23

24 To this day classical null hypothesis significance testing remains the dominant approach for  
25 inferring the validity of an observed result in the psychological and life sciences. It rests on the  
26 probability ('p-value') that the observed effect, or a more extreme one, could have occurred  
27 under the assumption that there is no population effect. If the p-value is sufficiently low, this null  
28 hypothesis is rejected. However, p-values are frequently misinterpreted by researchers,  
29 uninformative about the evidence *for* an experimental hypothesis, highly susceptible to biased  
30 data sampling strategies, and generally prone to false positives [1–11]. The most devastating  
31 effect of p-values may be that they encourage an artificial dichotomy between significant and  
32 non-significant results[5].

33

34 The scatter plots in Fig. 1 illustrate the problems with p-values. Fig. 1A shows an almost perfect  
35 correlation between two measures ( $r=0.995$ ,  $p<0.0004$ ). However, there are only five  
36 observations. In contrast, the data in Fig. 1B are clearly correlated even though the correlation is  
37 weaker. Yet the p-value is similar ( $r=0.535$ ,  $p<0.0004$ ) because the sample size is much larger.  
38 Surely the evidence for a correlation in the second example is more compelling and more likely to  
39 replicate?



40  
 41 **Fig. 1.** Scatter plots showing examples of correlation analysis. A-B. Correlated Gaussian data with  $n=5$   
 42 (A) and  $n=40$  (B). C-D. Uncorrelated Gaussian data with  $n=5$  (C) and  $n=40$  (D). E-F. Severely  
 43 heteroscedastic data with  $n=40$  (E) and  $n=200$  (F). Each black dot is one observation. The grey  
 44 shading denotes the Mahalanobis distance from the bivariate mean. Above each panel the Pearson's  
 45 correlation coefficient, the default Bayes factor [26]  $BF_{10}$  comparing the alternative and the null  
 46 hypothesis, and the bootstrapped evidence  $\varepsilon$  are given.

47  
 48 One journal went so far as to ban the use of classical inference completely while proposing no  
 49 viable alternative [8]. Others proposed guidelines to focus on effect size estimation and  
 50 confidence intervals instead [5,12]. However, the use of confidence intervals is also fraught with  
 51 problems [13] and may simply become a new significance testing procedure in disguise [5,14].  
 52 Moreover, like p-values, confidence intervals are frequently misinterpreted [14] and may perform  
 53 inadequately [15]. Evidence for a hypothesis should *compare* an experimental (alternative)  
 54 hypothesis to a baseline (null) hypothesis. Bayesian hypothesis tests using Bayes factors can  
 55 achieve that but are often difficult to apply and rely on the choice of a prior, which can result in  
 56 considerable debate (see e.g. [3,16–19]).

57  
 58 Here I present the *bootstrapped evidence* (BSE) test. It makes minimal assumptions and is  
 59 applicable to a wide range of situations. Crucially, it quantifies the evidence for either the  
 60 alternative or the null hypothesis non-dichotomously. Yet unlike Bayesian methods it is based  
 61 only on the existing data without any question about prior distributions. It works by  
 62 bootstrapping the effect size distributions under the two hypotheses by using different  
 63 resampling strategies for each. For instance, when quantifying evidence for a linear correlation  
 64 between two variables as in Fig. 1, under the null hypothesis data are resampled without respect

65 to how individual observations have been paired. In contrast, under the alternative hypothesis  
66 the pairing is held intact but pairs are resampled to estimate the strength of the correlation. The  
67 evidence measure quantifies how distinct the distributions for the two hypotheses are and also  
68 incorporates the precision of the estimates. Thus, rather than determining probabilities as most  
69 inferential methods, it provides a *signal-to-noise ratio* for the hypothesized effect. Large evidence  
70 suggests a strong effect relative to the uncertainty with which it can be estimated.

71

72 Simulations demonstrate the test's efficacy and robustness and compare it to classical frequentist  
73 and Bayesian inferential methods. Further, I apply this test to several concrete examples to show  
74 its advantages in practice. The MATLAB source code and example data are available for download  
75 (<http://dx.doi.org/10.6084/m9.figshare.1342798>) and we also plan to publish a standalone  
76 interface-based version and source code for R at a later date.

77

## 78 **Methods**

79

80 I will refer to the distributions for the effect size under null and alternative hypothesis as  $\theta_0$  and  
81  $\theta_1$ , respectively (Fig. 2, left panels, blue and red curves). To quantify the *similarity* of these two  
82 distributions at each bootstrapping step the difference between the effects under the two  
83 hypotheses is calculated, that is,  $\theta_{1i} - \theta_{0i}$  where  $i$  denotes the  $i$ -th bootstrap step. This produces a  
84 *distribution of differences*,  $\Delta$ , between the effects expected under the two hypotheses (Fig. 2,  
85 right panels). This is necessary because either  $\theta_0$  or  $\theta_1$  can be skewed separately by anomalous  
86 data. The  $\Delta$  distribution captures either of these distortions.

87

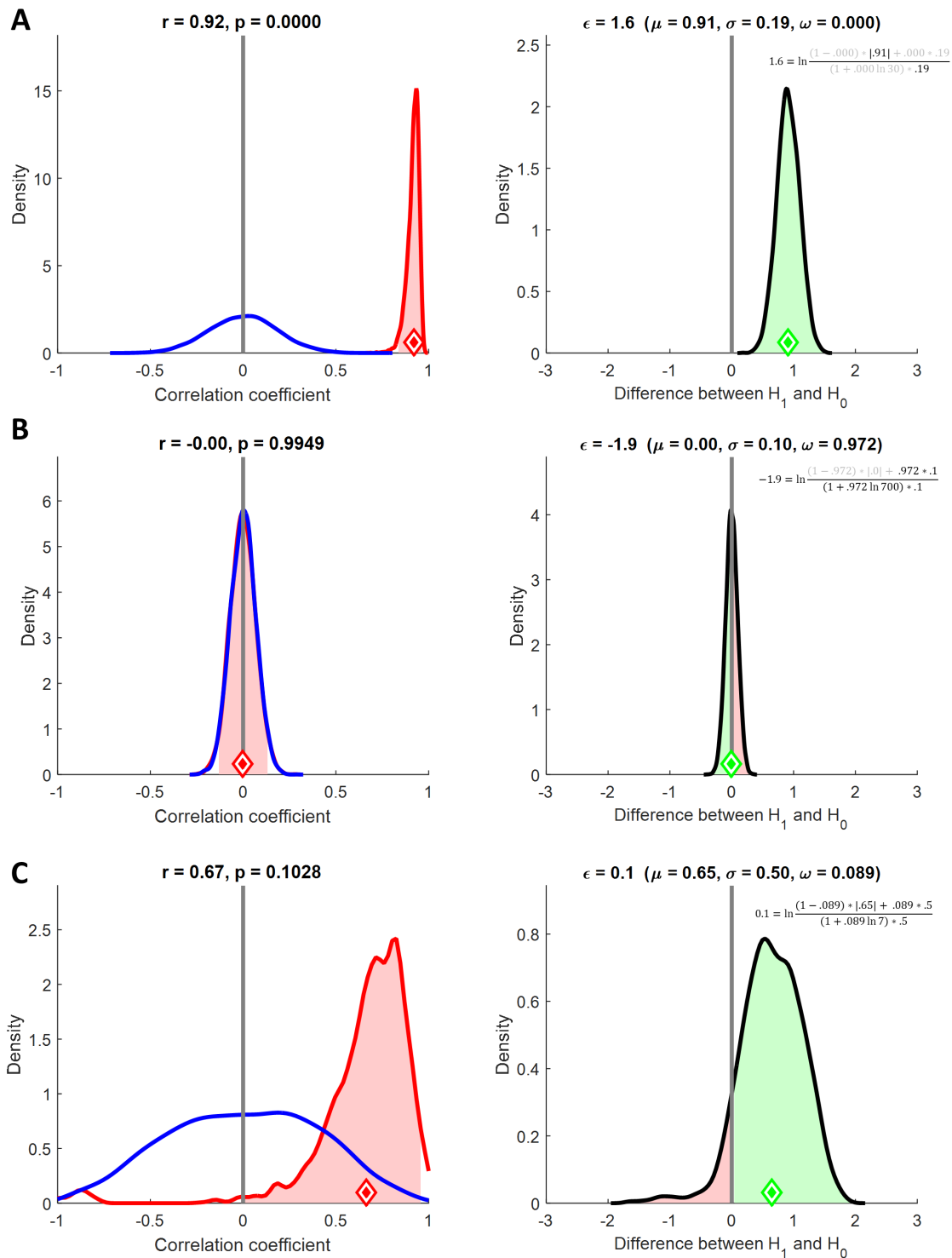
88 Theoretically, the evidence for  $H_1$  is a standardized score based on the  $\Delta$  distribution:

89

$$z = \frac{|\mu|}{\sigma} \tag{1}$$

90

91 where  $\mu$  and  $\sigma$ , respectively, denote the arithmetic mean and standard deviation of  $\Delta$ . This  
92 effectively normalizes the shift of this distribution relative to zero by its dispersion. When  $z$  is  
93 greater than 1 this implies that the sample effect size, and thus our estimate of the true  
94 population effect, is larger than the *uncertainty* of the estimate (Fig. 2A). This provides evidence  
95 supporting the alternative hypothesis. The more data we collect and the sample size,  $n$ , becomes  
96 larger, the more accurate is the estimate of the population effect,  $\mu$ , and the smaller is the  
97 uncertainty,  $\sigma$ .



98  
 99  
 100  
 101  
 102  
 103  
 104  
 105  
 106

**Fig. 2.** Examples of the bootstrapped evidence procedure. *Left panels* show the effect size distribution under the null (blue) and alternative (red) hypothesis estimated by bootstrapping. The red diamond denotes the observed effect size. The red shaded region denotes though the 95% confidence interval of the estimated effect. *Right panels* show the distribution of differences,  $\Delta$ , between the alternative and null distributions (see Methods). The green diamond denotes the mean,  $\mu$ . The ratio of the red and green areas under the curve is  $\omega$  and quantifies the overlap between the null and alternative distributions. The standard deviation of  $\Delta$  is denoted by  $\sigma$ . The bootstrapped evidence,  $\epsilon$ , incorporates these three parameters and the sample size  $n$  and is expressed on a

107 logarithmic scale (see equation). A. Strongly correlated data with  $n=30$ . Evidence compellingly  
108 supports  $H_1$ . B. Uncorrelated data with  $n=700$ . Evidence compellingly supports  $H_0$ . C. Uncorrelated  
109 data with  $n=7$ . Evidence is inconclusive.

110

111 Unfortunately, this only holds when the alternative hypothesis is true. When the null hypothesis  
112 is true instead, that is, when the population effect is zero, the parameters of  $\Delta$  do not reflect the  
113 evidence for  $H_0$ . While increasing  $n$  reduces the uncertainty,  $\sigma$ , on average it also reduces  
114 estimates of the population effect size,  $\mu$ . It follows that the ratio of these parameters,  $z$ , remains  
115 more or less constant.

116

117 Thus, in order to quantify how strongly the data support *either*  $H_1$  or  $H_0$  we must weight the ratio  
118 of  $\mu$  and  $\sigma$  according to how much evidence is available. This is achieved by calculating a third  
119 parameter describing the overlap of  $\theta_0$  and  $\theta_1$ . This is given by:

120

$$\omega = \frac{P(\Delta \operatorname{sgn} \mu \leq 0)}{P(\Delta \operatorname{sgn} \mu > 0)} \quad (2)$$

121

122

123 This divides the proportion of bootstraps in  $\Delta$  that have the opposite sign as  $\mu$  or zero (red area  
124 under the curves in the right panels of Fig. 2) by the proportion of iterations with the same sign  
125 (green area under the curves in the right panels of Fig. 2). The strength of evidence,  $\varepsilon$ , for or  
126 against  $H_1$  is then calculated as:

127

$$\varepsilon = \ln \frac{(1 - \omega)|\mu| + \omega\sigma}{(1 + \omega \ln n)\sigma} \quad (3)$$

128

129

130 While this equation may seem complex, it is essentially the same as  $z$ , the standardized score in  
131 equation 1, but it is moderated by the sample size and the overlap between  $\theta_1$  and  $\theta_0$ . When the  
132 data clearly support the alternative hypothesis, the  $\Delta$  distribution is shifted far away from zero  
133 and thus  $\omega$  (the ratio of the red and green areas under the curve) is very small (Fig. 2A). In fact,  
134 because it is based on the proportion of bootstraps it may equal 0. Under these circumstances  
135 the denominator is close to  $\sigma$ , the numerator is close to  $|\mu|$ , and thus the evidence is  
136 approximately  $z$  from equation 1. However, when the null hypothesis is true, or the effect size is  
137 too small to clearly support the alternative, the  $\Delta$  distribution is centered near zero and thus  $\omega$  is  
138 near 1. In this situation the denominator is a multiple of  $\sigma$ , growing ever larger as the sample size  
139 increases. This in turn ensures that the ratio in equation 3 becomes ever smaller and evidence for  
140 the null hypothesis grows (Fig. 2B).

141

142 The numerator is also moderated by the strength of the evidence. When  $H_0$  is true and  $\omega$  is near  
143 1, the numerator is close to  $\sigma$ . This reflects the fact that when the data provide only weak  
144 evidence, there is substantial uncertainty as to whether the estimate of the effect size is accurate.  
145 When the sample size is large, this means that the denominator dominates the equation and  $\varepsilon$   
146 becomes very small. However, when the sample size is *small*, the ratio in equation 3 is close to 1

147 and thus the data support neither  $H_1$  nor  $H_0$  very clearly – this means the evidence is inconclusive  
148 (Fig. 2C).

149

150 To recap, the bootstrapped evidence,  $\epsilon$ , measures how confident one can be of the effect size  
151 estimate given the uncertainty in the data. If the observed effect is relatively large, the  
152 uncertainty will decrease as sample size grows and thus support for the alternative hypothesis  
153 also grows. However, when the effect remains considerably smaller than the uncertainty, a larger  
154 sample instead provides greater evidence for the null hypothesis.

155

156 Finally, as equation 3 shows,  $\epsilon$  is the *natural logarithm* of this ratio. Therefore, when the ratio is  
157 near 1 and the evidence is inconclusive,  $\epsilon$  is approximately 0. Positive  $\epsilon$  indicates evidence for  $H_1$ ,  
158 while negative  $\epsilon$  indicates evidence for  $H_0$ . The bootstrapped evidence thus provides a non-  
159 dichotomous measure of the evidence for either hypothesis, similar to a Bayes factor [2,20,21].  
160 However, while a Bayes factor is a measure of how much one should update the prior odds due to  
161 the observed evidence, the bootstrapped evidence is in essence a signal-to-noise ratio. A strong  
162 “signal” implies strong evidence for  $H_1$ , while a negligible signal with a lot of data provides strong  
163 evidence for  $H_0$ . The only prior assumptions this procedure makes pertain to how the data are  
164 resampled under the two hypotheses.

165

#### 166 *Bootstrapped correlations*

167

168 To bootstrap the evidence for a linear correlation data are resampled *with replacement* and on  
169 each step the correlation coefficient is computed. To derive the null distribution ( $\theta_0$ ) data are  
170 resampled without restriction as would be expected if the effect occurred by chance, that is,  
171 observations for the two variables are no longer paired but intermixed randomly. This is  
172 essentially standard procedure for non-parametric resampling methods in the classical  
173 frequentist framework (although for this purpose permutation analysis where resampling is  
174 performed *without* replacement is more common). A classical one-tailed p-value could be  
175 calculated by determining the proportion of bootstraps  $\theta_{0i}$  that are at least as large as the  
176 observed effect size – that is, the area under the blue curve to the right of the red diamond.

177

178 However, to derive the alternative distribution ( $\theta_1$ ), quantifying the reliability of the observed  
179 effect, we restrict the resampling strategy on the alternative hypothesis that there is a  
180 correlation. In this case the pairing of data points in each variable is preserved so many resamples  
181 will show a positive linear relationship.

182

183 As described, we next calculate  $\Delta$ , the distribution of differences between  $\theta_1$  and  $\theta_0$  (Fig. 2B). Its  
184 standard deviation is the uncertainty,  $\sigma$ . For all of the examples given in this article, I used 10,000  
185 bootstrap iterations, except in the interest of time for lengthy simulations and curve fitting  
186 examples I only used 1,000 iterations. Reducing the number of iterations only changes the  
187 precision of the estimate of  $\epsilon$  but does not alter the general conclusions substantially.

188

189 The further apart the two distributions for  $\theta_0$  and  $\theta_1$  are, the farther  $\Delta$  is from zero and the  
190 greater is the evidence for  $H_1$ . This is quantified by  $\epsilon$ . When  $\epsilon$  is very negative, the evidence favors  
191  $H_0$  because it means the two distributions for  $\theta_0$  and  $\theta_1$  overlap considerably which means that  $\Delta$

192 is centered near zero. Intuitively this indicates that the effect size estimate under  $H_1$  could very  
193 likely have been smaller than that under  $H_0$ . When the evidence is  $-0.5 < \epsilon < 0.5$  this provides  
194 inconclusive support for the either hypothesis. This region is somewhat arbitrary but it reflects  
195 the fact that while there is overlap between the distributions for  $\theta_0$  and  $\theta_1$ , there is not sufficient  
196 data to be confident that there is no subtle effect. This corresponds to the range of  $\epsilon$  one typically  
197 obtains with small sample sizes when the null hypothesis is true.

198

### 199 *Bootstrapping differences*

200

201 Naturally, the same procedure can be applied to other statistical comparisons in addition to tests  
202 of correlation. For instance, when *comparing the means of two independent samples* the effect  
203 size is the difference between the sample means. To estimate the null distribution,  $\theta_0$ ,  
204 observations are resampled with replacement and divided into new samples of the same size as  
205 the original samples. To estimate the distribution for the alternative hypothesis,  $\theta_1$ , the  
206 segregation of the two samples is maintained and resampling is done *within* each sample. In all  
207 other respects, the procedure is identical to the correlation test already described.

208

209 When testing the *difference between two repeated measures* the same underlying principle  
210 applies. Here the effect size is the mean over the *differences* in each pair of observations. We  
211 estimate  $\theta_1$  by maintaining the pairing but resampling with replacement. For estimating  $\theta_0$  the  
212 pairing is also kept intact because what matters is only the variance across repeated measures.  
213 However, under the null hypothesis the order of measures is irrelevant so the resampling  
214 *randomizes the sign* for each observation. This corresponds to scrambling the order of  
215 observations in a repeated measures design.

216

217 Similarly, the BSE can also be used to test the *difference between two correlations*. Again the data  
218 need to be resampled based on the assumptions of the null as well as the alternative hypothesis.  
219 In this case the null hypothesis resamples data ignoring how variables are paired while the  
220 alternative hypothesis preserves the pairing. The estimated effect size is the difference between  
221 the two correlation coefficients.

222

223 The situation becomes more complicated for testing the *difference of one sample from a fixed*  
224 *value* (e.g., when comparing a normative sample to a patient case-study). Conceptually, it is  
225 possible to use the same resampling strategy as for a repeated measures design: the observations  
226 are the differences from the fixed value and for resampling the null distribution  $\theta_0$  we randomize  
227 the sign of each observation. However, this approach is probably not sufficiently conservative.  
228 While it is conceptually correct for repeated measures designs to assume a mean difference of 0  
229 under the null hypothesis, in many other situations fixed values are themselves subject to  
230 variability and/or measurement error. For instance, a measurement in a case study is subject to  
231 within-subject variability and chance performance in a behavioral task follows a probability  
232 distribution. Using a similar approach one can also incorporate variability in individual  
233 observations to calculate a group mean. In each round of bootstrapping we can simulate a new  
234 sample by drawing random data using the individual means and variances for every observation.  
235 This approach may be particularly suitable for meta-analysis.

236

237 *Bootstrapping tests against chance*

238

239 Simulating the null distribution is also suitable for testing binomial processes, such as whether a  
240 coin is fair or whether a participant performed better than chance at a behavioral task. As usual,  
241 in each bootstrap step the observed data (e.g. a series of 1s and 0s for heads or tails) are  
242 resampled to obtain the alternative distribution  $\theta_1$ . However, for estimating  $\theta_0$  we instead  
243 generate a *new* set of 1s and 0s of the same number as the observed trials using the chance  
244 probability (i.e. 0.5 or whatever the chance level is). Alternatively, one can also permute the raw  
245 trial data in each resampling step and recalculate the accuracy. The latter approach is advised  
246 when the design is unbalanced or if there is any suspicion that chance may not have a binomial  
247 distribution. In all other ways the procedure works as described.

248

249 A very similar approach for simulating a chance distribution based on the assumptions underlying  
250 the null hypothesis can also be used for other problems, for example comparing the performance  
251 of a group of participants against chance. In this case, we can simulate  $\theta_0$  by generating a new set  
252 of 'chance' participants at each resampling step under the same conditions as the actual  
253 experiment (same chance probability, number of trials, and number of participants).

254

255 *Bootstrapping curve fits*

256

257 The bootstrapped evidence procedure also affords itself easily for curve fitting or regression  
258 analyses. The estimated effect size in this case is the coefficient of determination,  $R^2$  (or  
259 goodness-of-fit). Otherwise the procedure works in much the same way as for calculating  
260 correlations. Under the null hypothesis the observations for the dependent and independent  
261 variables are scrambled randomly with replacement. Under the alternative hypothesis, the  
262 pairing is kept intact but observations are resampled with replacement.

263

264 **Results**

265

266 The principle underlying the BSE test is that under both the alternative ( $H_1$ ) and the null  
267 hypothesis ( $H_0$ ) the results follow a probability distribution (Fig. 2, left panels). In classical  
268 statistics, the p-value reflects the distance of the observed effect,  $\theta$  (e.g. a correlation coefficient),  
269 from the center of the null distribution. The one-tailed p-value is the area under the blue curve  
270 (null distribution) to the right of the red diamond, which denotes the observed effect.

271

272 While the null distribution depends on the sample variance, it nonetheless fails to take the  
273 *variability of the effect under  $H_1$*  into account. The BSE estimates how distinct these two  
274 distributions are from one another by bootstrapping both  $H_0$  and  $H_1$  and quantifying  $\Delta$ , the  
275 distribution of differences (Fig. 2, right panels), between them (see Methods). If the distribution is  
276 narrow and shifted away from zero this constitutes evidence for  $H_1$  (Fig. 2A). However, when the  
277 distribution is narrow but centered on zero this is instead evidence for  $H_0$  (Fig. 2B). A wide  $\Delta$   
278 distribution provides only inconclusive evidence (Fig. 2C).

279

280 The BSE is expressed by  $\epsilon$ , which is effectively a *signal-to-noise ratio on a logarithmic scale*. It is  
281 the ratio of the observed effect,  $\mu$ , and the uncertainty,  $\sigma$ , with which it can be estimated (Fig. 2



282 and Methods). When  $\mu$  is smaller than  $\sigma$  and thus the  $\Delta$  distribution overlaps zero (as quantified  
283 by  $\omega$ ),  $\epsilon$  decreases as sample size,  $n$ , increases. Since  $\epsilon$  is the logarithm of this ratio, it is positive  
284 when the data support  $H_1$  and negative when they favor  $H_0$ . If  $\epsilon$  is near zero ( $-0.5 < \epsilon < 0.5$ ) the  
285 evidence is inconclusive.

286

287 For the near perfect correlation with  $n=5$  (Fig. 1A) the evidence is only  $\epsilon=0.7$ . In contrast, for the  
288 modest correlation with  $n=40$  the evidence  $\epsilon=1$  (Fig. 1B). The data in Fig. 1C,D are uncorrelated  
289 and neither correlation would reach classical significance. However, for a small sample size of  $n=5$   
290 (Fig. 1C) the evidence is inconclusive ( $\epsilon=-0.4$ ) while for a large sample (Fig. 1D) the evidence  
291 compellingly favors the null hypothesis ( $\epsilon=-1.6$ ).

292

293 The assumption made by parametric tests of normally distributed errors is often violated as in Fig.  
294 1E. Even though the classical  $p$ -value is highly significant ( $r=-0.45$ ,  $p=0.0034$ ), the evidence for  $H_1$   
295 is fairly weak ( $\epsilon=0.6$ ). The reason is that the data are heteroscedastic and thus skew classical  
296 Pearson's correlation: the residuals of a linear fit depend on  $x$ . While  $y$  is chosen from a random  
297 normal distribution each point is also multiplied by the absolute magnitude of its paired value in  $x$   
298 [15]. Such situations can readily occur in real experimental data: for instance, the proliferation of  
299 cell growth or the mean firing rate of neurons may also be accompanied by greater variability in  
300 those measures. This in turn could skew any correlations between these measures and an  
301 independent variable.

302

303 Notably, even increasing the sample size does not alleviate this problem. The data in Fig. 1F were  
304 drawn from the same heteroscedastic population but the sample size is five times larger ( $n=200$ )  
305 and the correlation is again highly significant ( $r=0.22$ ,  $p=0.0014$ ). Even robust significance tests,  
306 including those specifically developed to control for heteroscedasticity, do not fare any better  
307 (skipped correlation [22–24]:  $t=4.03$ ,  $t_{\text{critical}}=2.35$ ; permutation test:  $r=0.22$ ,  $p=0.0013$ ; Spearman's  
308 rho:  $\rho=0.18$ ,  $p=0.0105$ ; Kendall's tau  $\tau=0.14$ ,  $p=0.0032$ ; percentage bend correlation [22]:  $r=0.2$ ,  
309  $p=0.0042$ ; Shepherd's pi [25]:  $\pi=-0.23$ ,  $p=0.0034$ ). In contrast, the BSE test suggests only  
310 inconclusive evidence for  $H_1$  ( $\epsilon=0.4$ ) because bootstrapped distributions for  $H_1$  and  $H_0$  overlap  
311 substantially. In comparison, a homoscedastic data set with the same effect and sample size  
312 would produce more convincing evidence for  $H_1$  ( $\epsilon=0.7$ ).

313

#### 314 *Performance on simulated data*

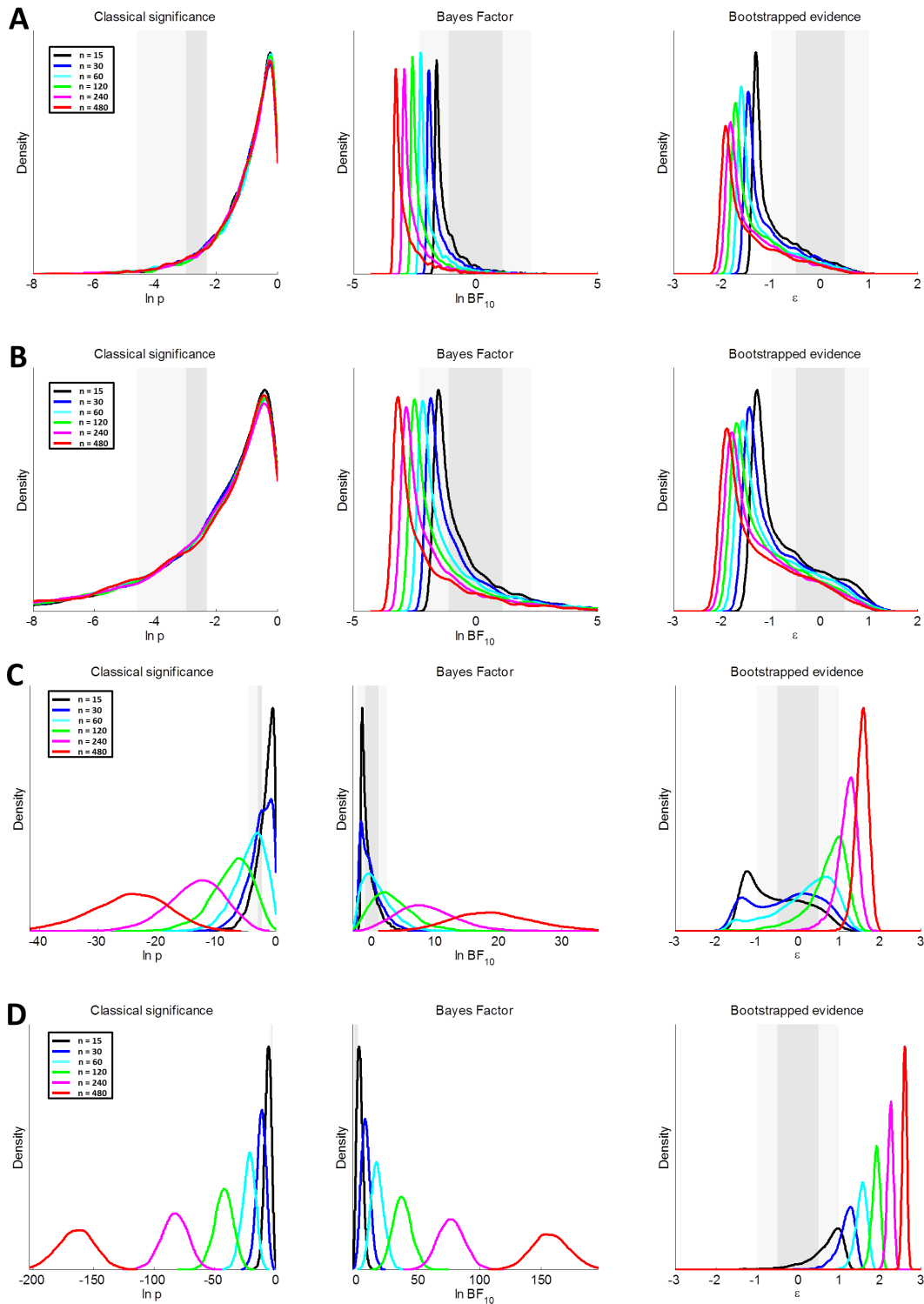
315

316 For a more objective evaluation of the method I ran a series of simulations. For several sample  
317 sizes ( $n=15, 30, 60, 120, 240$ , and  $480$ ) I generated 5,000 data sets each drawn from two different  
318 distributions in which  $H_0$  is true: an uncorrelated bivariate Gaussian and the same heteroscedastic  
319 distribution underlying Figures 1E,F. For each simulated data set I calculated the bootstrapped  
320 evidence,  $\epsilon$ , the classical parametric  $p$ -value and a default Bayes factor for  $H_1$  over  $H_0$  ( $BF_{10}$ ) [26].

321

322 Fig. 3A shows the distribution of these inferential statistics for uncorrelated Gaussian data with  
323 the various sample sizes denoted by different colors. As sample size increases the distributions for  
324 the default Bayes factor and bootstrapped evidence become increasingly shifted towards  
325 negative numbers, indicating increasing support for  $H_0$ . In contrast, the distributions for classical  
326  $p$ -values remain the same irrespective of sample size because under  $H_0$  the distribution of  $p$ -

327 values is uniform (because the x-axis is logarithmic this manifests as a long leftwards tail). This  
328 ensures that, provided the assumptions of the test are met, the false positive rate in classical  
329 statistics is constant across sample sizes when  $H_0$  is true. This illustrates why classical  $p$ -values *can*  
330 *never provide evidence for the null hypothesis* [20]. When  $H_0$  is true, a proportion of tests given by  
331 the  $\alpha$  level will be false positives. Because the estimated effect size with large sample sizes is  
332 typically very small (i.e. close to the truth of zero effect), trivially tiny effects may thus become  
333 statistically significant.



334

335 **Fig. 3.** Distributions of statistical evidence in 5,000 simulations of uncorrelated Gaussian data (A),  
336 severely heteroscedastic data (B), weakly correlated data with  $\rho=0.3$  (C) or strongly correlated data  
337 with  $\rho=0.7$  (D). Six sample sizes were tested (see color code). Each panel shows distributions for the  
338 classical p-value (left), default Bayes factor [26]  $BF_{10}$  (middle), and bootstrapped evidence  $\epsilon$  (right).  
339 The dark shaded regions denote “inconclusive” results (see text). The light shaded regions denote  
340 results that pass a basic criterion but provide no strong evidence for a given hypothesis.

341

342 The grey shaded regions in each panel indicate the boundaries of commonly used criterion levels.  
343 For classical statistics the dark grey region corresponds to p-values between 0.05-0.1, sometimes  
344 called “marginally significant.” The light grey region denotes the range between 0.01-0.05. Any p-  
345 value to the left of the light grey region would constitute a significant result. For Bayes factors  
346 and the bootstrapped evidence the regions are symmetric around 0. The dark grey region  
347 corresponds to inconclusive evidence that supports neither  $H_1$  nor  $H_0$  (i.e.  $\frac{1}{3}$ -3 for  $BF_{10}$ , -0.5-0.5  
348 for  $\epsilon$ ). The light grey regions refer to evidence that passed the criterion but which is still relatively  
349 weak (i.e.  $BF_{10}$  between  $10^{-1}$  and  $\frac{1}{3}$  or 3 and 10;  $\epsilon$  between -1 to -0.5 or 0.5 to 1). The proportion  
350 of these statistics to the right of the criterion becomes smaller as sample size increases.

351

352 Fig. 3B shows results from simulations using the uncorrelated but severely heteroscedastic  
353 distribution. For all sample sizes the distributions for classical p-values are biased so that false  
354 positives are drastically inflated. The same is also somewhat true for the default Bayes factor and  
355 bootstrapped evidence. However, for the BSE the skew is at worst modest, while the proportion  
356 of Bayes factors exceeding “strong” evidence for  $H_1$  is larger. This is because the default Bayes  
357 factor is a function of the Pearson’s correlation coefficient and the sample size. It is therefore  
358 skewed by heteroscedasticity in the same way as classical statistics. However, since the BSE is  
359 based on non-parametric resampling it is less affected by violations of parametric assumptions.

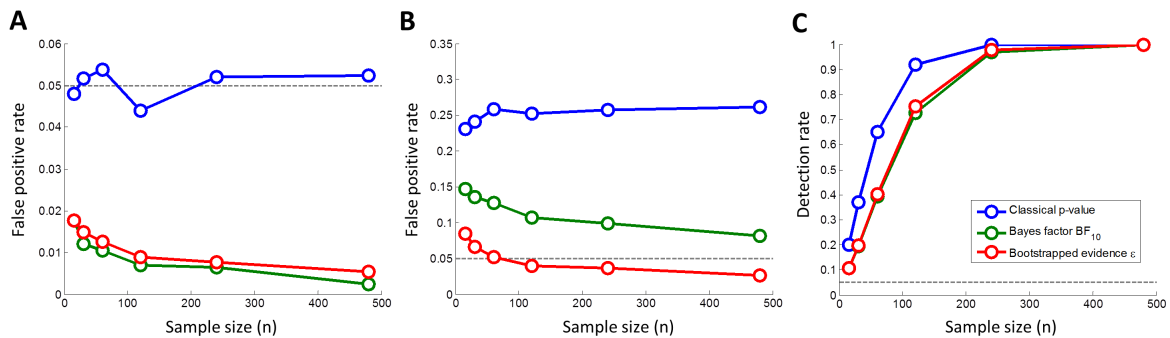
360

361 Next I performed a sensitivity analysis determining how well the BSE test detects true effects. I  
362 repeated the same kind of simulation but now data were chosen from a Gaussian bivariate  
363 distribution with population correlations of  $\rho=0.3$  or  $\rho=0.7$ . Unsurprisingly, for all of the three  
364 procedures the evidence for  $H_1$  becomes stronger as the sample size increases. For  $\rho=0.3$  the  
365 evidence passes criterion only for larger samples sizes (Fig. 3C) while for most data both Bayesian  
366 and bootstrapped evidence remains inconclusive. Classical statistics are less conservative as the  
367 peak of the distribution with  $n=120$  (green curve) is already below the  $p<0.01$  threshold. For  
368  $\rho=0.7$  the evidence with most sample sizes passes criterion (Fig. 3D).

369

370 I summarized the false positive and correct detection rates as a function of sample size. As  
371 expected, for classical statistics the false positive rate remains constant near the nominal level of  
372 5% across all sample sizes, if data are Gaussian and homoscedastic (Fig. 4A). For  $BF_{10}$  and  $\epsilon$ , false  
373 positive rates for standard criteria ( $BF_{10}>3$  and  $\epsilon>0.5$ , respectively) decrease as sample size  
374 increases. For either method the false positive rate is already below 5% even at the smallest  
375 sample size ( $n=15$ ). When heteroscedasticity is present, false positives are dramatically inflated:  
376 approximately one in four tests are positive at  $p<0.05$  (Fig. 4B). Both evidence-based methods  
377 also show some inflation; however, false positives for the BSE are only about half that for the  
378 default Bayes factor. For the smallest sample size ( $n=15$ ) the worst false positive rate for  $\epsilon>0.5$  is  
379  $\sim 8.5\%$  compared to  $\sim 14.7\%$  for  $BF_{10}>3$ . When there is a real effect ( $\rho=0.3$ ) the detection rate rises

380 steeply and then saturates for all three methods but Bayes factors and BSE are more conservative  
 381 than classical p-values (Fig. 4C).  
 382



383  
 384 **Fig. 4.** Detection rates from the simulations in Fig. 3 plotted against sample size for classical  $p < 0.05$   
 385 (blue), default Bayes factor [26]  $BF_{10}$  (green), and the bootstrapped evidence  $\epsilon$  (red). A. Uncorrelated  
 386 Gaussian data. B. Severely heteroscedastic data. C. Correlated data with  $\rho = 0.3$ .  
 387

388 While the conclusions one would draw from all three approaches are usually similar, one notable  
 389 difference is evident between classical and Bayesian inference and the bootstrapped evidence:  
 390 distributions for  $\epsilon$  tend to become *narrower as sample/effect size increase*. In contrast, the  
 391 distributions for p-values and  $BF_{10}$  become wider. Note that all of these plots are on logarithmic  
 392 scales (log-transformation is inherent to the calculation of  $\epsilon$ ; see Fig. 2 and Methods). Despite  
 393 this, the distributions for p-values and Bayes factors display extraordinary variability, e.g. the  
 394 distribution for  $\rho = 0.3$  at the largest sample size of  $n = 480$  (Fig. 3C, red curves). Here 95% of  
 395 simulated p-values are between  $8.5 \times 10^{-18}$  and  $1.8 \times 10^{-06}$ . All are highly significant at  $p < 0.001$  but  
 396 this range spans many orders of magnitude. The default Bayes factor behaves similarly. The  
 397 equivalent range spans  $BF_{10}$  between 3,212 and 378 quadrillion. Any of these would constitute  
 398 “decisive” evidence of  $BF_{10} > 100$  [26,27]. But from a pragmatic view, how much more confident  
 399 should we be of the highest Bayes factor in this range compared to the lowest? In comparison,  
 400 the equivalent range for the bootstrapped evidence is between 1.2 and 1.9. Again, these are well  
 401 above even a strict criterion of  $\epsilon > 1$  but there is no stark discrepancy between the weakest and  
 402 strongest evidence. This is because rather than determining probabilities it reflects the precision  
 403 with which the population effect size can be estimated. The precision increases with sample size.  
 404 Thus, replicate experiments will produce very consistent bootstrapped evidence for  $H_1$ .  
 405

406 Naturally, Bayesian analysis depends on the choice of a prior but typically with a range of default  
 407 priors the outcome usually does not vary qualitatively [3]. Nonetheless, choosing a prior could  
 408 theoretically also lead to substantial analytic flexibility, thus inflating the “researcher degrees of  
 409 freedom” [28]. The BSE test on the other hand makes no assumptions beyond the resampling  
 410 strategy needed for either hypothesis.  
 411

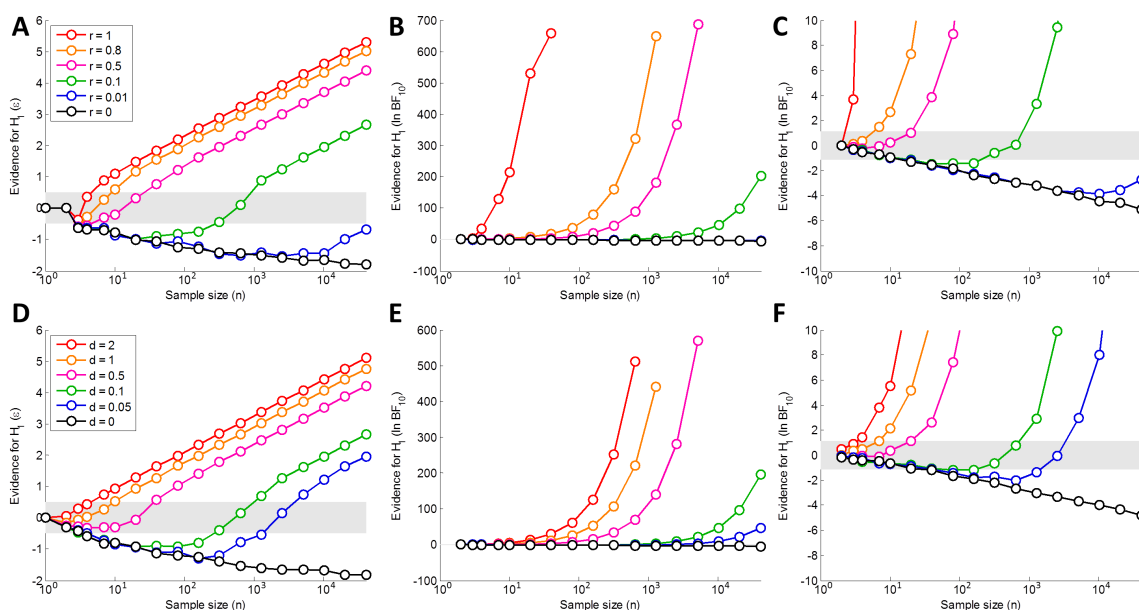
412 *Evidence as a function of sample size*

413  
 414 Both classical statistics and the default Bayes factor also place undue confidence on strong effects  
 415 when sample sizes are small as in Fig. 1A. The Bayes factor is rather large  $BF_{10} = 90.2$  while the BSE  
 416 is modest ( $\epsilon = 0.7$ ). The Bayes factor reflects how much more probable the data are under  $H_1$  than  
 417  $H_0$  [29]. However, from a pragmatic perspective this could nonetheless be problematic Combining

418 publication bias towards positive findings with underpowered experiments, high Bayes factors  
 419 may thus be misinterpreted as strong evidence for the alternative hypothesis. Given the problems  
 420 with spurious results and reproducibility in the scientific literature [30] this could be problematic.

421  
 422 Fig. 5 plots the evidence for a range of effect sizes against sample size. In most situations, the  
 423 conclusions we would draw from bootstrapped evidence (Fig. 5A) and the default Bayes factor  
 424 (Fig. 5B,C) are largely the same. For strong effects, the evidence rises continuously beyond the  
 425 inconclusive region, while for weaker effects the evidence starts off as indistinguishable from the  
 426 situation when the null hypothesis is true (black curves) until it departs and also rises. This  
 427 behavior is natural because if the true effect is weaker than what could be meaningfully detected  
 428 given the data at hand this constitutes support for the null hypothesis.

429



430  
 431 **Fig. 5.** Statistical evidence for H<sub>1</sub> plotted against sample size for a range of effect sizes (see color  
 432 code). A-C. Correlation analysis. D-F. Comparing the means of two samples. Individual panels show  
 433 the bootstrapped evidence  $\epsilon$  (A,D) or the default Bayes factor [26]  $BF_{10}$  (B-C, E-F). Panels C and F plot  
 434 the Bayes factor with y-axis zoomed in on zero. The shaded grey region denotes “inconclusive”  
 435 evidence (i.e.  $-0.5 < \epsilon < 0.5$  or  $1/3 < BF_{10} < 3$ , respectively). For the bootstrapped evidence (A,D) these data  
 436 represent the mean across 100 simulations.

437  
 438 The slopes of the curves for the bootstrapped evidence are far less steep. Thus it is possible to see  
 439 the behavior for the full range of conditions within the same plot. There is however one  
 440 considerable difference: for a perfect correlation ( $\rho=1$ ) the default Bayes factor immediately rises  
 441 even at tiny sample sizes (Fig. 5C). At  $n=3$  the  $BF_{10}$  is already 48.8. In contrast, the BSE for this  
 442 point is low ( $\epsilon=-0.4$ ) and inconclusive. As sample size increases, so does the bootstrapped  
 443 evidence. At  $n=4$  it is still inconclusive but favoring H<sub>1</sub> ( $\epsilon=0.4$ ). At  $n=7$  it clearly supports H<sub>1</sub> ( $\epsilon=0.9$ )  
 444 and it continues to rise as sample size increases. This behavior is more intuitive than that of the  
 445 default Bayes factor and also the classical p-value, which would be extremely significant in all  
 446 these situations. Compare this to the earlier example of a strong correlation ( $r=0.9953$ ) with a  
 447 small sample size of  $n=5$  (Fig. 1A). Classical inference would be extremely significant ( $p < 0.001$ )  
 448 and the default Bayes factor would yield “very strong” evidence for H<sub>1</sub> ( $BF_{10}=90.2$ ). The BSE is

449 however only fairly modest, especially given the strong effect ( $\epsilon=0.7$ ). It is above the criterion for  
450 conclusive evidence but it does not instill undue confidence in  $H_1$ .

451

452 The data in this example were in fact drawn from an uncorrelated Gaussian distribution so the  
453 null hypothesis was true. The bootstrapped evidence provides an intuitive measure of the  
454 weakness of the evidence in such situations and should thus be a safeguard against weak or  
455 inconclusive results.

456

#### 457 *Simulations of optional stopping*

458

459 The BSE test has further advantages over classical inference based on significance thresholds. In  
460 classical statistics, even when there is no true effect, it is theoretically possible to reach an  
461 arbitrarily significant p-value, if data collection continues until the p-value passes the significance  
462 threshold. This is known as “optional stopping”, which is an incorrect but possibly widespread use  
463 of classical statistics [4,31]. Under the classical framework one should first define the expected  
464 effect size *a priori*, perform a power analysis to see how large a sample is needed to detect this  
465 effect with sufficiently high probability, and then collect those data without stopping until the  
466 sample is complete. However, typically this is not realistic as one can often only make a vague  
467 guess about the expected effect size.

468

469 The bootstrapped evidence does not suffer from this conundrum. First, even if a dubious optional  
470 stopping strategy is used, the false positive rate is not inflated substantially. I simulated this 1,000  
471 times by drawing data repeatedly from an uncorrelated bivariate Gaussian distribution thus  
472 successively increasing the sample size by 1, starting with a minimal sample of  $n=5$ . At each step I  
473 applied classical statistics, the default Bayes hypothesis test [26], and the BSE test. The first  
474 instance one of these tests passed the criterion level, that is  $p<0.05$  for classical statistics,  $BF_{10}>3$   
475 or  $BF_{10}<\frac{1}{3}$  for Bayes factors, and  $\epsilon>0.5$  or  $\epsilon<-0.5$  for the bootstrapped evidence, I recorded the  
476 measure of evidence. In addition, I also performed this procedure on the bootstrapped  
477 uncertainty and recorded the first instance that  $\sigma<0.2$ . If none of the measures reached criterion  
478 simulated data collection would cease at  $n=150$ .

479

480 Under the assumptions of classical statistics this false positive rate should be near 5%, however,  
481 the actual probability was much greater, 36.4%, illustrating the considerable problems optional  
482 stopping can cause in classical inference. In contrast, the false positive rates of the evidence-  
483 based methods was much lower and well below the classical  $\alpha$  level (Bayes factor: 2.9%;  
484 bootstrapped evidence: 2.8%).

485

486 I repeated the same simulation but this time drawing from the heteroscedastic distribution used  
487 in previous examples (Figures 1E,F). Now classical statistics massively inflated support for the  
488 alternative hypothesis with a false positive rate of 84.6%. Default Bayes factors fared a lot better  
489 but are nonetheless strongly skewed (20.6%). For the bootstrapped evidence on the other hand  
490 the false positive rates were only half that (10.3%), again reflecting the fact that it is based on a  
491 non-parametric procedure that takes into account the anomalous distribution of the data.

492

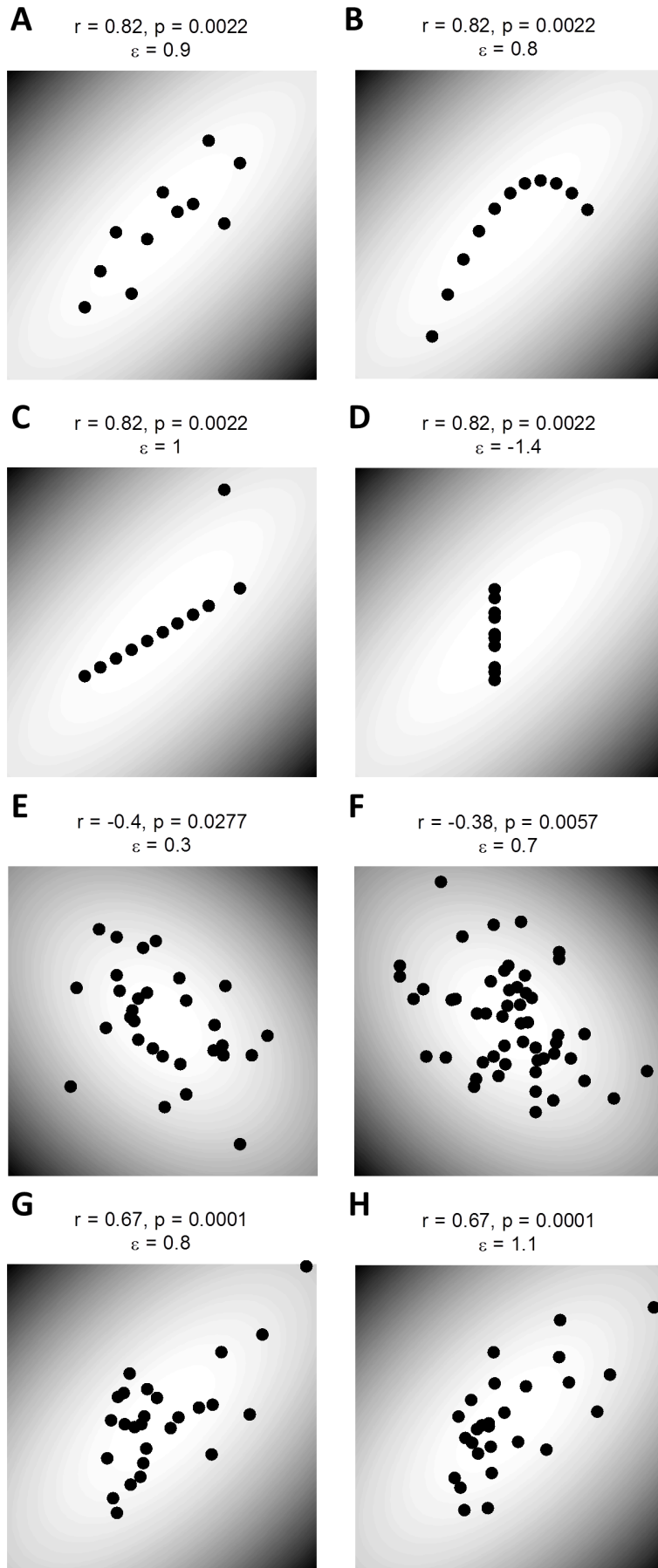
493 This demonstrates that optional stopping based on symmetric evidence is far less problematic  
494 than for classical statistics. In particular, sequential analysis until the bootstrapped evidence  
495 reaches conclusive support for either  $H_1$  or  $H_0$  results in only minimal false positive rates even in  
496 extreme situations. However, there is an even better optional stopping strategy that could be  
497 employed in the bootstrapped evidence framework. When data collection continued until the  
498 bootstrapped uncertainty,  $\sigma$ , was 0.2, the false positive rate for using  $\epsilon > 0.5$  in the first scenario  
499 (homoscedastic Gaussian data) was only 1.3%, while for the heteroscedastic data it was 7.9%.  
500 This suggests that using a criterion uncertainty level is the most optimal strategy for minimizing  
501 spurious findings in sequential analysis.

502

503 *Example 1: Anscombe's quartet*

504

505 Simulations are crucial for testing a method's performance because the ground truth is known.  
506 However, for illustration I also apply the method to Anscombe's quartet [32] a famous  
507 demonstration of the pitfalls of correlation analysis. It consists of four data sets, each comprising  
508 11 pairs of variables, in which Pearson's correlation produces (almost) identical results ( $r=0.82$ ,  
509  $p=0.002$ ). Applying the BSE test reveals that while the data afford low but sufficient confidence  
510 for the correlation in the first three data sets (Fig. 6A-C), in the final example (Fig. 6D) the  
511 evidence clearly supports  $H_0$  ( $\epsilon=-1.4$ ) because one influential outlier drives the correlation but the  
512 remaining data are uncorrelated.





514 **Fig. 6.** Example data sets. A-D. Anscombe's quartet: Each data sets has approximately the same  
515 Pearson's correlation ( $r=0.82$ ,  $p=0.0022$ ) and thus all have a default Bayes factor [26]  $BF_{10}\approx 23$ . Panels  
516 show Typical Gaussian data (A), data showing a perfect non-linear relationship (B), a perfect  
517 correlation contaminated by one influential outlier (C) and uncorrelated data contaminated by one  
518 influential outlier (D). E-H. Experimental data showing correlations between visual cortical surface  
519 area and perceptual function. Correlations between V1 area and Ebbinghaus illusion strength from  
520 [33] (E) and [34] (F). Correlations from [35] between travelling wave speed in binocular rivalry and  
521 the surface areas of V1 (G) and V2 (H). All other conventions as in Fig. 1.

522

523 Interestingly, the confidence in  $H_1$  is actually subtly greater ( $\epsilon=1$ ) for the third example (Fig. 6C)  
524 than the first ( $\epsilon=0.9$ , Fig. 6A). This is because a single outlier contaminates the perfect correlation  
525 in this example, whereas the first example contains noisy but normally distributed data.

526

527 The BSE does not distinguish strongly between the first and second examples (Fig. 6A,B). The  
528 second example contains a perfect relationship between  $x$  and  $y$ ; however, it does not conform to  
529 the linear relationship assumed by Pearson's correlation. Curve fitting can also be implemented in  
530 the BSE framework (see Methods). Here we could compare a simple linear fit to polynomial  
531 curves. The evidence for  $H_1$  with a second-order polynomial is considerably greater ( $\epsilon=1.6$ ) than  
532 for a standard linear model ( $\epsilon=1$ ). Interestingly, the BSE is also robust to overfitting more complex  
533 models: the evidence for higher-order polynomials is weaker than for the second-order (third-  
534 order:  $\epsilon=1.3$ ; fourth-order:  $\epsilon=1$ ).

535

536 *Example 2: Links between visual cortex surface area and perceptual function*

537

538 I further applied the BSE test to published experimental data that showed correlations between  
539 the size of early visual areas and perceptual function. These studies hypothesized that the  
540 transmission speed/strength of lateral connections running tangential to the cortical surface is  
541 reduced for individuals with larger cortical surface areas. In the first two studies, this should  
542 manifest as an anti-correlation between the strength of the Ebbinghaus illusion and V1 surface  
543 area [33,34]. Classical statistics confirmed this hypothesis in both studies (Fig. 6E,F). However,  
544 according to the BSE the findings of the initial study were inconclusive ( $r=-0.4$ ,  $p=0.028$ ,  $\epsilon=0.3$ ).  
545 The second study used a more sophisticated design producing more compelling evidence for this  
546 link ( $r=-0.38$ ,  $p=0.006$ ,  $\epsilon=0.7$ ; note, however, that this study also normalized V1 area by the whole  
547 cortical surface area to control for non-linearity issues and other confounds. For the sake of  
548 consistency with the other findings I chose not to apply this correction here).

549

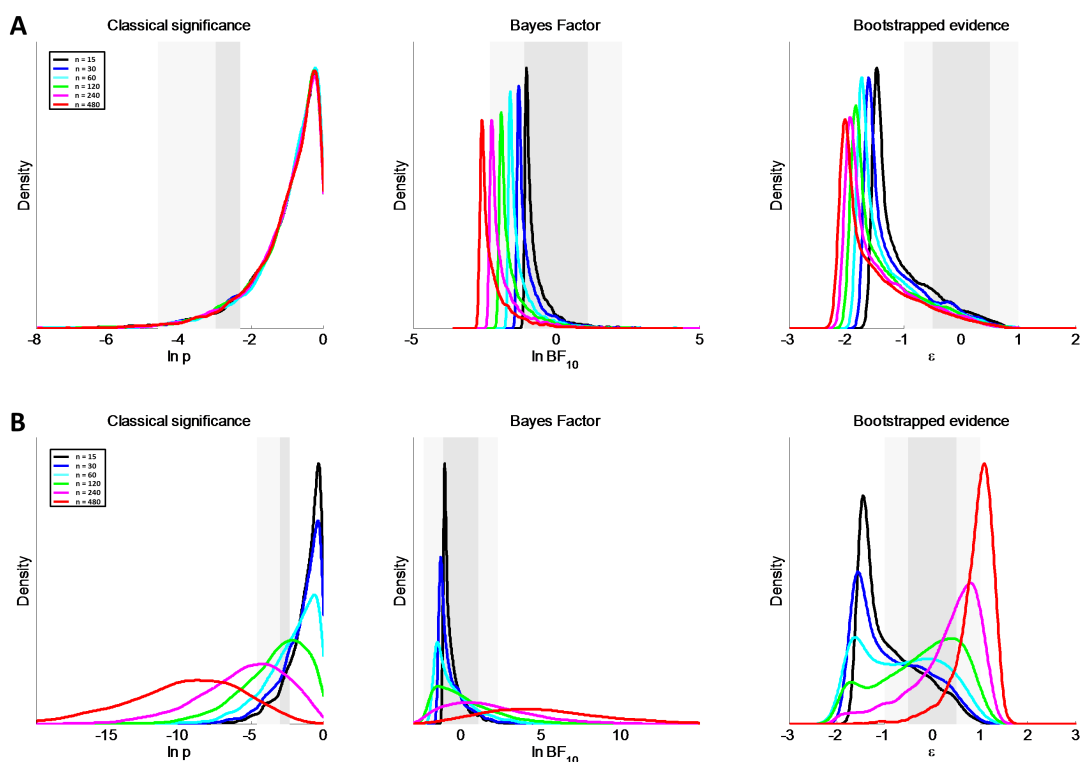
550 The third study [35] reported a linear relationship between the speed of travelling waves in  
551 binocular rivalry and the surface areas of V1 and V2. Classical statistics were very similar for both  
552 regions ( $r=0.67$ ,  $p=0.0001$ ). However, the bootstrapped evidence was in fact lower for V1 (Fig. 6G;  
553  $\epsilon=0.8$ ) than V2 (Fig. 6H;  $\epsilon=1.1$ ), possibly because influential outliers affected the former.

554

555 *Other statistical tests*

556

557 The BSE test can also address many other questions. One simply needs to change how the effect  
558 size is calculated and how data are resampled during bootstrapping (see Methods). For example I  
559 also ran simulations for comparing the means of two samples (Fig. 5 and Supplementary Fig. 1).



560  
 561 **Supplementary Figure 1.** Distributions as shown in Figure 4 but for comparing the means of two  
 562 samples. A. No difference, i.e. normally distributed data with unit standard deviation and a  
 563 population mean of 0. B. A weak effect with a population difference of 0.25 and unit standard  
 564 deviation. All other conventions as in Figure 4.

565  
 566 *Example 3: reassessing evidence for precognition*

567  
 568 A few years ago a psychology study reported experimental evidence for the proposition that  
 569 participants had precognitive abilities [16]. This study was criticized as an illustration of the  
 570 shortcomings of classical inference: Bayesian reanalysis found little evidence in favor of  
 571 precognition [3]. However, a Bayesian analysis by the original author supported his claims [17].  
 572 The conclusions are evidently dependent on the exact prior chosen [3,18]. A cautious prior seems  
 573 advisable when evaluating such unusual claims but the debate illustrates how Bayesian inference  
 574 can appear to lack objectivity.

575  
 576 Here I reevaluate claims from a more recent study on precognitive abilities [36] with the BSE test  
 577 (see Methods). It is a perfect test of the method because these results seem biologically and  
 578 physically implausible but both classical and default Bayesian inference [3] nonetheless support  
 579 the alternative hypothesis for several of the studies (Table 1). In contrast, the BSE does not  
 580 provide convincing support for precognition. It is noteworthy that despite a large sample size of  
 581 1222 participants, even the web-based study 3 only produced inconclusive evidence ( $\epsilon=-0.1$ ). This  
 582 illustrates that in comparison to classical inference, the BSE is far less susceptible to the inflation  
 583 of “significant” findings with large sample sizes. Taken together this suggests that the available  
 584 data do not provide conclusive evidence for either  $H_1$  or  $H_0$ .

585

Experiment	Classical statistics	Bayes Factor, $BF_{10}$	Bootstrapped evidence, $\epsilon$
Maier Study 1	<b>t(110)=-2.65, p=0.0092</b>	2.9 (anecdotal $H_1$ )	0.2 (inconclusive)
Maier Study 2	<b>t(200)=-2.99, p=0.0032</b>	<b>5.9 (moderate <math>H_1</math>)</b>	0.3 (inconclusive)
Maier Study 3	<b>t(1221)=-2.36, p=0.0186</b>	0.5 (anecdotal $H_0$ )	-0.1 (inconclusive)
Maier Unsuccessful 1	t(62)=0.16, p=0.8700	0.14 (moderate $H_0$ )	-1.8 (compelling $H_0$ )
Maier Unsuccessful 2	t(405)=0.44, p=0.6587	0.06 (strong $H_0$ )	-1.8 (compelling $H_0$ )
Maier Unsuccessful 3	t(639)=-1.3, p=0.1943	0.1 (strong $H_0$ )	-0.3 (inconclusive)
Maier Study 4	t(326)=-1.81, p=0.0717	0.3 (strong $H_0$ )	-0.2 (inconclusive)
Aspirin study*	<b>t(21998)=5.19, p&lt;10<sup>-6</sup></b>	<b>10561.9 (decisive <math>H_1</math>)</b>	<b>1.3 (compelling <math>H_1</math>)</b>

586

587 **Table 1.** Reanalysis of data purportedly showing precognitive abilities [36] and a simulated example of a  
588 clinical trial. For each experiment this shows the result using classical statistics, a default Bayes factor  
589 [20] and the bootstrapped evidence  $\epsilon$ . For Bayes factors and bootstrapped evidence verbal descriptions  
590 of the strength of evidence for  $H_1$  and  $H_0$  are also given. Cells with bold font denote tests formally  
591 supporting  $H_1$ , that is, if evidence for the alternative hypothesis is above criterion (i.e.  $p<0.05$ ,  $BF_{10}>3$  or  
592  $\epsilon>0.5$ ). All these statistics are based on my own analysis of the raw data. The asterisks indicate that the  
593 data for these studies were simulated based on the reported effect and sample sizes because the raw  
594 data were not available to me.

595

596 Is the BSE test simply too conservative to reveal these rather subtle effects? To test whether the  
597 low evidence in my reanalysis could be due to a lack in sensitivity, I also used the BSE test to  
598 evaluate a small effect size in the clinical literature. In this study, the effect of aspirin on cardiac  
599 health was tested in a large sample ( $n\approx 22,000$ ) [37]. The effect was minute (Cohen's  $d\approx 0.07$ ) but  
600 highly significant, and it was thus deemed to be of practical value in promoting health. In the  
601 absence of raw data I simulated this data set using the reported effect and sample sizes. Here the  
602 BSE test agrees with classical statistics and Bayesian inference: the observed data strongly  
603 support ( $\epsilon=1.3$ ) the efficacy of the drug (Table 1). The sample size in this case is more than  
604 sufficient to conclude that this effect is *small but real*. In contrast, to provide conclusive evidence  
605 for the tiny precognition effects, the sample sizes would have to be orders of magnitude larger  
606 (Fig. 5D).

607

## 608 Discussion

609

610 Here I outlined the bootstrapped evidence test, which makes minimal assumptions and is easily  
611 applicable to numerous situations. It quantifies non-dichotomously the evidence *for the*  
612 *alternative or the null hypothesis* rather than whether a particular statistic passes an arbitrary  
613 significance threshold. A result does not stand or fall based on its exact value. Rather it allows us  
614 to express how convincing a result is. This also allows for the use of sequential sampling  
615 strategies, which can greatly benefit research practice especially when there is uncertainty about  
616 the expected effect size.

617

618 A possible criticism of non-dichotomous alternatives to null hypothesis significance testing is that  
619 any measures could be subject to the same thresholding dilemma as classical p-values. If the  
620 consensus emerges that  $\epsilon>0.5$  is sufficient evidence for  $H_1$  then are we not merely shifting the

621 problem from p-values to a different measure? However, even in the classical framework many  
622 researchers regularly make non-dichotomous judgments based on p-values. “Marginally  
623 significant” findings are often reported that do not quite pass reach  $p < 0.05$ . Many researchers are  
624 probably more convinced by  $p = 10^{-17}$  than  $p = 0.049$ . The bootstrapped evidence test directly  
625 quantifies the reliability of the data in drawing conclusions about the two competing hypotheses.  
626 Being a new measure it will make it easier for researchers to adopt non-dichotomous thinking  
627 than with p-values.

628

629 It remains unclear in how far researchers need dichotomous thresholds for statistical inference  
630 and whether labels for the strength of evidence, such as those employed for Bayes factors  
631 [26,27], are necessary. I would argue that they are not and deliberately refrained from proposing  
632 a categorical stratification of  $\epsilon$ . For replication attempts or incremental experiments for which  
633 clear predictions can be made,  $\epsilon$  between 0.7-1 can already be very convincing. In contrast, for  
634 more extraordinary claims even  $\epsilon = 1$  is still low.

635

636 Another problem with most commonly used statistics is that they are based on parametric  
637 assumptions that may not hold. Anomalous data can skew the default Bayes factor to a similar  
638 extent as inferences based on classical p-values, because it is based on the same effect size  
639 calculations. Robustness to outliers and heteroscedasticity could be incorporated into Bayesian  
640 hypothesis testing, e.g. by outlier removal procedures. Within the framework of classical  
641 inference robust hypothesis testing suggest following a complicated tree of tests for the presence  
642 of outliers and heteroscedasticity, and then applying the appropriate robust test depending on  
643 the situation [22,23]. This is often accompanied by remonstrations that there is a “statistical  
644 toolbox” that should be employed and that one method is not best for every situation [9].

645

646 While it is doubtless true that inference should never be made without thought [10], it is  
647 nevertheless advisable to seek a method that serves a universal purpose because in practice  
648 many researchers are not statisticians. The idea of a “statistical toolbox” is fraught with danger. It  
649 must inevitably result in increasing the underreported flexibility in the range of methods  
650 employed by published studies [28,38].

651

652 Of course the BSE test does not preclude the use of other statistical methods. In particular, I view  
653 it as a *complement to Bayesian inference* rather than an alternative. It provides a common  
654 starting point for inference. It is an objective method with minimal assumptions that in principle  
655 works exactly the same for almost any situation. The BSE test is particularly useful when there is  
656 substantial uncertainty about the expected effect size, when there are violations of parametric  
657 assumptions, and generally whenever application of Bayesian inference is difficult. It is also  
658 especially useful for exploratory analyses. The bootstrapped distribution of the effect size  
659 estimate could inform a prior used for Bayesian analysis of subsequent replication attempts.

660

661 The BSE framework also directly encourages sharing of raw data. While it can be applied to data  
662 simulated based on parameter estimates, as I have done for the Aspirin study in Table 1, this  
663 neglects additional information about potential data anomalies (Fig. 6). For the purpose of meta-  
664 analysis or even simply reanalyzing individual research findings, the full advantages of the BSE  
665 test become apparent when it is used on raw data.

666

667 Finally, it is important to remember that no statistical procedure can replace scientific scrutiny.  
668 Statistics do not confirm or refute theories. No amount of statistical evidence can prove whether  
669 particular phenomena, be it social priming, brain-behavior correlations or even far-fetched claims  
670 like precognition, exist. They can only provide support that the predictions a particular hypothesis  
671 makes are likely not to have occurred by chance or another, more trivial explanation.

672

### 673 **Acknowledgements**

674

675 I thank Ged Ridgway, Benjamin de Haas, and Micah Allen for comments on previous versions of  
676 this manuscript.

677

### 678 **References**

679

680 1. Cohen J. The Earth is round ( $p < .05$ ). *Am Psychol.* 1994;49: 997–1003.

681 2. Wagenmakers E-J. A practical solution to the pervasive problems of p values. *Psychon Bull*  
682 *Rev.* 2007;14: 779–804.

683 3. Wagenmakers E-J, Wetzels R, Borsboom D, van der Maas HLJ. Why psychologists must  
684 change the way they analyze their data: the case of psi: comment on Bem (2011). *J Pers Soc*  
685 *Psychol.* 2011;100: 426–432. doi:10.1037/a0022790

686 4. Masicampo EJ, Lalande DR. A peculiar prevalence of p values just below .05. *Q J Exp Psychol*  
687 2006. 2012;65: 2271–2279. doi:10.1080/17470218.2012.711335

688 5. Cumming G. The new statistics: why and how. *Psychol Sci.* 2014;25: 7–29.  
689 doi:10.1177/0956797613504966

690 6. Colquhoun D. An investigation of the false discovery rate and the misinterpretation of p-  
691 values. *R Soc Open Sci.* 2014;1: 140216. doi:10.1098/rsos.140216

692 7. Halsey LG, Curran-Everett D, Vowler SL, Drummond GB. The fickle P value generates  
693 irreproducible results. *Nat Methods.* 2015;12: 179–185. doi:10.1038/nmeth.3288

694 8. Trafimow D, Marks M. Editorial. *Basic Appl Soc Psychol.* 2015;37: 1–2.  
695 doi:10.1080/01973533.2015.1012991

696 9. Gigerenzer G, Marewski JN. Surrogate Science The Idol of a Universal Method for Scientific  
697 Inference. *J Manag.* 2015;41: 421–440. doi:10.1177/0149206314547522

698 10. Gigerenzer G. Mindless statistics. *J Socio-Econ.* 2004;33: 587–606.  
699 doi:10.1016/j.socec.2004.09.033

700 11. Nuzzo R. Scientific method: statistical errors. *Nature.* 2014;506: 150–152.  
701 doi:10.1038/506150a

702 12. Psychological Science. 2014 Submission Guidelines - Association for Psychological Science  
703 [Internet]. 2014 [cited 2 Apr 2014]. Available:  
704 [http://www.psychologicalscience.org/index.php/publications/journals/psychological\\_science/ps-submissions](http://www.psychologicalscience.org/index.php/publications/journals/psychological_science/ps-submissions)  
705

- 706 13. Morey RD, Rouder JN, Verhagen J, Wagenmakers E-J. Why Hypothesis Tests Are Essential for  
707 Psychological Science: A Comment on Cumming (2014). *Psychol Sci.* 2014;  
708 doi:10.1177/0956797614525969
- 709 14. Hoekstra R, Morey RD, Rouder JN, Wagenmakers E-J. Robust misinterpretation of confidence  
710 intervals. *Psychon Bull Rev.* 2014; doi:10.3758/s13423-013-0572-3
- 711 15. Wilcox RR, Muska J. Inferences about correlations when there is heteroscedasticity. *Br J Math*  
712 *Stat Psychol.* 2001;54: 39–47.
- 713 16. Bem DJ. Feeling the future: experimental evidence for anomalous retroactive influences on  
714 cognition and affect. *J Pers Soc Psychol.* 2011;100: 407–425. doi:10.1037/a0021524
- 715 17. Bem DJ, Utts J, Johnson WO. Must psychologists change the way they analyze their data? *J*  
716 *Pers Soc Psychol.* 2011;101: 716–719. doi:10.1037/a0024777
- 717 18. Wagenmakers E-J, Wetzels R, Borsboom D, Kievit R, van der Maas HLJ. Yes, psychologists  
718 must change the way they analyze their data: Clarifications for Bem, Utts, & Johnson  
719 [Internet]. 2011. Available:  
720 <http://dl.dropbox.com/u/1018886/ClarificationsForBemUttsJohnson.pdf>
- 721 19. Savalei V, Dunn E. Is the call to abandon p-values the red herring of the replicability crisis?  
722 *Cognition.* 2015;6: 245. doi:10.3389/fpsyg.2015.00245
- 723 20. Rouder JN, Speckman PL, Sun D, Morey RD, Iverson G. Bayesian t tests for accepting and  
724 rejecting the null hypothesis. *Psychon Bull Rev.* 2009;16: 225–237.  
725 doi:10.3758/PBR.16.2.225
- 726 21. Dienes Z. Using Bayes to get the most out of non-significant results. *Quant Psychol Meas.*  
727 2014;5: 781. doi:10.3389/fpsyg.2014.00781
- 728 22. Wilcox RR. Introduction to robust estimation and hypothesis testing. Academic Press; 2005.
- 729 23. Pernet CR, Wilcox R, Rousselet GA. Robust correlation analyses: false positive and power  
730 validation using a new open source matlab toolbox. *Front Psychol.* 2012;3: 606.  
731 doi:10.3389/fpsyg.2012.00606
- 732 24. Rousselet GA, Pernet CR. Improving standards in brain-behaviour correlation analyses. *Front*  
733 *Hum Neurosci.* 2012;6: 119. doi:10.3389/fnhum.2012.00119
- 734 25. Schwarzkopf DS, De Haas B, Rees G. Better ways to improve standards in brain-behavior  
735 correlation analysis. *Front Hum Neurosci.* 2012;6: 200. doi:10.3389/fnhum.2012.00200
- 736 26. Wetzels R, Wagenmakers E-J. A default Bayesian hypothesis test for correlations and partial  
737 correlations. *Psychon Bull Rev.* 2012;19: 1057–1064. doi:10.3758/s13423-012-0295-x
- 738 27. Jeffreys H. Theory of probability. Oxford, UK: Oxford University Press; 1961.
- 739 28. Simmons JP, Nelson LD, Simonsohn U. False-positive psychology: undisclosed flexibility in  
740 data collection and analysis allows presenting anything as significant. *Psychol Sci.* 2011;22:  
741 1359–1366. doi:10.1177/0956797611417632

- 742 29. Rouder JN. Optional stopping: No problem for Bayesians. *Psychon Bull Rev.* 2014;  
743 doi:10.3758/s13423-014-0595-4
- 744 30. Ioannidis JPA. Why most published research findings are false. *PLoS Med.* 2005;2: e124.  
745 doi:10.1371/journal.pmed.0020124
- 746 31. Lakens D. What p-hacking really looks like [Internet]. 2015 [cited 4 Feb 2015]. Available:  
747 <https://osf.io/ycag9/>
- 748 32. Anscombe FJ. Graphs in Statistical Analysis. *Am Stat.* 1973;27: 17–21. doi:10.2307/2682899
- 749 33. Schwarzkopf DS, Song C, Rees G. The surface area of human V1 predicts the subjective  
750 experience of object size. *Nat Neurosci.* 2011;14: 28–30. doi:10.1038/nn.2706
- 751 34. Schwarzkopf DS, Rees G. Subjective size perception depends on central visual cortical  
752 magnification in human v1. *PloS One.* 2013;8: e60550. doi:10.1371/journal.pone.0060550
- 753 35. Genç E, Bergmann J, Singer W, Kohler A. Surface Area of Early Visual Cortex Predicts  
754 Individual Speed of Traveling Waves During Binocular Rivalry. *Cereb Cortex N Y N* 1991.  
755 2014; doi:10.1093/cercor/bht342
- 756 36. Maier MA, Buechner VL, Kuhbandner C, Pflitsch M, Fernandez-Capo M, Gamiz-Sanfeliu M.  
757 Feeling the Future Again: Retroactive Avoidance of Negative Stimuli. *J Conscious Stud.*  
758 2014;21: 121–152.
- 759 37. Young F, Nightingale S, Temple R. The preliminary report of the findings of the aspirin  
760 component of the ongoing physicians' health study: The fda perspective on aspirin for the  
761 primary prevention of myocardial infarction. *JAMA.* 1988;259: 3158–3160.  
762 doi:10.1001/jama.1988.03720210048028
- 763 38. Wagenmakers E-J, Wetzels R, Borsboom D, Maas HLJ van der, Kievit RA. An Agenda for Purely  
764 Confirmatory Research. *Perspect Psychol Sci.* 2012;7: 632–638.  
765 doi:10.1177/1745691612463078
- 766