

¹ Indices of Effect Existence and Significance in the Bayesian Framework

² Dominique Makowski¹, Mattan S. Ben-Shachar², Daniel Lüdecke³, & S.H. Annabel Chen^{1,4}

³ ¹ Nanyang Technological University, Singapore

⁴ ² Ben-Gurion University of the Negev, Israel

⁵ ³ University Medical Center Hamburg-Eppendorf, Germany

⁶ ⁴ Centre for Research and Development in Learning (CRADLE), Singapore

⁷ Author Note

⁸ Daniel Lüdecke and S.H. Annabel Chen share senior authorship.

⁹ Correspondence concerning this article should be addressed to Dominique Makowski,
¹⁰ HSS 04-18, 48 Nanyang Avenue, Singapore. E-mail: dmakowski@ntu.edu.sg

11

Abstract

12 Turmoil has engulfed psychological science. Causes and consequences of the reproducibility
13 crisis are in dispute. With the hope of addressing some of its aspects, Bayesian methods are
14 gaining increasing attention in psychological science. Some of their advantages, as opposed
15 to the frequentist framework, are the ability to describe parameters in probabilistic terms
16 and explicitly incorporate prior knowledge about them into the model. These issues are
17 crucial in particular regarding the current debate about statistical significance. Bayesian
18 methods are not necessarily the only remedy against incorrect interpretations or wrong
19 conclusions, but there is an increasing agreement that they are one of the keys to avoid such
20 fallacies. Nevertheless, its flexible nature is its power and weakness, for there is no agreement
21 about what indices of “significance” should be computed or reported. This lack of a
22 consensual index or guidelines, such as the frequentist p -value, further contributes to the
23 unnecessary opacity that many non-familiar readers perceive in Bayesian statistics. Thus,
24 this study describes and compares several Bayesian indices, provide intuitive visual
25 representation of their “behavior” in relationship with common sources of variance such as
26 sample size, magnitude of effects and also frequentist significance. The results contribute to
27 the development of an intuitive understanding of the values that researchers report, allowing
28 to draw sensible recommendations for Bayesian statistics description, critical for the
29 standardization of scientific reporting.

30

Keywords: Bayesian, significance, NHST, p-value, Bayes factors

31

Word count: 6293

32 Indices of Effect Existence and Significance in the Bayesian Framework

33 **Introduction**

34 The Bayesian framework is quickly gaining popularity among psychologists and
35 neuroscientists (Andrews & Baguley, 2013). Reasons to prefer this approach are reliability,
36 better accuracy in noisy data, better estimation for small samples, less proneness to type I
37 errors, the possibility of introducing prior knowledge into the analysis and the intuitiveness
38 and straightforward interpretation of results (Dienes & Mcclatchie, 2018; Etz &
39 Vandekerckhove, 2016; Kruschke, 2010; Kruschke, Aguinis, & Joo, 2012; Wagenmakers et al.,
40 2018; Wagenmakers, Morey, & Lee, 2016). On the other hand, the frequentist approach has
41 been associated with the focus on *p*-values and null hypothesis significance testing (NHST).
42 The misinterpretation and misuse of *p*-values, so called “p-hacking” (Simmons, Nelson, &
43 Simonsohn, 2011), has been shown to critically contribute to the reproducibility crisis in
44 psychological science (Chambers, Feredoes, Muthukumaraswamy, & Etchells, 2014; Szucs &
45 Ioannidis, 2016). Not only are *p*-values used to draw inappropriate inferences from noisy
46 data, but even when used properly, effects are drastically overestimated, sometimes even in
47 the wrong direction, when estimation is tied to statistical significance in highly variable data
48 (Gelman, 2018). In response, there is a general agreement that the generalization and
49 utilization of the Bayesian framework is one way of overcoming these issues (Benjamin et al.,
50 2018; Etz & Vandekerckhove, 2016; Halsey, 2019; Marasini, Quatto, & Ripamonti, 2016;
51 Maxwell, Lau, & Howard, 2015; Wagenmakers et al., 2017).

52 The tenacity and resilience of the *p*-value as an index of significance is remarkable,
53 despite the long-lasting criticism and discussion about its misuse and misinterpretation
54 (Anderson, Burnham, & Thompson, 2000; Cohen, 2016; Fidler, Thomason, Cumming, Finch,
55 & Leeman, 2004; Finch et al., 2004; Gardner & Altman, 1986). This endurance might be
56 informative on how such indices, and the accompanying heuristics applied to interpret them
57 (e.g., assigning thresholds like .05, .01 and .001 to certain levels of significance), are useful

58 and necessary for researchers to gain an intuitive (although possibly simplified)
59 understanding of the interactions and structure of their data. Moreover, the utility of such
60 an index is most salient in contexts where decisions must be made and rationalized (e.g., in
61 medical settings). Unfortunately, these heuristics can become severely rigidified, and meeting
62 significance has become a goal unto itself rather than a tool for understanding the data
63 (Cohen, 2016; Kirk, 1996). This is particularly problematic given that *p*-values can only be
64 used to reject the null hypothesis and not to accept it as true, because a statistically
65 non-significant result does not mean that there is no difference between groups or no effect of
66 a treatment (Amrhein, Greenland, & McShane, 2019; Wagenmakers, 2007).

67 While significance testing (and its inherent categorical interpretation heuristics) might
68 have its place as a complementary perspective to effect estimation, it does not preclude the
69 fact that drastic improvements are needed. For instance, one possible advance could focus on
70 improving the mathematical understanding (e.g., through a new simpler index) of the values
71 being used (as opposed to the obscure mathematical definition of the *p*-value, contributing to
72 its common misinterpretation). Another improvement could be found in providing an
73 intuitive understanding (e.g., by visual means) of the behavior of the indices in relationship
74 with main sources of variance, such as sample size, noise or effect presence. Such better
75 overall understanding of the indices would hopefully act as a barrier against their mindless
76 reporting by allowing the users to nuance the interpretations and conclusions that they draw.

77 The Bayesian framework offers several alternative indices for the *p*-value. To better
78 understand these indices, it is important to point out one of the core differences between
79 Bayesian and frequentist methods. From a frequentist perspective, the effects are fixed (but
80 unknown) and data are random. On the other hand, instead of having single estimates of
81 some “true effect” (for instance, the “true” correlation between *x* and *y*), Bayesian methods
82 compute the probability of different effects values *given* the observed data (and some prior
83 expectation), resulting in a distribution of possible values for the parameters, called the

84 posterior distribution. The description of the posterior distribution (e.g., through its
85 centrality, dispersion, etc.) allows to draw conclusions from Bayesian analyses.

86 Bayesian “significance” testing indices could be roughly grouped into three overlapping
87 categories: Bayes factors, posterior indices and Region of Practical Equivalence
88 (ROPE)-based indices. Bayes factors are a family of indices of relative evidence of one model
89 over another (e.g., the null *vs.* the alternative hypothesis; Jeffreys, 1998; Ly, Verhagen, &
90 Wagenmakers, 2016). They provide many advantages over the *p*-value by having a
91 straightforward interpretation as well as allowing to quantify evidence in favor of the null
92 hypothesis (Dienes, 2014; Jarosz & Wiley, 2014). However, its use for parameters description
93 in complex models is still a matter of debate (Heck, 2019; Wagenmakers, Lodewyckx,
94 Kuriyal, & Grasman, 2010), being highly dependent on the specification of priors (Etz, Haaf,
95 Rouder, & Vandekerckhove, 2018; Kruschke & Liddell, 2018). On the contrary, “posterior
96 indices” reflect objective characteristics of the posterior distribution, for instance the
97 proportion of strictly positive values. While the simplicity of their computation and
98 interpretation is an asset, it might also limit the information that they provide. Finally,
99 ROPE-based indices are related to the redefinition of the null hypothesis from the classic
100 point-null hypothesis to a range of values considered negligible or too small to be of any
101 practical relevance (the Region of Practical Equivalence - ROPE; Kruschke, 2014; Lakens,
102 2017; Lakens, Scheel, & Isager, 2018), usually spread equally around 0 (e.g., [-0.1; 0.1]). It is
103 interesting to note that this perspective unites significance testing with the focus on effect
104 size (involving a discrete separation between at least two categories: negligible and
105 non-negligible), which finds an echo in recent statistical recommendations (Ellis & Steyn,
106 2003; Simonsohn, Nelson, & Simmons, 2014; Sullivan & Feinn, 2012).

107 Despite the richness provided by the Bayesian framework and the availability of
108 multiple indices, no consensus has yet emerged on which ones to be used. Literature
109 continues to bloom in a raging debate, often polarized between proponents of the Bayes

factor as the supreme index and its detractors (Robert, 2014, 2016; Spanos, 2013; Wagenmakers, Lee, Rouder, & Morey, 2019), with strong theoretical arguments being developed on both sides. Yet no practical, empirical and direct comparison between these indices has been done. This might be a deterrent for scientists interested in adopting the Bayesian framework. Moreover, this grey area can increase the difficulty of readers or reviewers unfamiliar with the Bayesian framework to follow the assumptions and conclusions, which could in turn generate unnecessary doubt upon an entire study. While we think that such indices of significance and their interpretation guidelines (in the form of rules of thumb) are useful in practice, we also strongly believe that they should be accompanied with the understanding of their “behavior” in relationship with major sources of variance, such as sample size, noise or effect presence. This knowledge is important for people to implicitly and intuitively appraise the meaning and implication of the mathematical values they report. Such an understanding could prevent the crystallization of the possible heuristics and categories derived from such indices, as has unfortunately occurred for the p -values.

Thus, based on the simulation of linear and logistic regressions (arguably some of the most widely used models in the psychological sciences), the present work aims at comparing several indices of effect “significance”, provide visual representations of the “behavior” of such indices in relationship with sample size, noise and effect presence, as well as their relationship to frequentist p -values (an index which, beyond its many flaws, is well known and could be used as a reference for Bayesian neophytes), and finally draw recommendations for Bayesian statistics reporting.

Methods

Data Simulation

We simulated datasets suited for linear and logistic regression and started by simulating an independent, normally distributed x variable (with mean 0 and SD 1) of a given sample size. Then, the corresponding y variable was added, having a perfect correlation

136 (in the case of data for linear regressions) or as a binary variable perfectly separated by x .
137 The case of no effect was simulated by creating a y variable that was independent of (i.e. not
138 correlated to) x . Finally, a Gaussian noise was added to the x variable (the error).

139 The simulation aimed at modulating the following characteristics: *outcome type* (linear
140 or logistic regression), *sample size* (from 20 to 100 by steps of 10), *null hypothesis* (original
141 regression coefficient from which data is drawn prior to noise addition, 1 - presence of “true”
142 effect, or 0 - absence of “true” effect) and *noise* (Gaussian noise applied to the predictor with
143 SD uniformly spread between 0.666 and 6.66, with 1000 different values), which is directly
144 related to the absolute value of the coefficient (i.e., the effect size). We generated a dataset
145 for each combination of these characteristics, resulting in a total of 36,000 (2 model types * 2
146 presence/absence of effect * 9 sample sizes * 1,000 noise variations) datasets. The code used
147 for data generation is available on GitHub ([https://github.com/easystats/easystats/tree/
148 master/publications/makowski_2019_bayesian/data](https://github.com/easystats/easystats/tree/master/publications/makowski_2019_bayesian/data)). Note that it takes usually several
149 days/weeks for the generation to complete.

150 Indices

151 For each of these datasets, Bayesian and frequentist regressions were fitted to predict y
152 from x as a single unique predictor. We then computed the following seven indices from all
153 simulated models (see **Figure 1**), related to the effect of x .

154 **Frequentist *p*-value.** This was the only index computed by the frequentist version
155 of the regression. The *p*-value represents the probability that for a given statistical model,
156 when the null hypothesis is true, the effect would be greater than or equal to the observed
157 coefficient (Wasserstein, Lazar, & others, 2016).

158 **Probability of Direction (*pd*).** The *Probability of Direction* (*pd*) varies between
159 50% and 100% and can be interpreted as the probability that a parameter (described by its
160 posterior distribution) is strictly positive or negative (whichever is the most probable). It is
161 mathematically defined as the proportion of the posterior distribution that is of the median’s

162 sign (Makowski, Ben-Shachar, & Lüdecke, 2019).

163 **MAP-based *p*-value.** The *MAP-based p-value* is related to the odds that a
164 parameter has against the null hypothesis (Mills, 2017; Mills & Parent, 2014). It is
165 mathematically defined as the density value at 0 divided by the density at the Maximum A
166 Posteriori (MAP), *i.e.*, the equivalent of the mode for continuous distributions.

167 **ROPE (95%).** The *ROPE (95%)* refers to the percentage of the 95% HDI that lies
168 within the ROPE. As suggested by Kruschke (2014), the Region of Practical Equivalence
169 (ROPE) was defined as range from -0.1 to 0.1 for linear regressions and its equivalent, -0.18
170 to 0.18, for logistic models (based on the $\pi/\sqrt{3}$ formula to convert log odds ratios to
171 standardized differences; Cohen, 1988).

172 **ROPE (full).** The *ROPE (full)* is similar to *ROPE (95%)*, with the exception that
173 it refers to the percentage of the *whole* posterior distribution that lies within the ROPE.

174 **Bayes factor (*vs.* 0).** The Bayes Factor (*BF*) used here is based on prior and
175 posterior distributions of a single parameter. In this context, the Bayes factor indicates the
176 degree by which the mass of the posterior distribution has shifted further away from or closer
177 to the null value (0), relative to the prior distribution, thus indicating if the null hypothesis
178 has become less or more likely given the observed data. The *BF* was computed as a
179 Savage-Dickey density ratio, which is also an approximation of a Bayes factor comparing the
180 marginal likelihoods of the model against a model in which the tested parameter has been
181 restricted to the point-null (Wagenmakers et al., 2010).

182 **Bayes factor (*vs.* ROPE).** The *Bayes factor (vs. ROPE)* is similar to the *Bayes*
183 *factor (vs. 0)*, but instead of a point-null, the null hypothesis is a range of negligible values
184 (defined here same as for the ROPE indices). The *BF* was computed by comparing the prior
185 and posterior odds of the parameter falling within vs. outside the ROPE (see
186 *Non-overlapping Hypotheses* in Morey & Rouder, 2011). This measure is closely related to
187 the *ROPE (full)*, as it can be formally defined as the ratio between the *ROPE (full)* odds for
188 the posterior distribution and the *ROPE (full)* odds for the prior distribution:

$$BF_{rope} = \frac{odds(ROPE_{\text{full posterior}})}{odds(ROPE_{\text{full prior}})}$$

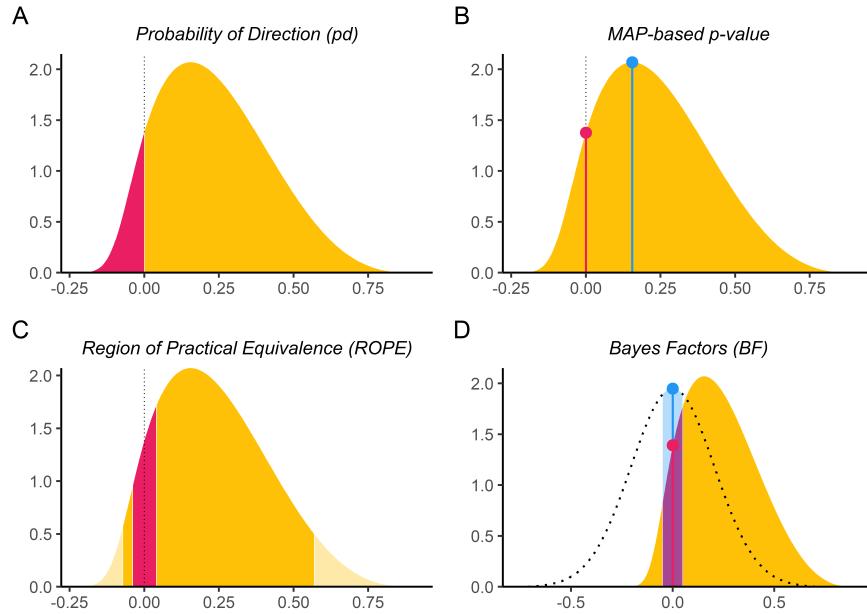


Figure 1. Bayesian indices of effect existence and significance. (A) The Probability of Direction (*pd*) is defined as the proportion of the posterior distribution that is of the median's sign (the size of the yellow area relative to the whole distribution). (B) The MAP-based *p*-value is defined as the density value at 0, - the height of the red lollipop, divided by the density at the Maximum A Posteriori (MAP), - the height of the blue lollipop. (C) The percentage in ROPE corresponds to the red area relative to the distribution (with or without tails for ROPE (*full*) and ROPE (*95%*), respectively). (D) The Bayes factor (vs. 0) corresponds to the point-null density of the prior (the blue lollipop on the dotted distribution) divided by that of the posterior (the red lollipop on the yellow distribution), and the Bayes factor (vs. ROPE) is calculated as the odds of the prior falling within vs. outside the ROPE (the blue area on the dotted distribution) divided by that of the posterior (the red area on the yellow distribution).

189 Data Analysis

190 In order to achieve the two-fold aim of this study; 1) comparing Bayesian indices and
191 2) provide visual guides for an intuitive understanding of the numeric values in relation to a
192 known frame of reference (the frequentist *p*-value), we will start by 1) presenting the
193 relationship between these indices and main sources of variance, such as sample size, noise
194 and null hypothesis (true if absence of effect, false if presence of effect). We will then 2)
195 compare Bayesian indices with the frequentist *p*-value and its commonly used thresholds (.05,
196 .01, .001). Finally, we will show the mutual relationship between 3 recommended Bayesian
197 candidates. Taken together, these results will help us outline guides to ease the reporting
198 and interpretation of the indices.

199 In order to provide an intuitive understanding of values, data processing will focus on
200 creating clear visual figures to help the user grasp the patterns and variability that exists
201 when computing the investigated indices. Nevertheless, we decided to also mathematically
202 test our claims in cases where the graphical representation begged for a deeper investigation.
203 Thus, we fitted two regression models to assess the impact of sample size and noise,
204 respectively. To ensure that any differences between the indices are not due to differences in
205 their scale or distribution, we converted all indices to the same scale by normalizing the
206 indices between 0 and 1 (note that *BFs* were transformed to posterior probabilities,
207 assuming uniform prior odds) and reversing the *p*-values, the MAP-based *p*-values and the
208 ROPE indices so that a higher value corresponds to stronger “significance”.

209 The statistical analyses were conducted using R (R Core Team, 2019). Computations of
210 Bayesian models were done using the *rstanarm* package (Goodrich, Gabry, Ali, & Brilleman,
211 2019), a wrapper for Stan probabilistic language (Carpenter et al., 2017). We used Markov
212 Chain Monte Carlo sampling (in particular, Hamiltonian Monte Carlo; Gelman et al., 2014)
213 with 4 chains of 2000 iterations, half of which used for warm-up. Mildly informative priors (a
214 normal distribution with mean 0 and SD 1) were used for the parameter in all models. The

²¹⁵ Bayesian indices were calculated using the *bayestestR* package (Makowski et al., 2019).

²¹⁶

Results

²¹⁷ Impact of Sample Size

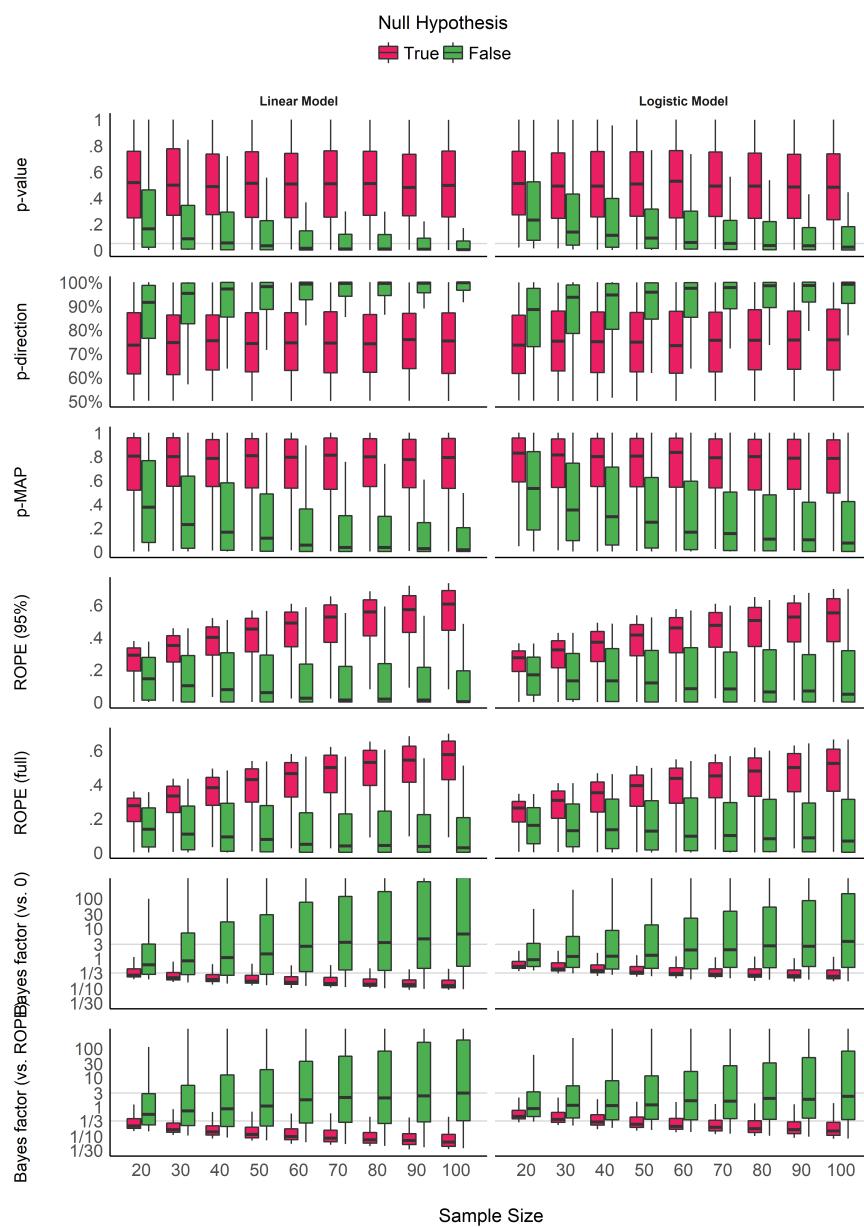


Figure 2. Impact of Sample Size on the different indices, for linear and logistic models, and when the null hypothesis is true or false. Grey vertical lines for *p*-values and Bayes factors represent commonly used thresholds.

Table 1

Sensitivity to sample size. This table shows the standardized coefficient between the sample size and the value of each index, adjusted for error, and stratified by model type and presence of true effect. The stronger the coefficient is, the stronger the relationship with sample size.

Index	Linear Models /	Linear Models /	Logistic Models /	Logistic Models /
	Presence of Effect	Absence of Effect	Presence of Effect	Absence of Effect
p-value	0.17	0.01	0.16	0.02
p-direction	0.17	0.01	0.15	0.02
p-MAP	0.24	0.00	0.24	0.03
ROPE (95%)	0.03	0.36	0.01	0.31
ROPE (full)	0.03	0.36	0.02	0.31
Bayes factor (vs. 0)	0.20	0.12	0.12	0.14
Bayes factor (vs. ROPE)	0.15	0.14	0.08	0.18

218 **Figure 2** shows the sensitivity to sample size of the indices. The *p*-value, the *pd* and
 219 the MAP-based *p*-value are sensitive to sample size only in case of the presence of a true
 220 effect (when the null hypothesis is false). When the null hypothesis is true, all three indices
 221 are unaffected by sample size. In other words, these indices reflect the amount of observed
 222 evidence (the sample size) for the presence of an effect (i.e., against the null hypothesis being
 223 true), but not for the absence of an effect. The *ROPE* indices, however, appear as strongly
 224 modulated by the sample size when there is no effect, suggesting their sensitivity to the
 225 amount of evidence for the absence of effect. Finally, the figure suggests that *BFs* are
 226 sensitive to sample size for both presence and absence of true effect.

227 Consistently with **Figure 2**, the model investigating the sensitivity of sample size on
 228 the different indices suggests that *BF* indices are sensitive to sample size both when an effect
 229 is present (null hypothesis is false) and absent (null hypothesis is true). *ROPE* indices are
 230 particularly sensitive to sample size when the null hypothesis is true, while *p*-value, *pd* and
 231 MAP-based *p*-value are only sensitive to sample size when the null hypothesis is false, in

which case they are more sensitive than *ROPE* indices. These findings can be related to the concept of consistency: as the number of data points increases, the statistic converges toward some “true” value. Here, we observe that *p*-value, *pd* and the MAP-based *p*-value are consistent only when the null hypothesis is false. In other words, as sample size increases, they tend to reflect more strongly that the effect is present. On the other hand, *ROPE* indices appear as consistent when the effect is absent. Finally, *BFs* are consistent both when the effect is absent and when it is present. Note also that *BF* (*vs.* *ROPE*), compared to *BF* (*vs.* 0), is more sensitive to sample size when the null hypothesis is true, and *ROPE* (*full*) is overall slightly more consistent than *ROPE* (95%).

Impact of Noise

Figure 3 shows the indices’ sensitivity to noise. Unlike the patterns of sensitivity to sample size, the indices display more similar patterns in their sensitivity to noise (or magnitude of effect). All indices are unidirectional impacted by noise: as noise increases, the observed coefficients decrease in magnitude, and the indices become less “pronounced” (respectively to their direction). However, it is interesting to note that the variability of the indices seems differently impacted by noise. For the *p*-values, the *pd* and the *ROPE* indices, the variability increases as the noise increases. In other words, small variation in small observed coefficients can yield very different values. On the contrary, the variability of *BFs* decreases as the true effect tends toward 0. For the MAP-based *p*-value, the variability appears to be the highest for moderate amount of noise. This behavior seems consistent across model types.

Consistently with **Figure 3**, the model investigating the sensitivity of noise when an effect is present (as there is only noise in the absence of effect), adjusted for sample size, suggests that *BFs* (especially *vs.* *ROPE*), followed by the MAP-based *p*-value and percentages in *ROPE*, are the most sensitive to noise. As noise is a proxy of effect size (linearly related to the absolute value of the coefficient of the parameter), this result

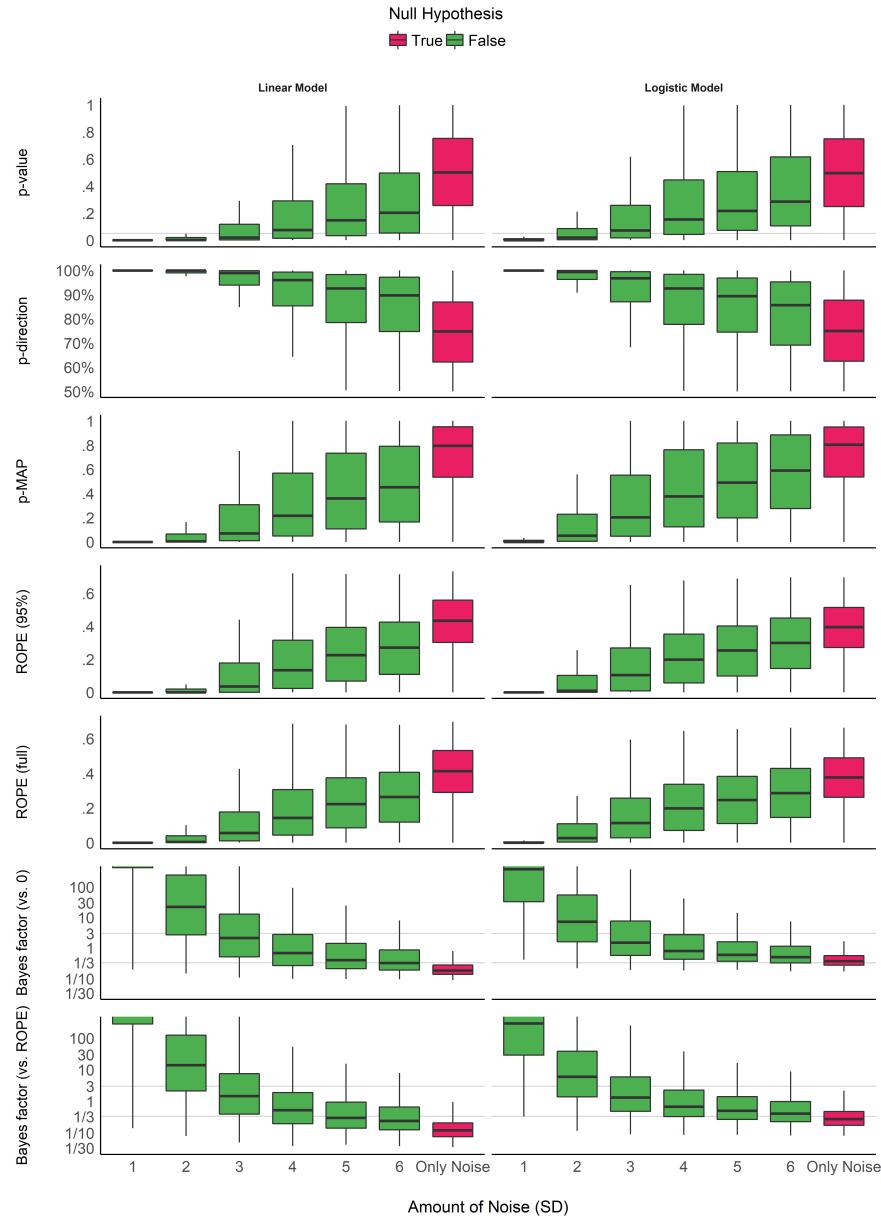


Figure 3. Impact of Noise. The noise corresponds to the standard deviation of the Gaussian noise that was added to the generated data. It is related to the magnitude the parameter (the more noise there is, the smaller the coefficient). Grey vertical lines for $*p*$ -values and Bayes factors represent commonly used thresholds. The scale is capped for the Bayes factors as these extend to infinity.

Table 2

Sensitivity to noise. This table shows the standardized coefficient between noise and the value of each index when the true effect is present, adjusted for sample size and stratified by model type. The stronger the coefficient is, the stronger the relationship with noise.

Index	Linear Models /	Logistic Models /
	Presence of Effect	Presence of Effect
p-value	0.35	0.40
p-direction	0.36	0.40
p-MAP	0.55	0.60
ROPE (95%)	0.45	0.45
ROPE (full)	0.46	0.45
Bayes factor (vs. 0)	0.79	0.65
Bayes factor (vs. ROPE)	0.81	0.67

258 highlights the fact that these indices are sensitive to the magnitude of the effect. For
 259 example, as noise increases, evidence for an effect becomes weak, and data seems to support
 260 the absence of an effect (or at the very least the presence of a negligible effect), which is
 261 reflected in *BFs* being consistently smaller than 1. On the other hand, as the *p*-value and
 262 the *pd* quantify evidence only for the presence of an effect, as noise increases, they are
 263 become more dependent on larger sample size to be able to detect the presence of an effect.

264 Relationship with the frequentist *p*-value

265 **Figure 4** suggests that the *pd* has a 1:1 correspondence with the frequentist *p*-value
 266 (through the formula $p_{two-sided} = 2 * (1 - p_d)$). *BF* indices still appear as having a severely
 267 non-linear relationship with the frequentist index, mostly due to the fact that smaller
 268 *p*-values correspond to stronger evidence in favor of the presence of an effect, but the reverse
 269 is not true. *ROPE*-based percentages appear to be only weakly related to *p*-values.
 270 Critically, their relationship seems to be strongly dependent on sample size.

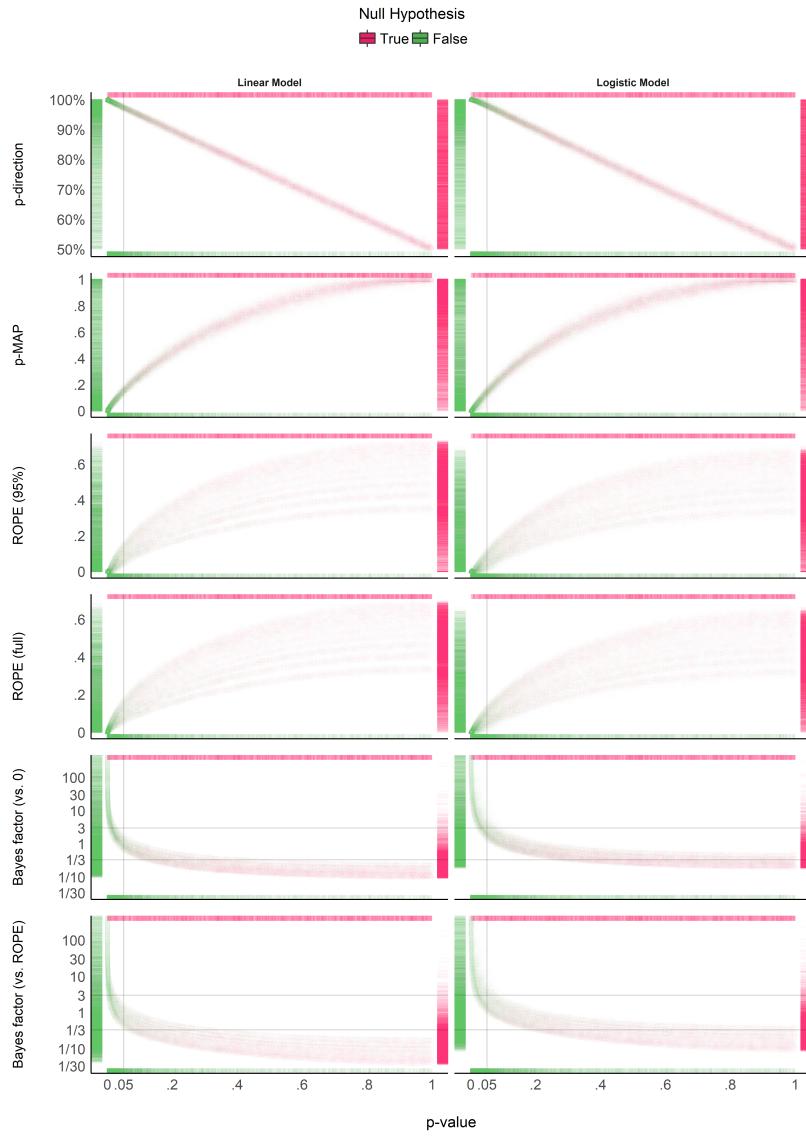


Figure 4. Relationship with the frequentist $*p*$ -value. In each plot, the $*p*$ -value densities are visualized by the marginal top (absence of true effect) and bottom (presence of true effect) markers, whereas on the left (presence of true effect) and right (absence of true effect), the markers represent the density of the index of interest. Different point shapes, representing different sample sizes, specifically illustrate its impact on the percentages in ROPE, for which each "curve line" is associated with one sample size (the bigger the sample size, the higher the percentage in ROPE).

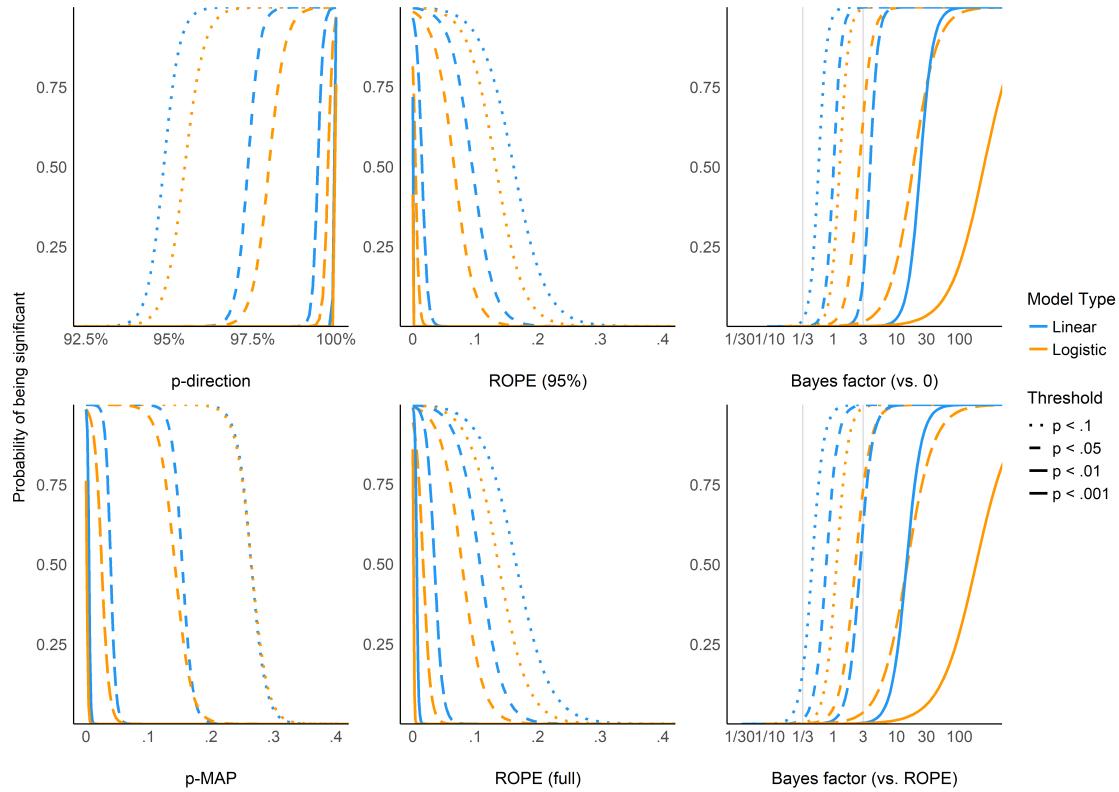


Figure 5. The probability of reaching different p -value based significance thresholds (.1, .05, .01, .001 for solid, long-dashed, short-dashed and dotted lines, respectively) for different values of the corresponding Bayesian indices.

271 **Figure 5** shows equivalence between p -value thresholds (.1, .05, .01, .001) and the

272 Bayesian indices. As expected, the p -direction has the sharpest thresholds (95%, 97.5%,
 273 99.5% and 99.95%, respectively). For logistic models, these threshold points appear as more
 274 conservative (i.e., Bayesian indices have to be more “pronounced” to reach the same level of
 275 significance). This sensitivity to model type is the strongest for BFs (which is possibly
 276 related to the difference in the prior specification for these two types of models).

277 **Relationship between ROPE (full), pd and BF (vs. ROPE)**

278 **Figure 6** suggests that the relationship between the $ROPE (full)$ and the pd might be

279 strongly affected by the sample size, and subject to differences across model types. This
 280 seems to echo the relationship between $ROPE (full)$ and p -value, the latter having a 1:1

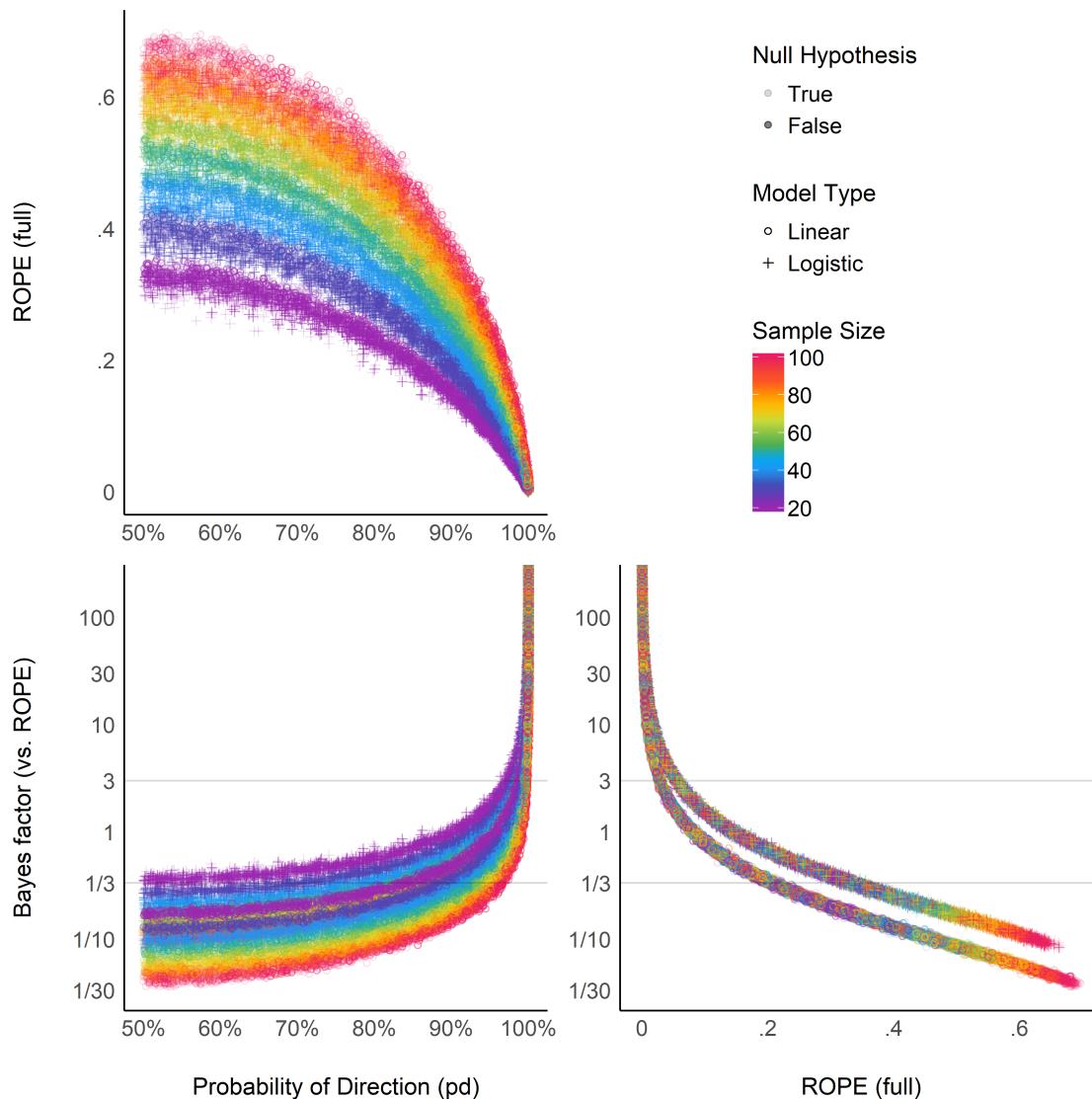


Figure 6. Relationship between three Bayesian indices: The Probability of Direction (*pd*), the percentage of the full posterior distribution in the ROPE, and the Bayes factor (*vs.* ROPE).

correspondence with pd . On the other hand, the $ROPE$ (*full*) and the BF (*vs.* $ROPE$) seem very closely related within the same model type, reflecting their formal relationship (see definition of BF (*vs.* $ROPE$) above). Overall, these results help to demonstrate $ROPE$ (*full*) and BF (*vs.* $ROPE$)'s consistency both in case of presence and absence of a true effect, whereas the pd , being equivalent to the p -value, is only consistent when the true effect is

286 absent.

287

Discussion

288 Based on the simulation of linear and logistic models, the present work aimed at
289 comparing several Bayesian indices of effect “significance” (see **Table 3**), providing visual
290 representations of the “behavior” of such indices in relationship with important sources of
291 variance such as sample size, noise and effect presence, as well as comparing them with the
292 well-known and widely used frequentist *p*-value and its arbitrary interpretation thresholds.

293 The results tend to suggest that the investigated indices could be separated into two
294 categories. The first group, including the *pd* and the MAP-based *p*-value, presents similar
295 properties to those of the frequentist *p*-value: they are sensitive to the amount of evidence
296 for the alternative hypothesis only (i.e., when an effect is truly present). In other words,
297 these indices are not able to reflect the amount of evidence in favor of the null hypothesis
298 (Rouder & Morey, 2012; Rouder, Speckman, Sun, Morey, & Iverson, 2009). A high value
299 suggests that the effect exists, but a low value indicates *uncertainty* about its existence (but
300 not certainty that it is non-existent). The second group, including ROPE and Bayes factors,
301 seem sensitive to both presence and absence of effect, accumulating evidence as the sample
302 size increases. However, the ROPE seems particularly suited to provide evidence in favor of
303 the null hypothesis. Consistently with this, combining Bayes factors with the ROPE (BF *vs.*
304 ROPE), as compared to Bayes factors against the point-null (BF *vs.* 0), leads to a higher
305 sensitivity to null-effects (Morey & Rouder, 2011; Rouder & Morey, 2012).

306 We also showed that besides sharing similar properties, the *pd* has a 1:1
307 correspondence with the frequentist *p*-value, being its Bayesian equivalent. On the contrary
308 Bayes factors appear as having a severely non-linear relationship with the frequentist index,
309 which is to be expected from their mathematical definition and their sensitivity when the
310 null hypothesis is true. This in turn can lead to surprising conclusions. For instance, Bayes

311 factors lower than 1, which are considered as providing evidence *against* the presence of an
312 effect, can still correspond to a “significant” frequentist *p*-value (see **Figures 3 and 4**).
313 ROPE indices are more closely related to the *p*-value, as their relationship appears
314 dependent on another factor, the sample size. This suggests that the ROPE encapsulates
315 additional information about the strength of evidence.

316 What is the point of comparing Bayesian indices with the frequentist *p*-value,
317 especially after having pointed out to its many flaws? While this comparison may seem
318 counter-intuitive (as Bayesian thinking is intrinsically different from the frequentist
319 framework), we believe that this juxtaposition is interesting for didactic reasons. The
320 frequentist *p*-value “speaks” to many and can thus be seen as a reference and a way to
321 facilitate the shift toward the Bayesian framework. Thus, pragmatically documenting such
322 bridges can only foster the understanding of the methodological issues that our field is facing,
323 and in turn act against dogmatic adherence to a framework. This does not preclude,
324 however, that a change in the general paradigm of significance seeking and “p-hacking” is
325 necessary, and that Bayesian indices are fundamentally different from the frequentist *p*-value,
326 rather than mere approximations or equivalents.

327 Critically, while the purpose of these indices was solely referred to as *significance* until
328 now, we would like to emphasize the nuanced perspective of the existence-significance testing
329 as a dual-framework for parameters description and interpretation. The idea supported here
330 is that there is a conceptual and practical distinction, and possible dissociation to be made,
331 between an effect’s existence *and* significance. In this context, *existence* is simply defined as
332 the consistency of an effect in one particular direction (i.e., positive or negative), without
333 any assumptions or conclusions as to its size, importance, relevance or meaning. It is an
334 objective feature of an estimate (tied to its uncertainty). On the other hand, *significance*
335 would be here re-framed following its original literally definition such as “being worthy of
336 attention” or “importance”. An effect can be considered significant if its magnitude is higher

Table 3

Summary of Bayesian Indices of Effect Existence and Significance.

Index	Interpretation	Definition	Strengths	Limitations
Probability of Direction (pd)	Probability that an effect is of the same sign as the median's.	Proportion of the posterior distribution of the same sign than the median's.	Straightforward computation and interpretation. Objective property of the posterior distribution. 1:1 correspondence with the frequentist p-value.	Limited information favoring the null hypothesis.
MAP-based p-value	Relative odds of the presence of an effect against 0.	Density value at 0 divided by the density value at the mode of the posterior distribution.	Straightforward computation. Objective property of the posterior distribution	Limited information favoring the null hypothesis. Relates on density approximation. Indirect relationship between mathematical definition and interpretation.
ROPE (95%)	Probability that the credible effect values are not negligible.	Proportion of the 95% CI inside of a range of values defined as the ROPE.	Provides information related to the practical relevance of the effects.	A ROPE range needs to be arbitrarily defined. Sensitive to the scale (the unit) of the predictors. Not sensitive to highly significant effects.
ROPE (full)	Probability that the effect possible values are not negligible.	Proportion of the posterior distribution inside of a range of values defined as the ROPE.	Provides information related to the practical relevance of the effects.	A ROPE range needs to be arbitrarily defined. Sensitive to the scale (the unit) of the predictors.
Bayes factor (vs. 0)	The degree by which the probability mass has shifted away from or towards the null value, after observing the data.	Ratio of the density of the null value between the posterior and the prior distributions.	An unbounded continuous measure of relative evidence. Allows statistically supporting the null hypothesis.	Sensitive to selection of prior distribution shape, location and scale.
Bayes factor (vs. ROPE)	The degree by which the probability mass has into or outside of the null interval (ROPE), after observing the data.	Ratio of the odds of the posterior vs the prior distribution falling inside of the range of values defined as the ROPE.	An unbounded continuous measure of relative evidence. Allows statistically supporting the null hypothesis. Compared to the BF (vs 0), evidence is accumulated faster for the null when the null is true.	Sensitive to selection of prior distribution shape, location and scale. Additionally, a ROPE range needs to be arbitrarily defined, which is sensitive to the scale (the unit) of the predictors.

337 than some given threshold. This aspect can be explored, to a certain extent, in an objective
338 way with the concept of *practical equivalence* (Kruschke, 2014; Lakens, 2017; Lakens et al.,
339 2018), which suggests the use of a range of values assimilated to the absence of an effect (the
340 ROPE). If the effect falls within this range, it is considered as non-significant *for practical*
341 *reasons*: the magnitude of the effect is likely to be too small to be of high importance in
342 real-world scenarios or applications. Nevertheless, *significance* also withdraws a more
343 subjective aspect, corresponding to its contextual meaningfulness and relevance. This,
344 however, is usually dependent on the literature, priors, novelty, context or field, and thus
345 cannot be objectively or neutrally assessed with a statistical index alone.

346 While indices of existence and significance can be numerically related (as shown in our
347 results), the former is conceptually independent from the latter. For example, an effect for
348 which the whole posterior distribution is concentrated within the [0.0001, 0.0002] range
349 would be considered as positive with a high certainty (and thus, *existing* in a that direction),
350 but also not significant (i.e., too small to be of any practical relevance). Acknowledging the
351 distinction and complementary of these two aspects can in turn enrich the information and
352 usefulness of the results reported in psychological science (for practical reasons, the
353 implementation of this dual-framework of existence-significance testing is made
354 straightforward through the *bayestestR* open-source package for R; Makowski et al., 2019).
355 In this context, the *pd* and the MAP-based *p*-value appear as indices of effect existence,
356 mostly sensitive to the certainty related to the direction of the effect. ROPE-based indices
357 and Bayes factors are indices of effect significance, related to the magnitude and the amount
358 of evidence in favor of it (see also a similar discussion of statistical significance vs. effect size
359 in the frequentist framework; e.g., Cohen, 2016)

360 The inherent subjectivity related to the assessment of significance is one of the
361 practical limitation the ROPE-based indices (although being, conceptually, an asset,
362 allowing for contextual nuance in the interpretation), as they require an explicit definition of

363 the non-significant range (the ROPE). Although default values were reported in the
364 literature (for instance, half of a “negligible” effect size reference value; Kruschke, 2014), it is
365 critical for the reproducibility and transparency that the researcher’s choice is explicitly
366 stated (and, if possible, justified). Beyond being arbitrary, this range also has hard bounds
367 (for instance, contrary to a value of 0.0499, a value of 0.0501 would be considered as
368 non-negligible if the range ends at 0.05). This reinforces a categorical and clustered
369 perspective of what is by essence a continuous space of possibilities. Importantly, as this
370 range is fixed to the scale of the response (it is expressed in the unit of the response), ROPE
371 indices are sensitive to changes in the scale of the predictors. For instance, negligible results
372 may change into non-negligible results when predictors are scaled up (e.g. express reaction
373 times in seconds instead of milliseconds), which one inattentive or malicious researcher could
374 misleadingly present as “significant” (note that indices of existence, such as the pd , would
375 not be affected). Finally, the ROPE definition is also dependent on the model type, and
376 selecting a consistent or homogeneous range for all the families of models is not
377 straightforward. This can make comparisons between model types difficult, and an
378 additional burden when interpreting ROPE-based indices. In summary, while a well-defined
379 ROPE can be a powerful tool to give a different and new perspective, it also requires extra
380 caution from the authors and the readers.

381 As for the difference between ROPE (95%) and ROPE (full), we suggest reporting the
382 latter (i.e., the percentage of the whole posterior distribution that falls within the ROPE
383 instead of a given proportion of CI). This bypass the usage of another arbitrary range (95%)
384 and appears to be more sensitive to delineate highly significant effects). Critically, rather
385 than using the percentage in ROPE as a dichotomous, all-or-nothing decision criterion, such
386 as suggested by the original equivalence test (Kruschke, 2014), we recommend using the
387 percentage as a continuous index of significance (with explicitly specified cut-off points if
388 categorization is needed, for instance 5% for significance and 95% for non-significance).

389 Our results underline Bayes factor as an interesting index, able to provide evidence in
390 favor or against the presence of an effect. Moreover, its easy interpretation in terms of odds
391 in favor, or against, one or the other hypothesis makes it a compelling index for
392 communication. Nevertheless, one of the main critiques of Bayes factors, is its sensitivity to
393 priors (shown in our results here through its sensitivity to model types, as priors' odds for
394 logistic and linear models are different). Moreover, while the BF against a ROPE appears as
395 even better than the BF against a point-null, it also carries all the limitations related to the
396 ROPE specification mentioned above. Thus, we recommend using Bayes factors
397 (preferentially *vs.* a ROPE) if the user has explicitly specified (and have a rationale for)
398 informative priors (often called “subjective” priors; Wagenmakers, 2007). In the end, there is
399 a relative proximity between Bayes factors (*vs.* ROPE) and the percentage in ROPE (full),
400 consistently with their mathematical relationship.

401 Being quite different from the Bayes factors and the ROPE indices, the Probability of
402 Direction (pd) is an index of effect existence representing the certainty with which an effect
403 goes in a particular direction (i.e., is positive or negative). Beyond its simplicity of
404 interpretation, understanding and computation, this index also presents other interesting
405 properties. It is independent from the model, i.e., it is solely based on the posterior
406 distributions and does not require any additional information from the data or the model.
407 Contrary to ROPE-based indices, it is robust to the scale of both the response variable and
408 the predictors. Nevertheless, this index also presents some limitations. Most importantly, the
409 pd is not relevant to assess size or importance of the effect and is not able to give
410 information *in favor* of the null hypothesis. In other words, a high pd suggests the presence
411 of an effect but a small pd does not give us any information about how much the null
412 hypothesis is plausible, suggesting that this index can only be used to eventually reject the
413 null hypothesis (which is consistent with the interpretation of the frequentist p -value). On
414 the contrary, the BFs (and to some extent the percentage in ROPE) increase or decrease as
415 the evidence becomes stronger (more data points), in both directions.

Much of the strengths of the pd also apply to the MAP-based p -value. Although possibly showing some superiority in terms of sensitivity as compared to it, it also presents an important limitation. Indeed, the MAP is mathematically dependent on the density at 0 and at the mode. However, the density estimation of a continuous distribution is a statistical problem on its own and many different methods exist. It is possible that changing the density estimation might impact the MAP-based p -value with unknown results. The pd , however, has a linear relationship with the frequentist p -value, which is in our opinion an asset.

After all the criticism regarding the frequentist p -value, it might appear as contradictory to suggest the usage of its Bayesian empirical equivalent. The subtler perspective that we support is that the p -value is not an intrinsically bad, or wrong, index. Instead, it is its misuse, misunderstanding and misinterpretation that fuels the decay of the situation into the crisis. Interestingly, the proximity between the pd and the p -value suggests that the latter is more an index of effect *existence* than *significance* (as in “worth of interest”; Cohen, 2016). Addressing this confusion, the Bayesian equivalent has an intuitive meaning and interpretation, contributing to making more obvious the fact that all thresholds and heuristics are arbitrary. In summary, its mathematical and interpretative transparency of the pd , and its conceptualization as an index of effect existence, offers a valuable insight into the characterization of Bayesian results, and its practical proximity with the frequentist p -value makes it a perfect metric to ease the transition of psychological research into the adoption of the Bayesian framework.

Our study has some limitations. First, our simulations were based on simple linear and logistic regression models. Although these models are widely spread, the behavior of the presented indices for other model families or types, like count models or mixed effects models, still needs to be explored. Furthermore, we only tested continuous predictors. The indices might behave differently when varying the type of predictor (binary, ordinal) as well. Finally, we limited our simulations to small sample sizes, for reasons that data is particularly

442 noisy in small samples, and experiments in psychology often include only a limited number
443 of subjects. However, it is possible that the indices converge (or diverge), for larger samples.
444 Importantly, before being able to draw a definitive conclusion about the qualities of these
445 indices, further studies need to investigate the robustness of these indices to sampling
446 characteristics (*e.g.*, sampling algorithm, number of iterations, chains, warm-up) and the
447 impact of prior specification (Kass & Raftery, 1995; Kruschke, 2011; Vanpaemel, 2010), all of
448 which are important parameters of Bayesian statistics.

449 **Reporting Guidelines**

450 How can the current observations be used to improve statistical good practices in
451 psychological science? Based on the present comparison, we can start outlining the following
452 guidelines. As *existence* and *significance* are complementary perspectives, we suggest using
453 at minimum one index of each category. As an objective index of effect existence, the *pd*
454 should be reported, for its simplicity of interpretation, its robustness and its numeric
455 proximity to the well-known frequentist *p*-value; As an index of significance either the *BF*
456 (*vs.* *ROPE*) or the *ROPE* (*full*) should be reported, for their ability to discriminate between
457 presence and absence of effect (De Santis, 2007), and the information they provide related to
458 evidence of the size of the effect. Selection between the *BF* (*vs.* *ROPE*) or the *ROPE* (*full*)
459 should depend on the informativeness of the priors used - when uninformative priors are
460 used, and there is little prior knowledge regarding the expected size of the effect, the *ROPE*
461 (*full*) should be reported as it reflects only the posterior distribution, and is not sensitive to
462 the width of a wide-range of prior scales (Rouder, Haaf, & Vandekerckhove, 2018). On the
463 other hand, in cases where informed priors are used, reflecting prior knowledge regarding the
464 expected size of the effect, *BF* (*vs.* *ROPE*) should be used.

465 Defining appropriate heuristics to help the interpretation is beyond the scope of this
466 paper, as it would require testing them on more natural datasets. Nevertheless, if we take
467 the frequentist framework and the existing literature as a reference point, it seems that 95%,

468 97% and 99% might be relevant reference points (i.e., easy-to-remember values) for the pd
469 and 3, 10 and 30 (weak evidence) appropriate for the BF. A concise, standardized, reference
470 template sentence to describe the parameter of a model including an index of point-estimate,
471 uncertainty, existence, significance and effect size (Cohen, 1988) could be, in the case of pd
472 and BF :

473 “There is moderate evidence ($BF_{ROPE} = 3.44$) [BF (*vs.* $ROPE$)] in favor of the
474 presence of effect of X, which has a probability of 98.14% [pd] of being negative
475 ($Median = -5.04$, $89\%CI[-8.31., 0.12]$), and can be considered as small
476 ($Std.Median = -0.29$) [*standardized coefficient*]”

477 And if the user decides to use the percentage in ROPE instead of the BF :

478 “The effect of X has a probability of 98.14% [pd] of being negative ($Median = -5.04$,
479 $89\%CI[-8.31, 0.12]$), and can be considered as small ($Std.Median = -0.29$) [*standardized*
480 *coefficient*] and significant (0.82% in $ROPE$) [$ROPE$ (*full*)]”.

481 Data Availability

482 In the spirit of open and honest science, the full R code used for data generation, data
483 processing, figures creation and manuscript compiling is available on GitHub at https://github.com/easystats/easystats/tree/master/publications/makowski_2019_bayesian.

485 Ethics Statement

486 No human participants, but the authors of the present manuscript, were used to
487 produce the current study. The latter verbally reported being endowed with a feeling of
488 free-will at the moment of writing.

Author Contributions

DM conceived and coordinated the study. DM, MSB and DL participated in the study design, statistical analysis, data interpretation and manuscript drafting. DL supervised the manuscript drafting. AC performed a critical review of the manuscript, assisted with manuscript drafting and provided funding for publication. All authors read and approved the final manuscript.

Conflict of Interest Statement

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Acknowledgments

This study was made possible by the development of the **bayestestR** package, itself part of the *easystats* ecosystem (Lüdecke, Waggoner, & Makowski, 2019), an open-source and collaborative project created to facilitate the usage of R. Thus, there is substantial evidence in favor of the fact that we thank the masters of easystats and all the other padawan following the way of the Bayes.

504

References

- 505 Amrhein, V., Greenland, S., & McShane, B. (2019). Scientists rise up against statistical
506 significance. *Nature*, 567(7748), 305–307. doi:10.1038/d41586-019-00857-9
- 507 Anderson, D. R., Burnham, K. P., & Thompson, W. L. (2000). Null hypothesis testing:
508 Problems, prevalence, and an alternative. *The Journal of Wildlife Management*,
509 912–923.
- 510 Andrews, M., & Baguley, T. (2013). Prior approval: The growth of bayesian methods in
511 psychology. *British Journal of Mathematical and Statistical Psychology*, 66(1), 1–7.
- 512 Benjamin, D. J., Berger, J. O., Johannesson, M., Nosek, B. A., Wagenmakers, E.-J., Berk,
513 R., . . . others. (2018). Redefine statistical significance. *Nature Human Behaviour*,
514 2(1), 6.
- 515 Carpenter, B., Gelman, A., Hoffman, M. D., Lee, D., Goodrich, B., Betancourt, M., . . .
516 Riddell, A. (2017). Stan: A Probabilistic Programming Language. *Journal of
517 Statistical Software*, 76(1). doi:10.18637/jss.v076.i01
- 518 Chambers, C. D., Feredoes, E., Muthukumaraswamy, S. D., & Etchells, P. (2014). Instead of
519 ‘playing the game’ it is time to change the rules: Registered reports at aims
520 neuroscience and beyond. *AIMS Neuroscience*, 1(1), 4–17.
- 521 Cohen, J. (1988). Statistical power analysis for the social sciences.
- 522 Cohen, J. (2016). The earth is round ($p < .05$). In *What if there were no significance tests?*
523 (pp. 69–82). Routledge.
- 524 De Santis, F. (2007). Alternative bayes factors: Sample size determination and
525 discriminatory power assessment. *Test*, 16(3), 504–522.

- 526 Dienes, Z. (2014). Using bayes to get the most out of non-significant results. *Frontiers in*
527 *Psychology*, 5, 781.
- 528 Dienes, Z., & Mcclatchie, N. (2018). Four reasons to prefer bayesian analyses over significance
529 testing. *Psychonomic Bulletin & Review*, 25(1), 207–218.
- 530 Ellis, S., & Steyn, H. (2003). Practical significance (effect sizes) versus or in combination
531 with statistical significance (p-values): Research note. *Management Dynamics:*
532 *Journal of the Southern African Institute for Management Scientists*, 12(4), 51–53.
- 533 Etz, A., Haaf, J. M., Rouder, J. N., & Vandekerckhove, J. (2018). Bayesian inference and
534 testing any hypothesis you can specify. *Advances in Methods and Practices in*
535 *Psychological Science*, 2515245918773087.
- 536 Etz, A., & Vandekerckhove, J. (2016). A bayesian perspective on the reproducibility project:
537 Psychology. *PloS One*, 11(2), e0149794.
- 538 Fidler, F., Thomason, N., Cumming, G., Finch, S., & Leeman, J. (2004). Editors can lead
539 researchers to confidence intervals, but can't make them think: Statistical reform
540 lessons from medicine. *Psychological Science*, 15(2), 119–126.
- 541 Finch, S., Cumming, G., Williams, J., Palmer, L., Griffith, E., Alders, C., ... Goodman, O.
542 (2004). Reform of statistical inference in psychology: The case ofMemory & cognition.
543 *Behavior Research Methods, Instruments, & Computers*, 36(2), 312–324.
- 544 Gardner, M. J., & Altman, D. G. (1986). Confidence intervals rather than p values:
545 Estimation rather than hypothesis testing. *Br Med J (Clin Res Ed)*, 292(6522),
546 746–750.
- 547 Gelman, A. (2018). The Failure of Null Hypothesis Significance Testing When Studying
548 Incremental Changes, and What to Do About It. *Personality and Social Psychology*

- 549 *Bulletin*, 44(1), 16–23. doi:10.1177/0146167217729162
- 550 Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D. B., Vehtari, A., & Rubin, D. B. (2014).
551 *Bayesian data analysis*. (Third edition.). Boca Raton: CRC Press.
- 552 Goodrich, B., Gabry, J., Ali, I., & Brilleman, S. (2019). Rstanarm: Bayesian applied
553 regression modeling via Stan. Retrieved from <http://mc-stan.org/>
- 554 Halsey, L. G. (2019). The reign of the p-value is over: What alternative analyses could we
555 employ to fill the power vacuum? *Biology Letters*, 15(5), 20190174.
- 556 Heck, D. W. (2019). A caveat on the savage–dickey density ratio: The case of computing
557 bayes factors for regression parameters. *British Journal of Mathematical and
558 Statistical Psychology*, 72(2), 316–333.
- 559 Jarosz, A. F., & Wiley, J. (2014). What are the odds? A practical guide to computing and
560 reporting bayes factors. *The Journal of Problem Solving*, 7(1), 2.
- 561 Jeffreys, H. (1998). *The theory of probability*. OUP Oxford.
- 562 Kass, R. E., & Raftery, A. E. (1995). Bayes factors. *Journal of the American Statistical
563 Association*, 90(430), 773–795.
- 564 Kirk, R. E. (1996). Practical significance: A concept whose time has come. *Educational and
565 Psychological Measurement*, 56(5), 746–759.
- 566 Kruschke, J. (2014). *Doing bayesian data analysis: A tutorial with r, jags, and stan*.
567 Academic Press.
- 568 Kruschke, J. K. (2010). What to believe: Bayesian methods for data analysis. *Trends in
569 Cognitive Sciences*, 14(7), 293–300.
- 570 Kruschke, J. K. (2011). Bayesian assessment of null values via parameter estimation and

- 571 model comparison. *Perspectives on Psychological Science*, 6(3), 299–312.
- 572 Kruschke, J. K., Aguinis, H., & Joo, H. (2012). The time has come: Bayesian methods for
573 data analysis in the organizational sciences. *Organizational Research Methods*, 15(4),
574 722–752.
- 575 Kruschke, J. K., & Liddell, T. M. (2018). The bayesian new statistics: Hypothesis testing,
576 estimation, meta-analysis, and power analysis from a bayesian perspective.
577 *Psychonomic Bulletin & Review*, 25(1), 178–206.
- 578 Lakens, D. (2017). Equivalence tests: A practical primer for t tests, correlations, and
579 meta-analyses. *Social Psychological and Personality Science*, 8(4), 355–362.
- 580 Lakens, D., Scheel, A. M., & Isager, P. M. (2018). Equivalence testing for psychological
581 research: A tutorial. *Advances in Methods and Practices in Psychological Science*,
582 2515245918770963.
- 583 Lüdecke, D., Waggoner, P., & Makowski, D. (2019). Insight: A unified interface to access
584 information from model objects in r. *Journal of Open Source Software*, 4(38), 1412.
585 doi:10.21105/joss.01412
- 586 Ly, A., Verhagen, J., & Wagenmakers, E.-J. (2016). Harold jeffreys's default bayes factor
587 hypothesis tests: Explanation, extension, and application in psychology. *Journal of
588 Mathematical Psychology*, 72, 19–32.
- 589 Makowski, D., Ben-Shachar, M., & Lüdecke, D. (2019). bayestestR: Describing Effects and
590 their Uncertainty, Existence and Significance within the Bayesian Framework.
591 *Journal of Open Source Software*, 4(40), 1541. doi:10.21105/joss.01541
- 592 Marasini, D., Quatto, P., & Ripamonti, E. (2016). The use of p-values in applied research:
593 Interpretation and new trends. *Statistica*, 76(4), 315–325.

- 594 Maxwell, S. E., Lau, M. Y., & Howard, G. S. (2015). Is psychology suffering from a
595 replication crisis? What does “failure to replicate” really mean? *American*
596 *Psychologist*, 70(6), 487.
- 597 Mills, J. A. (2017). Objective bayesian precise hypothesis testing. *University of Cincinnati*
598 [*Original Version: 2007*].
- 599 Mills, J. A., & Parent, O. (2014). Bayesian mcmc estimation. In *Handbook of regional*
600 *science* (pp. 1571–1595). Springer.
- 601 Morey, R. D., & Rouder, J. N. (2011). Bayes factor approaches for testing interval null
602 hypotheses. *Psychological Methods*, 16(4), 406.
- 603 R Core Team. (2019). *R: A language and environment for statistical computing*. Vienna,
604 Austria: R Foundation for Statistical Computing. Retrieved from
605 <https://www.R-project.org/>
- 606 Robert, C. P. (2014). On the jeffreys-lindley paradox. *Philosophy of Science*, 81(2), 216–232.
- 607 Robert, C. P. (2016). The expected demise of the bayes factor. *Journal of Mathematical*
608 *Psychology*, 72, 33–37.
- 609 Rouder, J. N., Haaf, J. M., & Vandekerckhove, J. (2018). Bayesian inference for psychology,
610 part iv: Parameter estimation and bayes factors. *Psychonomic Bulletin & Review*,
611 25(1), 102–113.
- 612 Rouder, J. N., & Morey, R. D. (2012). Default bayes factors for model selection in regression.
613 *Multivariate Behavioral Research*, 47(6), 877–903.
- 614 Rouder, J. N., Speckman, P. L., Sun, D., Morey, R. D., & Iverson, G. (2009). Bayesian t
615 tests for accepting and rejecting the null hypothesis. *Psychonomic Bulletin & Review*,
616 16(2), 225–237.

- 617 Simmons, J. P., Nelson, L. D., & Simonsohn, U. (2011). False-Positive Psychology:
618 Undisclosed Flexibility in Data Collection and Analysis Allows Presenting Anything
619 as Significant. *Psychological Science*, 22(11), 1359–1366.
620 doi:10.1177/0956797611417632
- 621 Simonsohn, U., Nelson, L. D., & Simmons, J. P. (2014). P-curve and effect size: Correcting
622 for publication bias using only significant results. *Perspectives on Psychological
623 Science*, 9(6), 666–681.
- 624 Spanos, A. (2013). Who should be afraid of the jeffreys-lindley paradox? *Philosophy of
625 Science*, 80(1), 73–93.
- 626 Sullivan, G. M., & Feinn, R. (2012). Using effect size—or why the p value is not enough.
627 *Journal of Graduate Medical Education*, 4(3), 279–282.
- 628 Szucs, D., & Ioannidis, J. P. (2016). Empirical assessment of published effect sizes and power
629 in the recent cognitive neuroscience and psychology literature. *BioRxiv*, 071530.
- 630 Vanpaemel, W. (2010). Prior sensitivity in theory testing: An apologia for the bayes factor.
631 *Journal of Mathematical Psychology*, 54(6), 491–498.
- 632 Wagenmakers, E.-J. (2007). A practical solution to the pervasive problems of p values.
633 *Psychonomic Bulletin & Review*, 14(5), 779–804.
- 634 Wagenmakers, E.-J., Lee, M., Rouder, J., & Morey, R. (2019, August). *Another statistical
635 paradox*. Retrieved from
636 <http://www.ejwagenmakers.com/submitted/AnotherStatisticalParadox.pdf>
- 637 Wagenmakers, E.-J., Lodewyckx, T., Kuriyal, H., & Grasman, R. (2010). Bayesian
638 hypothesis testing for psychologists: A tutorial on the savage–dickey method.
639 *Cognitive Psychology*, 60(3), 158–189.

- 640 Wagenmakers, E.-J., Marsman, M., Jamil, T., Ly, A., Verhagen, J., Love, J., ... others.
641 (2018). Bayesian inference for psychology. Part i: Theoretical advantages and
642 practical ramifications. *Psychonomic Bulletin & Review*, 25(1), 35–57.
- 643 Wagenmakers, E.-J., Morey, R. D., & Lee, M. D. (2016). Bayesian benefits for the pragmatic
644 researcher. *Current Directions in Psychological Science*, 25(3), 169–176.
- 645 Wagenmakers, E.-J., Verhagen, J., Ly, A., Matzke, D., Steingroever, H., Rouder, J. N., &
646 Morey, R. D. (2017). The need for bayesian hypothesis testing in psychological
647 science. *Psychological Science Under Scrutiny: Recent Challenges and Proposed
648 Solutions*, 123–138.
- 649 Wasserstein, R. L., Lazar, N. A., & others. (2016). The asa's statement on p-values:
650 Context, process, and purpose. *The American Statistician*, 70(2), 129–133.