# Indices of Effect Existence and Significance in the Bayesian Framework

1

2

3    Dominique Makowski [1, *], Mattan S. Ben-Shachar [2], Daniel Lüdecke [3, †], & S.H. Annabel Chen [1,4, †]

4

5

6                    [1] Nanyang Technological University, Singapore

7                    [2] Ben-Gurion University of the Negev, Israel

8                    [3] University Medical Center Hamburg-Eppendorf, Germany

9            [4] Centre for Research and Development in Learning (CRADLE), Singapore

10

11

12

13    [*] Correspondence concerning this article should be addressed to Dominique Makowski, HSS 04-18,
14    48 Nanyang Avenue, Singapore. E-mail: dmakowski@ntu.edu.sg.

15    [†] Daniel Lüdecke and S.H. Annabel Chen share senior authorship.

16

17      **Abstract**

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

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

36      Word count: 6293

# Indices of Effect Existence and Significance in the Bayesian Framework

## 1    Introduction

38   The Bayesian framework is quickly gaining popularity among psychologists and neuroscientists
39   (Andrews & Baguley, 2013). Reasons to prefer this approach are reliability, better accuracy in noisy
40   data, better estimation for small samples, less proneness to type I errors, the possibility of introducing
41   prior knowledge into the analysis and the intuitiveness and straightforward interpretation of results
42   (Dienes & Mclatchie, 2018; Etz & Vandekerckhove, 2016; Kruschke, 2010; Kruschke, Aguinis, &
43   Joo, 2012; Wagenmakers et al., 2018; Wagenmakers, Morey, & Lee, 2016). On the other hand, the
44   frequentist approach has been associated with the focus on *p*-values and null hypothesis significance
45   testing (NHST). The misinterpretation and misuse of *p*-values, so called 'p-hacking' (Simmons,
46   Nelson, & Simonsohn, 2011), has been shown to critically contribute to the reproducibility crisis in
47   psychological science (Chambers, Feredoes, Muthukumaraswamy, & Etchells, 2014; Szucs &
48   Ioannidis, 2016). Not only are *p*-values used to draw inappropriate inferences from noisy data, but
49   even when used properly, effects are drastically overestimated, sometimes even in the wrong
50   direction, when estimation is tied to statistical significance in highly variable data (Gelman, 2018). In
51   response, there is a general agreement that the generalization and utilization of the Bayesian
52   framework is one way of overcoming these issues (Benjamin et al., 2018; Etz & Vandekerckhove,
53   2016; Halsey, 2019; Marasini, Quatto, & Ripamonti, 2016; Maxwell, Lau, & Howard, 2015;
54   Wagenmakers et al., 2017).

55   The tenacity and resilience of the *p*-value as an index of significance is remarkable, despite the long-
56   lasting criticism and discussion about its misuse and misinterpretation (Anderson, Burnham, &
57   Thompson, 2000; Cohen, 2016; Fidler, Thomason, Cumming, Finch, & Leeman, 2004; Finch et al.,
58   2004; Gardner & Altman, 1986). This endurance might be informative on how such indices, and the
59   accompanying heuristics applied to interpret them (e.g., assigning thresholds like .05, .01 and .001 to
60   certain levels of significance), are useful and necessary for researchers to gain an intuitive (although
61   possibly simplified) understanding of the interactions and structure of their data. Moreover, the utility
62   of such an index is most salient in contexts where decisions must be made and rationalized (e.g., in
63   medical settings). Unfortunately, these heuristics can become severely rigidified, and meeting
64   significance has become a goal unto itself rather than a tool for understanding the data (Cohen, 2016;
65   Kirk, 1996). This is particularly problematic given that *p*-values can only be used to reject the null
66   hypothesis and not to accept it as true, because a statistically non-significant result does not mean
67   that there is no difference between groups or no effect of a treatment (Amrhein, Greenland, &
68   McShane, 2019; Wagenmakers, 2007).

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

79     The Bayesian framework offers several alternative indices for the *p*-value. To better understand these
80     indices, it is important to point out one of the core differences between Bayesian and frequentist
81     methods. From a frequentist perspective, the effects are fixed (but unknown) and data are random.
82     On the other hand, instead of having single estimates of some "true effect" (for instance, the "true"
83     correlation between *x* and *y*), Bayesian methods compute the probability of different effects values
84     *given* the observed data (and some prior expectation), resulting in a distribution of possible values for
85     the parameters, called the posterior distribution. The description of the posterior distribution (e.g.,
86     through its centrality, dispersion, etc.) allows to draw conclusions from Bayesian analyses.

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

106    Despite the richness provided by the Bayesian framework and the availability of multiple indices, no
107    consensus has yet emerged on which ones to be used. Literature continues to bloom in a raging
108    debate, often polarized between proponents of the Bayes factor as the supreme index and its
109    detractors (Robert, 2014, 2016; Spanos, 2013; Wagenmakers, Lee, Rouder, & Morey, 2019), with
110    strong theoretical arguments being developed on both sides. Yet no practical, empirical and direct
111    comparison between these indices has been done. This might be a deterrent for scientists interested in
112    adopting the Bayesian framework. Moreover, this grey area can increase the difficulty of readers or
113    reviewers unfamiliar with the Bayesian framework to follow the assumptions and conclusions, which
114    could in turn generate unnecessary doubt upon an entire study. While we think that such indices of
115    significance and their interpretation guidelines (in the form of rules of thumb) are useful in practice,
116    we also strongly believe that they should be accompanied with the understanding of their "behavior"
117    in relationship with major sources of variance, such as sample size, noise or effect presence. This
118    knowledge is important for people to implicitly and intuitively appraise the meaning and implication
119    of the mathematical values they report. Such an understanding could prevent the crystallization of the
120    possible heuristics and categories derived from such indices, as has unfortunately occurred for the *p*-
121    values.

122    Thus, based on the simulation of linear and logistic regressions (arguably some of the most widely
123    used models in the psychological sciences), the present work aims at comparing several indices of
124    effect "significance", provide visual representations of the "behavior" of such indices in relationship

125  with sample size, noise and effect presence, as well as their relationship to frequentist *p*-values (an
126  index which, beyond its many flaws, is well known and could be used as a reference for Bayesian
127  neophytes), and finally draw recommendations for Bayesian statistics reporting.

## 2    Methods

### 2.1   Data Simulation

130  We simulated datasets suited for linear and logistic regression and started by simulating an
131  independent, normally distributed *x* variable (with mean 0 and SD 1) of a given sample size. Then,
132  the corresponding *y* variable was added, having a perfect correlation (in the case of data for linear
133  regressions) or as a binary variable perfectly separated by *x*. The case of no effect was simulated by
134  creating a *y* variable that was independent of (i.e. not correlated to) *x*. Finally, a Gaussian noise was
135  added to the *x* variable (the error).

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

### 2.2   Indices

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

#### 2.2.1 Frequentist *p*-value

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

#### 2.2.2 Probability of Direction (*pd*)

155  The *Probability of Direction (pd)* varies between 50% and 100% and can be interpreted as the
156  probability that a parameter (described by its posterior distribution) is strictly positive or negative
157  (whichever is the most probable). It is mathematically defined as the proportion of the posterior
158  distribution that is of the median's sign (Makowski, Ben-Shachar, & Lüdecke, 2019).

#### 2.2.3 MAP-based *p*-value

160  The *MAP-based p-value* is related to the odds that a parameter has against the null hypothesis (Mills,
161  2017; Mills & Parent, 2014). It is mathematically defined as the density value at 0 divided by the
162  density at the Maximum A Posteriori (MAP), i.e., the equivalent of the mode for continuous
163  distributions.

164 **2.2.4 ROPE (95%)**

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

170 **2.2.5 ROPE (full)**

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

173 **2.2.6 Bayes factor (*vs.* 0)**

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

181 **2.2.7 Bayes factor (*vs.* ROPE)**

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

188
$$BF_{rope} = \frac{odds(ROPE_{\text{full posterior}})}{odds(ROPE_{\text{full prior}})}$$
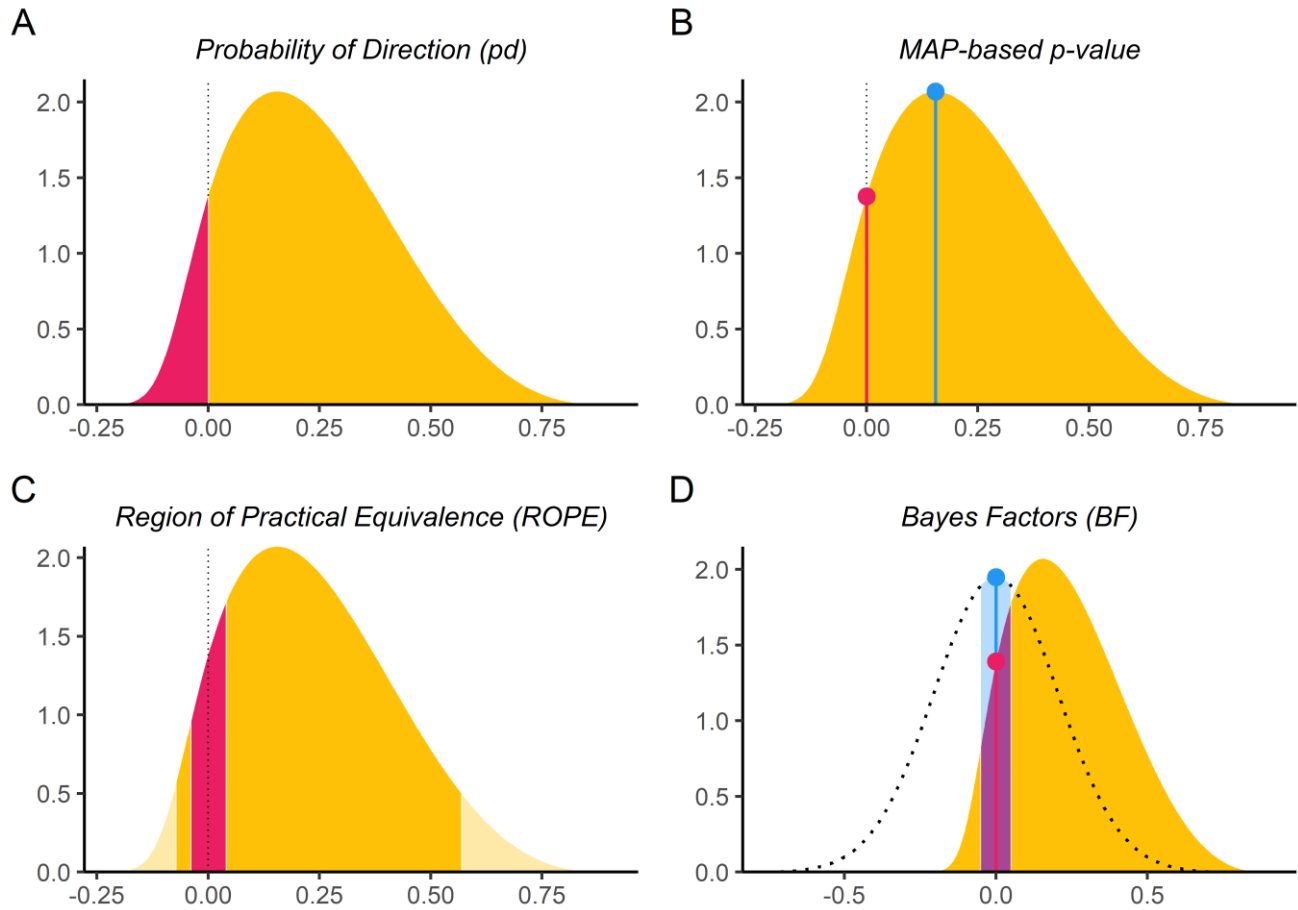
**Figure 1**. Bayesian indices of effect existence and significance. (A) The Probability of Direction (*pd*) is defined as the proportion of the posterior distribution that is of the median's sign (the size of the yellow area relative to the whole distribution). (B) The MAP-based *p*-value is defined 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).

## 2.3 Data Analysis

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

210    In order to provide an intuitive understanding of values, data processing will focus on creating clear
211    visual figures to help the user grasp the patterns and variability that exists when computing the
212    investigated indices. Nevertheless, we decided to also mathematically test our claims in cases where
213    the graphical representation begged for a deeper investigation. Thus, we fitted two regression models
214    to assess the impact of sample size and noise, respectively. For these models (but not for the figures),
215    to ensure that any differences between the indices are not due to differences in their scale or
216    distribution, we converted all indices to the same scale by normalizing the indices between 0 and 1
217    (note that *BF*s were transformed to posterior probabilities, assuming uniform prior odds) and
218    reversing the *p*-values, the MAP-based *p*-values and the ROPE indices so that a higher value
219    corresponds to stronger "significance".

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

227    **3    Results**
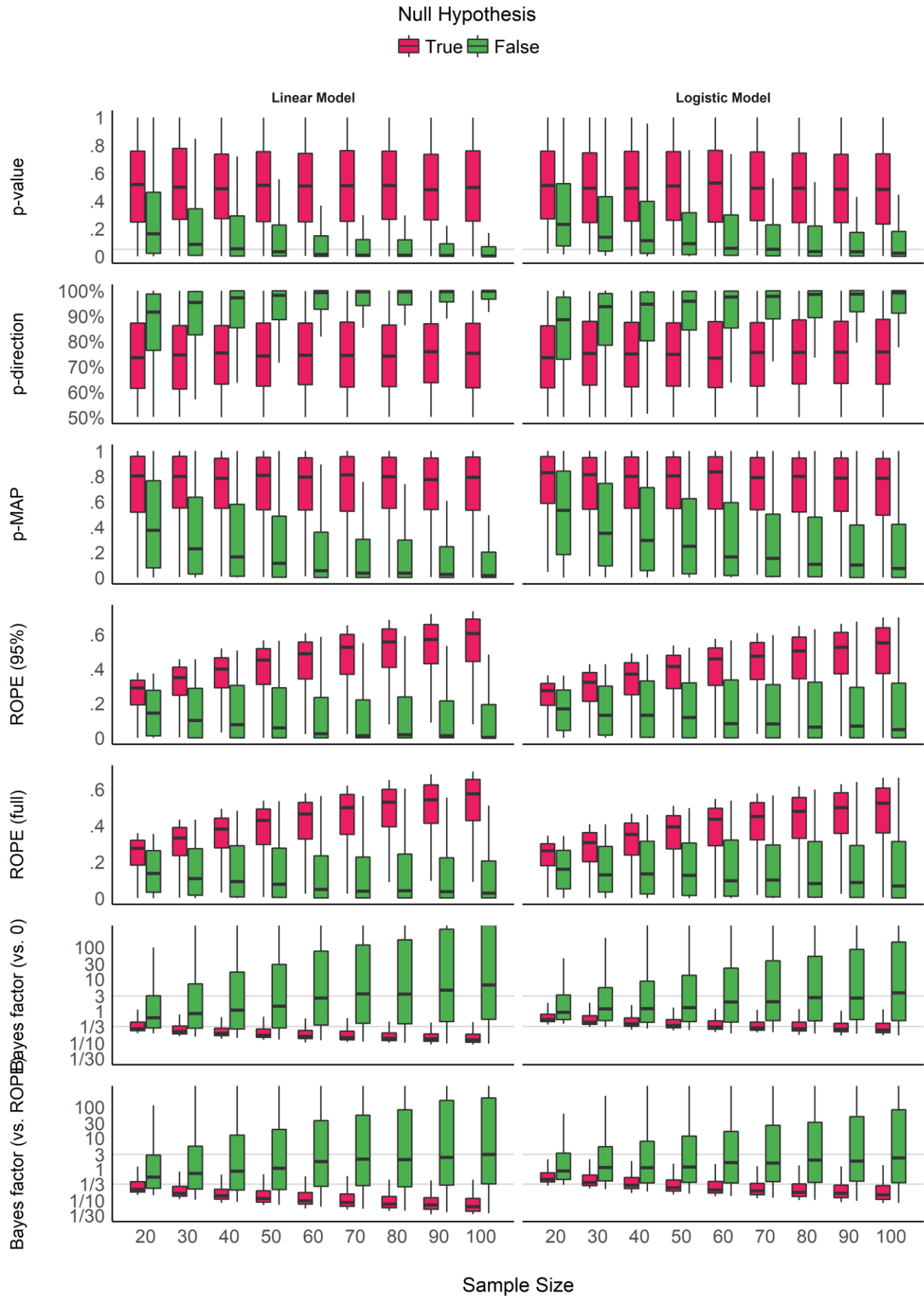
228    **3.1    Impact of Sample Size**

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

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

241 **Table 1**. Sensitivity to sample size. This table shows the standardized coefficient between the sample
242 size and the value of each index, adjusted for error, and stratified by model type and presence of true
243 effect. The stronger the coefficient is, the stronger the relationship with sample size.

| Index | Linear Models / Presence of Effect | Linear Models / Absence of Effect | Logistic Models / Presence of Effect | Logistic Models / Absence of Effect |
|---|---|---|---|---|
| *p*-value | 0.166 | 0.008 | 0.157 | 0.020 |
| *p*-direction | 0.171 | 0.013 | 0.154 | 0.024 |
| *p*-MAP | 0.239 | 0.002 | 0.238 | 0.032 |
| ROPE (95%) | 0.033 | 0.359 | 0.008 | 0.310 |
| ROPE (full) | 0.025 | 0.363 | 0.016 | 0.315 |
| Bayes factor (vs. 0) | 0.198 | 0.116 | 0.116 | 0.141 |
| Bayes factor (vs. ROPE) | 0.152 | 0.136 | 0.078 | 0.180 |

244 Consistently with **Figure 2**, the model investigating the sensitivity of sample size on the different
245 indices suggests that *BF* indices are sensitive to sample size both when an effect is present (null
246 hypothesis is false) and absent (null hypothesis is true). *ROPE* indices are particularly sensitive to
247 sample size when the null hypothesis is true, while *p*-value, *pd* and MAP-based *p*-value are only
248 sensitive to sample size when the null hypothesis is false, in which case they are more sensitive than
249 *ROPE* indices. These findings can be related to the concept of consistency: as the number of data
250 points increases, the statistic converges toward some "true" value. Here, we observe that *p*-value, *pd*
251 and the MAP-based *p*-value are consistent only when the null hypothesis is false. In other words, as
252 sample size increases, they tend to reflect more strongly that the effect is present. On the other hand,
253 *ROPE* indices appear as consistent when the effect is absent. Finally, *BFs* are consistent both when
254 the effect is absent and when it is present, and *BF (vs. ROPE)*, compared to *BF (vs. 0)*, is more
255 sensitive to sample size when the null hypothesis is true, and *ROPE (full)* is overall slightly more
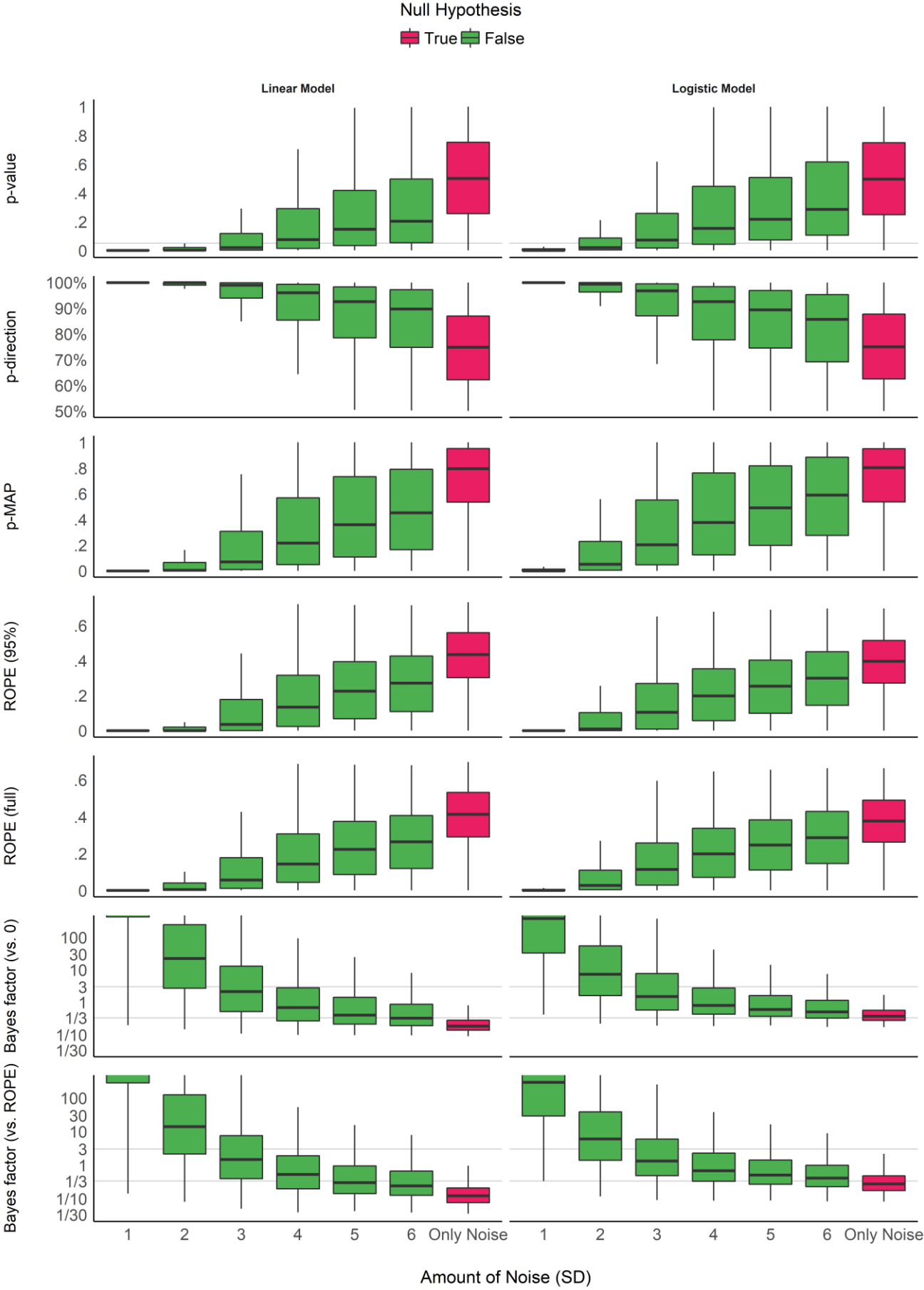256 consistent than *ROPE (95%)*.

257 **3.2 Impact of Noise**

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

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

273 **Table 2**. Sensitivity to noise. This table shows the standardized coefficient between noise and the
274 value of each index when the true effect is present, adjusted for sample size and stratified by model
275 type. The stronger the coefficient is, the stronger the relationship with noise.

| Index | Linear Models / Presence of Effect | Logistic Models / 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 |

276 Consistently with **Figure 3**, the model investigating the sensitivity of noise when an effect is present
277 (as there is only noise in the absence of effect), adjusted for sample size, suggests that BFs
278 (especially *vs.* ROPE), followed by the MAP-based *p*-value and percentages in *ROPE*, are the most
279 sensitive to noise. As noise is a proxy of effect size (linearly related to the absolute value of the
280 coefficient of the parameter), this result highlights the fact that these indices are sensitive to the
281 magnitude of the effect. For example, as noise increases, evidence for an effect becomes weak, and
282 data seems to support the absence of an effect (or at the very least the presence of a negligible effect),
283 which is reflected in *BF*s being consistently smaller than 1. On the other hand, as the *p*-value and the

284  *pd* quantify evidence only for the presence of an effect, as noise increases, they are become more
285  dependent on larger sample size to be able to detect the presence of an effect.

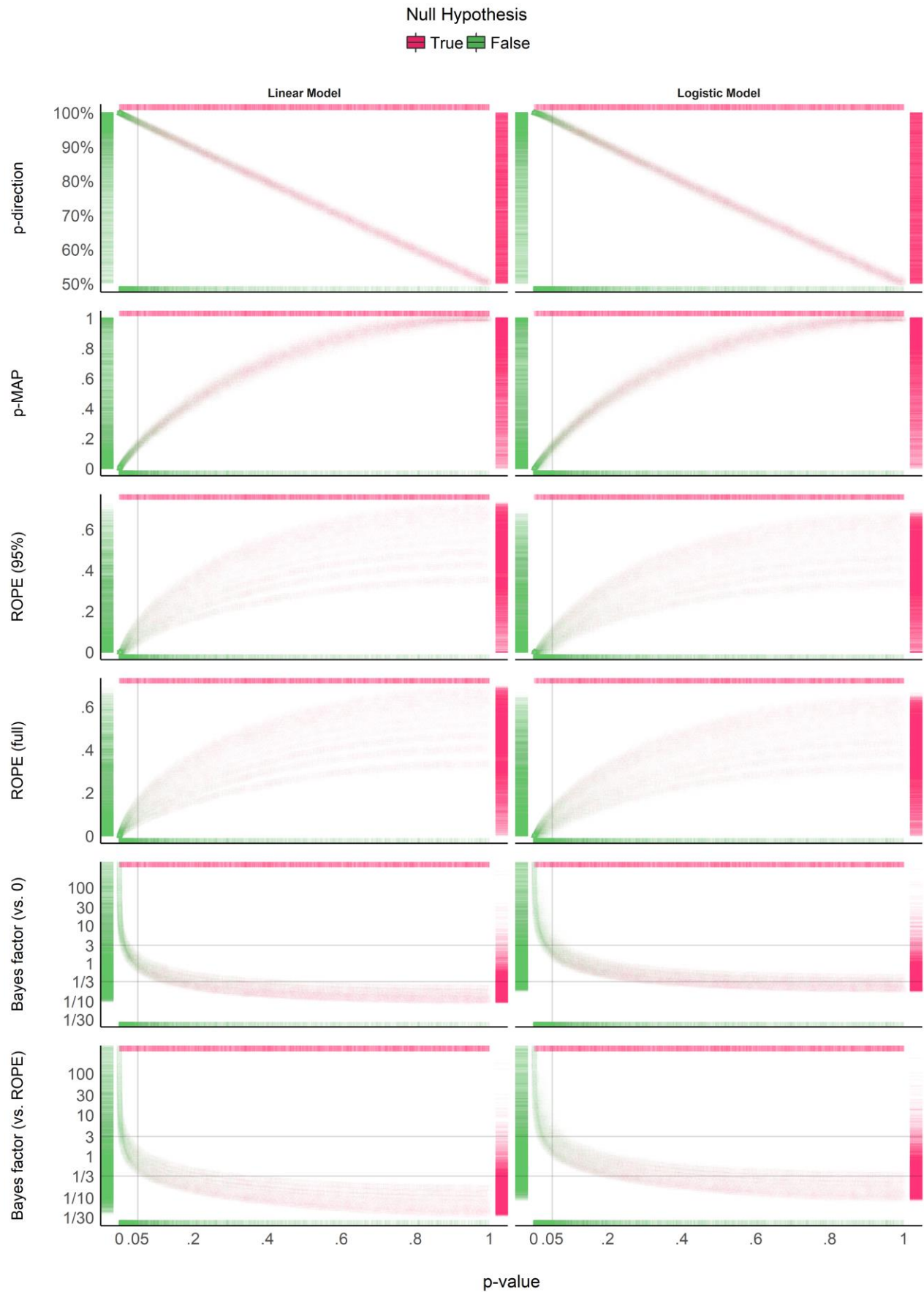## 3.3   Relationship with the frequentist *p*-value

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

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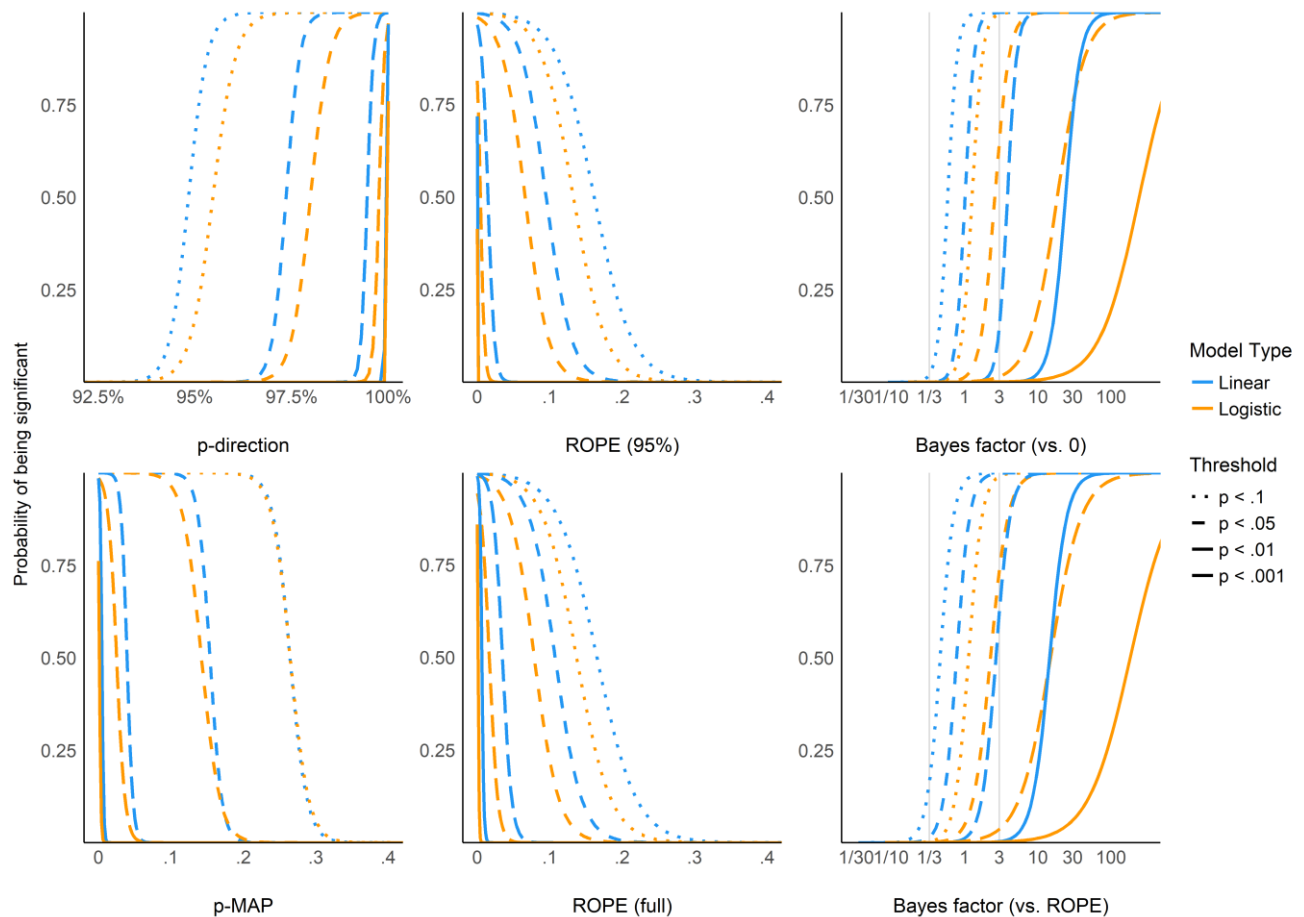299    strongly dependent on sample size.



300

301    **Figure 5**. The probability of reaching different *p*-value based significance thresholds (.1, .05, .01,
302    .001 for solid, long-dashed, short-dashed and dotted lines, respectively) for different values of the
303    corresponding Bayesian indices.

304    **Figure 5** shows equivalence between *p*-value thresholds (.1, .05, .01, .001) and the Bayesian indices.
305    As expected, the *pd* has the sharpest thresholds (95%, 97.5%, 99.5% and 99.95%, respectively). For
306    logistic models, these threshold points appear as more conservative (i.e., Bayesian indices have to be
307    more "pronounced" to reach the same level of significance). This sensitivity to model type is the

308 strongest for BFs (which is possibly related to the difference in the prior specification for these two
309 types of models).

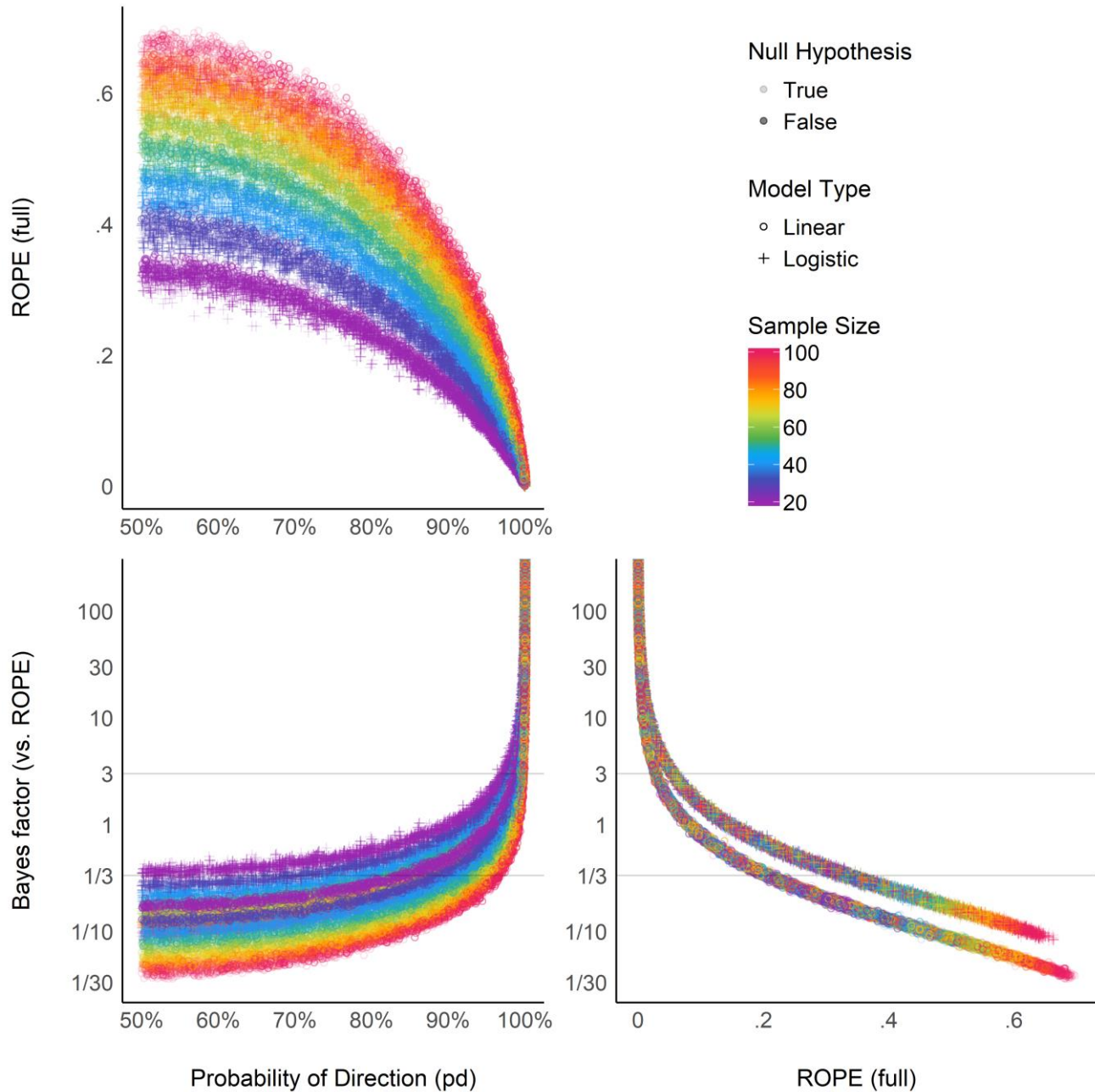310 **3.4    Relationship between ROPE (full), pd and BF (vs. ROPE)**



311

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

314 **Figure 6** suggests that the relationship between the *ROPE (full)* and the *pd* might be strongly affected
315 by the sample size, and subject to differences across model types. This seems to echo the relationship
316 between *ROPE (full)* and *p*-value, the latter having a 1:1 correspondence with *pd*. On the other hand,
317 the *ROPE (full)* and the *BF (vs. ROPE)* seem very closely related within the same model type,
318 reflecting their formal relationship (see definition of *BF (vs. ROPE)* above). Overall, these results

16

319  help to demonstrate *ROPE (full)* and *BF (vs. ROPE)*'s consistency both in case of presence and
320  absence of a true effect, whereas the *pd*, being equivalent to the *p*-value, is only consistent when the
321  true effect is absent.

## 4    Discussion

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

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

340  We also showed that besides sharing similar properties, the *pd* has a 1:1 correspondence with the
341  frequentist *p*-value, being its Bayesian equivalent. On the contrary Bayes factors appear as having a
342  severely non-linear relationship with the frequentist index, which is to be expected from their
343  mathematical definition and their sensitivity when the null hypothesis is true. This in turn can lead to
344  surprising conclusions. For instance, Bayes factors lower than 1, which are considered as providing
345  evidence *against* the presence of an effect, can still correspond to a "significant" frequentist *p*-value
346  (see **Figures 3 and 4**). ROPE indices are more closely related to the *p*-value, as their relationship
347  appears dependent on another factor, the sample size. This suggests that the ROPE encapsulates
348  additional information about the strength of evidence.

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

359  **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, | Ratio of the density of the null value between the posterior and | An unbounded continuous measure of relative evidence. Allows statistically supporting the null hypothesis. | Sensitive to selection of prior distribution shape, location and scale. |

| | | | | |
|---|---|---|---|---|
| | after observing the data. | the prior distributions. | | |
| 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. |

360 Critically, while the purpose of these indices was solely referred to as *significance* until now, we
361 would like to emphasize the nuanced perspective of the existence-significance testing as a dual-
362 framework for parameters description and interpretation. The idea supported here is that there is a
363 conceptual and practical distinction, and possible dissociation to be made, between an effect's
364 existence *and* significance. In this context, *existence* is simply defined as the consistency of an effect
365 in one particular direction (i.e., positive or negative), without any assumptions or conclusions as to its
366 size, importance, relevance or meaning. It is an objective feature of an estimate (tied to its
367 uncertainty). On the other hand, *significance* would be here re-framed following its original literally
368 definition such as "being worthy of attention" or "importance". An effect can be considered
369 significant if its magnitude is higher than some given threshold. This aspect can be explored, to a
370 certain extent, in an objective way with the concept of *practical equivalence* (Kruschke, 2014;
371 Lakens, 2017; Lakens et al., 2018), which suggests the use of a range of values assimilated to the
372 absence of an effect (the ROPE). If the effect falls within this range, it is considered as non-
373 significant *for practical reasons*: the magnitude of the effect is likely to be too small to be of high
374 importance in real-world scenarios or applications. Nevertheless, *significance* also withholds a more
375 subjective aspect, corresponding to its contextual meaningfulness and relevance. This, however, is
376 usually dependent on the literature, priors, novelty, context or field, and thus cannot be objectively or
377 neutrally assessed with a statistical index alone.

378 While indices of existence and significance can be numerically related (as shown in our results), the
379 former is conceptually independent from the latter. For example, an effect for which the whole
380 posterior distribution is concentrated within the [0.0001, 0.0002] range would be considered as
381 positive with a high certainty (and thus, *existing* in a that direction), but also not significant (i.e., too
382 small to be of any practical relevance). Acknowledging the distinction and complementary of these
383 two aspects can in turn enrich the information and usefulness of the results reported in psychological
384 science (for practical reasons, the implementation of this dual-framework of existence-significance
385 testing is made straightforward through the *bayestestR* open-source package for R; Makowski et al.,
386 2019). In this context, the *pd* and the MAP-based *p*-value appear as indices of effect existence,
387 mostly sensitive to the certainty related to the direction of the effect. ROPE-based indices and Bayes
388 factors are indices of effect significance, related to the magnitude and the amount of evidence in
389 favor of it (see also a similar discussion of statistical significance vs. effect size in the frequentist
390 framework; e.g., Cohen, 2016)

391     The inherent subjectivity related to the assessment of significance is one of the practical limitation
392     the ROPE-based indices (although being, conceptually, an asset, allowing for contextual nuance in
393     the interpretation), as they require an explicit definition of the non-significant range (the ROPE).
394     Although default values were reported in the literature (for instance, half of a "negligible" effect size
395     reference value; Kruschke, 2014), it is critical for the reproducibility and transparency that the
396     researcher's choice is explicitly stated (and, if possible, justified). Beyond being arbitrary, this range
397     also has hard bounds (for instance, contrary to a value of 0.0499, a value of 0.0501 would be
398     considered as non-negligible if the range ends at 0.05). This reinforces a categorical and clustered
399     perspective of what is by essence a continuous space of possibilities. Importantly, as this range is
400     fixed to the scale of the response (it is expressed in the unit of the response), ROPE indices are
401     sensitive to changes in the scale of the predictors. For instance, negligible results may change into
402     non-negligible results when predictors are scaled up (e.g. express reaction times in seconds instead of
403     milliseconds), which one inattentive or malicious researcher could misleadingly present as
404     "significant" (note that indices of existence, such as the *pd*, would not be affected). Finally, the
405     ROPE definition is also dependent on the model type, and selecting a consistent or homogeneous
406     range for all the families of models is not straightforward. This can make comparisons between
407     model types difficult, and an additional burden when interpreting ROPE-based indices. In summary,
408     while a well-defined ROPE can be a powerful tool to give a different and new perspective, it also
409     requires extra caution from the authors and the readers.

410     As for the difference between ROPE (95%) and ROPE (full), we suggest reporting the latter (i.e., the
411     percentage of the whole posterior distribution that falls within the ROPE instead of a given
412     proportion of CI). This bypass the usage of another arbitrary range (95%) and appears to be more
413     sensitive to delineate highly significant effects). Critically, rather than using the percentage in ROPE
414     as a dichotomous, all-or-nothing decision criterion, such as suggested by the original equivalence test
415     (Kruschke, 2014), we recommend using the percentage as a continuous index of significance (with
416     explicitly specified cut-off points if categorization is needed, for instance 5% for significance and
417     95% for non-significance).

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

429     Being quite different from the Bayes factors and the ROPE indices, the Probability of Direction (*pd*)
430     is an index of effect existence representing the certainty with which an effect goes in a particular
431     direction (i.e., is positive or negative). Beyond its simplicity of interpretation, understanding and
432     computation, this index also presents other interesting properties. It is independent from the model,
433     i.e., it is solely based on the posterior distributions and does not require any additional information
434     from the data or the model. Contrary to ROPE-based indices, it is robust to the scale of both the
435     response variable and the predictors. Nevertheless, this index also presents some limitations. Most
436     importantly, the *pd* is not relevant to assess size or importance of the effect and is not able to give

437    information *in favor* of the null hypothesis. In other words, a high *pd* suggests the presence of an
438    effect but a small *pd* does not give us any information about how much the null hypothesis is
439    plausible, suggesting that this index can only be used to eventually reject the null hypothesis (which
440    is consistent with the interpretation of the frequentist *p*-value). On the contrary, the BFs (and to some
441    extent the percentage in ROPE) increase or decrease as the evidence becomes stronger (more data
442    points), in both directions.

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

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

462    Our study has some limitations. First, our simulations were based on simple linear and logistic
463    regression models. Although these models are widely spread, the behavior of the presented indices
464    for other model families or types, like count models or mixed effects models, still needs to be
465    explored. Furthermore, we only tested continuous predictors. The indices might behave differently
466    when varying the type of predictor (binary, ordinal) as well. Finally, we limited our simulations to
467    small sample sizes, for reasons that data is particularly noisy in small samples, and experiments in
468    psychology often include only a limited number of subjects. However, it is possible that the indices
469    converge (or diverge), for larger samples. Importantly, before being able to draw a definitive
470    conclusion about the qualities of these indices, further studies need to investigate the robustness of
471    these indices to sampling characteristics (*e.g.*, sampling algorithm, number of iterations, chains,
472    warm-up) and the impact of prior specification (Kass & Raftery, 1995; Kruschke, 2011; Vanpaemel,
473    2010), all of which are important parameters of Bayesian statistics.

474    **5    Reporting Guidelines**

475    How can the current observations be used to improve statistical good practices in psychological
476    science? Based on the present comparison, we can start outlining the following guidelines. As
477    *existence* and *significance* are complementary perspectives, we suggest using at minimum one index
478    of each category. As an objective index of effect existence, the *pd* should be reported, for its
479    simplicity of interpretation, its robustness and its numeric proximity to the well-known frequentist *p*-
480    value; As an index of significance either the *BF (vs. ROPE)* or the *ROPE (full)* should be reported,

481    for their ability to discriminate between presence and absence of effect (De Santis, 2007), and the
482    information they provide related to evidence of the size of the effect. Selection between the *BF*
483    *(vs. ROPE)* or the *ROPE (full)* should depend on the informativeness of the priors used - when
484    uninformative priors are used, and there is little prior knowledge regarding the expected size of the
485    effect, the *ROPE (full)* should be reported as it reflects only the posterior distribution, and is not
486    sensitive to the width of a wide-range of prior scales (Rouder, Haaf, & Vandekerckhove, 2018). On
487    the other hand, in cases where informed priors are used, reflecting prior knowledge regarding the
488    expected size of the effect, *BF (vs. ROPE)* should be used.

489    Defining appropriate heuristics to help the interpretation is beyond the scope of this paper, as it
490    would require testing them on more natural datasets. Nevertheless, if we take the frequentist
491    framework and the existing literature as a reference point, it seems that 95%, 97% and 99% might be
492    relevant reference points (i.e., easy-to-remember values) for the *pd*. A concise, standardized,
493    reference template sentence to describe the parameter of a model including an index of point-
494    estimate, uncertainty, existence, significance and effect size (Cohen, 1988) could be, in the case of *pd*
495    and *BF*:

496    "There is moderate evidence ($BF_{ROPE} = 3.44$) [*BF (vs. ROPE)*] in favor of the presence of effect of
497    X, which has a probability of 98.14% [*pd*] of being negative ($Median = -5.04$,
498    $89\% CI[-8.31., 0.12]$), and can be considered as small ($Std. Median = -0.29$) [*standardized*
499    *coefficient*]"

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

501    "The effect of X has a probability of 98.14% [*pd*] of being negative ($Median = -5.04$,
502    $89\% CI[-8.31, 0.12]$), and can be considered as small ($Std. Median = -0.29$) [*standardized*
503    *coefficient*] and significant ($0.82\%$ in $ROPE$) [*ROPE (full)*]".

## 504    6    Data Availability

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

## 508    7    Ethics Statement

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

## 512    8    Author Contributions

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

## 517    9    Conflict of Interest Statement

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

## 11 References

526 Amrhein, V., Greenland, S., & McShane, B. (2019). Scientists rise up against statistical significance.
527 *Nature*, *567*(7748), 305–307. https://doi.org/10.1038/d41586-019-00857-9

528 Anderson, D. R., Burnham, K. P., & Thompson, W. L. (2000). Null hypothesis testing: Problems,
529 prevalence, and an alternative. *The Journal of Wildlife Management*, 912–923.

530 Andrews, M., & Baguley, T. (2013). Prior approval: The growth of bayesian methods in psychology.
531 *British Journal of Mathematical and Statistical Psychology*, *66*(1), 1–7.

532 Benjamin, D. J., Berger, J. O., Johannesson, M., Nosek, B. A., Wagenmakers, E.-J., Berk, R., …
533 others. (2018). Redefine statistical significance. *Nature Human Behaviour*, *2*(1), 6.

534 Carpenter, B., Gelman, A., Hoffman, M. D., Lee, D., Goodrich, B., Betancourt, M., … Riddell, A.
535 (2017). Stan: A Probabilistic Programming Language. *Journal of Statistical Software*, *76*(1).
536 https://doi.org/10.18637/jss.v076.i01

537 Chambers, C. D., Feredoes, E., Muthukumaraswamy, S. D., & Etchells, P. (2014). Instead of 'playing
538 the game' it is time to change the rules: Registered reports at aims neuroscience and beyond. *AIMS
539 Neuroscience*, *1*(1), 4–17.

540 Cohen, J. (1988). *Statistical power analysis for the social sciences*.

541 Cohen, J. (2016). The earth is round (p<. 05). In *What if there were no significance tests?* (pp. 69–
542 82). Routledge.

543 De Santis, F. (2007). Alternative bayes factors: Sample size determination and discriminatory power
544 assessment. *Test*, *16*(3), 504–522.

545 Dienes, Z. (2014). Using bayes to get the most out of non-significant results. *Frontiers in
546 Psychology*, *5*, 781.

547 Dienes, Z., & Mclatchie, N. (2018). Four reasons to prefer bayesian analyses over significance
548 testing. *Psychonomic Bulletin & Review*, *25*(1), 207–218.

549 Ellis, S., & Steyn, H. (2003). Practical significance (effect sizes) versus or in combination with
550 statistical significance (p-values): Research note. *Management Dynamics: Journal of the Southern
551 African Institute for Management Scientists*, *12*(4), 51–53.

552 Etz, A., Haaf, J. M., Rouder, J. N., & Vandekerckhove, J. (2018). Bayesian inference and testing any
553 hypothesis you can specify. *Advances in Methods and Practices in Psychological Science*,
554 2515245918773087.

555 Etz, A., & Vandekerckhove, J. (2016). A bayesian perspective on the reproducibility project:
556 Psychology. *PloS One*, *11*(2), e0149794.

557 Fidler, F., Thomason, N., Cumming, G., Finch, S., & Leeman, J. (2004). Editors can lead researchers
558 to confidence intervals, but can't make them think: Statistical reform lessons from medicine.
559 *Psychological Science*, *15*(2), 119–126.

560 Finch, S., Cumming, G., Williams, J., Palmer, L., Griffith, E., Alders, C., … Goodman, O. (2004).
561 Reform of statistical inference in psychology: The case ofMemory & cognition. *Behavior Research*
562 *Methods, Instruments, & Computers*, *36*(2), 312–324.

563 Gardner, M. J., & Altman, D. G. (1986). Confidence intervals rather than p values: Estimation rather
564 than hypothesis testing. *Br Med J (Clin Res Ed)*, *292*(6522), 746–750.

565 Gelman, A. (2018). The Failure of Null Hypothesis Significance Testing When Studying Incremental
566 Changes, and What to Do About It. *Personality and Social Psychology Bulletin*, *44*(1), 16–23.
567 https://doi.org/10.1177/0146167217729162

568 Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D. B., Vehtari, A., & Rubin, D. B. (2014). *Bayesian*
569 *data analysis.* (Third edition). Boca Raton: CRC Press.

570 Goodrich, B., Gabry, J., Ali, I., & Brilleman, S. (2019). *Rstanarm: Bayesian applied regression*
571 *modeling via Stan.* Retrieved from http://mc-stan.org/

572 Halsey, L. G. (2019). The reign of the p-value is over: What alternative analyses could we employ to
573 fill the power vacuum? *Biology Letters*, *15*(5), 20190174.

574 Heck, D. W. (2019). A caveat on the savage–dickey density ratio: The case of computing bayes
575 factors for regression parameters. *British Journal of Mathematical and Statistical Psychology*, *72*(2),
576 316–333.

577 Jarosz, A. F., & Wiley, J. (2014). What are the odds? A practical guide to computing and reporting
578 bayes factors. *The Journal of Problem Solving*, *7*(1), 2.

579 Jeffreys, H. (1998). *The theory of probability*. OUP Oxford.

580 Kass, R. E., & Raftery, A. E. (1995). Bayes factors. *Journal of the American Statistical Association*,
581 *90*(430), 773–795.

582 Kirk, R. E. (1996). Practical significance: A concept whose time has come. *Educational and*
583 *Psychological Measurement*, *56*(5), 746–759.

584 Kruschke, J. (2014). *Doing bayesian data analysis: A tutorial with r, jags, and stan*. Academic Press.

585 Kruschke, J. K. (2010). What to believe: Bayesian methods for data analysis. *Trends in Cognitive*
586 *Sciences*, *14*(7), 293–300.

587  Kruschke, J. K. (2011). Bayesian assessment of null values via parameter estimation and model
588  comparison. *Perspectives on Psychological Science*, *6*(3), 299–312.

589  Kruschke, J. K., Aguinis, H., & Joo, H. (2012). The time has come: Bayesian methods for data
590  analysis in the organizational sciences. *Organizational Research Methods*, *15*(4), 722–752.

591  Kruschke, J. K., & Liddell, T. M. (2018). The bayesian new statistics: Hypothesis testing, estimation,
592  meta-analysis, and power analysis from a bayesian perspective. *Psychonomic Bulletin & Review*,
593  *25*(1), 178–206.

594  Lakens, D. (2017). Equivalence tests: A practical primer for t tests, correlations, and meta-analyses.
595  *Social Psychological and Personality Science*, *8*(4), 355–362.

596  Lakens, D., Scheel, A. M., & Isager, P. M. (2018). Equivalence testing for psychological research: A
597  tutorial. *Advances in Methods and Practices in Psychological Science*, 2515245918770963.

598  Lüdecke, D., Waggoner, P., & Makowski, D. (2019). Insight: A unified interface to access
599  information from model objects in r. *Journal of Open Source Software*, *4*(38), 1412.
600  https://doi.org/10.21105/joss.01412

601  Ly, A., Verhagen, J., & Wagenmakers, E.-J. (2016). Harold jeffreys's default bayes factor hypothesis
602  tests: Explanation, extension, and application in psychology. *Journal of Mathematical Psychology*,
603  *72*, 19–32.

604  Makowski, D., Ben-Shachar, M., & Lüdecke, D. (2019). bayestestR: Describing Effects and their
605  Uncertainty, Existence and Significance within the Bayesian Framework. *Journal of Open Source*
606  *Software*, *4*(40), 1541. https://doi.org/10.21105/joss.01541

607  Marasini, D., Quatto, P., & Ripamonti, E. (2016). The use of p-values in applied research:
608  Interpretation and new trends. *Statistica*, *76*(4), 315–325.

609  Maxwell, S. E., Lau, M. Y., & Howard, G. S. (2015). Is psychology suffering from a replication
610  crisis? What does "failure to replicate" really mean? *American Psychologist*, *70*(6), 487.

611  Mills, J. A. (2017). Objective bayesian precise hypothesis testing. *University of Cincinnati [Original*
612  *Version: 2007]*.

613  Mills, J. A., & Parent, O. (2014). Bayesian mcmc estimation. In *Handbook of regional science* (pp.
614  1571–1595). Springer.

615  Morey, R. D., & Rouder, J. N. (2011). Bayes factor approaches for testing interval null hypotheses.
616  *Psychological Methods*, *16*(4), 406.

617  R Core Team. (2019). *R: A language and environment for statistical computing*. Retrieved from
618  https://www.R-project.org/

619  Robert, C. P. (2014). On the jeffreys-lindley paradox. *Philosophy of Science*, *81*(2), 216–232.

620  Robert, C. P. (2016). The expected demise of the bayes factor. *Journal of Mathematical Psychology*,
621  *72*, 33–37.

Rouder, J. N., Haaf, J. M., & Vandekerckhove, J. (2018). Bayesian inference for psychology, part iv: Parameter estimation and bayes factors. *Psychonomic Bulletin & Review*, *25*(1), 102–113.

Rouder, J. N., & Morey, R. D. (2012). Default bayes factors for model selection in regression. *Multivariate Behavioral Research*, *47*(6), 877–903.

Rouder, J. N., Speckman, P. L., Sun, D., Morey, R. D., & Iverson, G. (2009). Bayesian t tests for accepting and rejecting the null hypothesis. *Psychonomic Bulletin & Review*, *16*(2), 225–237.

Simmons, J. P., Nelson, L. D., & Simonsohn, U. (2011). False-Positive Psychology: Undisclosed Flexibility in Data Collection and Analysis Allows Presenting Anything as Significant. *Psychological Science*, *22*(11), 1359–1366. https://doi.org/10.1177/0956797611417632

Simonsohn, U., Nelson, L. D., & Simmons, J. P. (2014). P-curve and effect size: Correcting for publication bias using only significant results. *Perspectives on Psychological Science*, *9*(6), 666–681.

Spanos, A. (2013). Who should be afraid of the jeffreys-lindley paradox? *Philosophy of Science*, *80*(1), 73–93.

Sullivan, G. M., & Feinn, R. (2012). Using effect size—or why the p value is not enough. *Journal of Graduate Medical Education*, *4*(3), 279–282.

Szucs, D., & Ioannidis, J. P. (2016). Empirical assessment of published effect sizes and power in the recent cognitive neuroscience and psychology literature. *BioRxiv*, 071530.

Vanpaemel, W. (2010). Prior sensitivity in theory testing: An apologia for the bayes factor. *Journal of Mathematical Psychology*, *54*(6), 491–498.

Wagenmakers, E.-J. (2007). A practical solution to the pervasive problems ofp values. *Psychonomic Bulletin & Review*, *14*(5), 779–804.

Wagenmakers, E.-J., Lee, M., Rouder, J., & Morey, R. (2019, August). Another statistical paradox. Retrieved from http://www.ejwagenmakers.com/submitted/AnotherStatisticalParadox.pdf

Wagenmakers, E.-J., Lodewyckx, T., Kuriyal, H., & Grasman, R. (2010). Bayesian hypothesis testing for psychologists: A tutorial on the savage–dickey method. *Cognitive Psychology*, *60*(3), 158–189.

Wagenmakers, E.-J., Marsman, M., Jamil, T., Ly, A., Verhagen, J., Love, J., … others. (2018). Bayesian inference for psychology. Part i: Theoretical advantages and practical ramifications. *Psychonomic Bulletin & Review*, *25*(1), 35–57.

Wagenmakers, E.-J., Morey, R. D., & Lee, M. D. (2016). Bayesian benefits for the pragmatic researcher. *Current Directions in Psychological Science*, *25*(3), 169–176.

Wagenmakers, E.-J., Verhagen, J., Ly, A., Matzke, D., Steingroever, H., Rouder, J. N., & Morey, R. D. (2017). The need for bayesian hypothesis testing in psychological science. *Psychological Science Under Scrutiny: Recent Challenges and Proposed Solutions*, 123–138.

Wasserstein, R. L., Lazar, N. A., & others. (2016). The asa's statement on p-values: Context, process, and purpose. *The American Statistician*, *70*(2), 129–133.