

1 **Twenty steps towards an adequate inferential interpretation of p -values**

2

3 **Contributing authors**

4 *Norbert Hirschauer (Corresponding Author)*

5 Martin Luther University Halle-Wittenberg, Faculty of Natural Sciences III

6 Institute of Agricultural and Nutritional Sciences, Chair of Agribusiness Management

7 Karl-Freiherr-von-Fritsch-Str. 4, 06120 Halle (Saale), Germany

8 norbert.hirschauer@landw.uni-halle.de

9

10 *Oliver Mußhoff*

11 Georg August University Goettingen

12 Department for Agricultural Economics and Rural Development, Farm Management

13 Platz der Göttinger Sieben 5, 37073 Göttingen, Germany

14 Oliver.Musshoff@agr.uni-goettingen.de

15

16 *Claudia Becker*

17 Martin Luther University Halle-Wittenberg, Faculty of Law and Economics

18 Institute of Business Studies, Chair of Statistics

19 Große Steinstraße 73, 06099 Halle (Saale), Germany

20 claudia.becker@wiwi.uni-halle.de

21

22 *Sven Grüner*

23 Martin Luther University Halle-Wittenberg, Faculty of Natural Sciences III

24 Institute of Agricultural and Nutritional Sciences, Chair of Agribusiness Management

25 Karl-Freiherr-von-Fritsch-Str. 4, 06120 Halle (Saale), Germany

26 sven.gruener@googlemail.com

1

2

1 **Twenty steps towards an adequate inferential interpretation of p -values**

2 **Abstract:** We suggest twenty immediately actionable steps to reduce widespread inferential errors related
3 to “statistical significance testing.” Our propositions refer first to the theoretical preconditions for using p -
4 values. They furthermore include wording guidelines as well as structural and operative advice of how to
5 present results, especially in multiple regression analysis. Our propositions aim at fostering the logical
6 consistency of inferential arguments by avoiding false categorical reasoning. They are not aimed at dis-
7 pensing with p -values or completely replacing frequentist approaches by Bayesian statistics.

8 **Keywords:** statistical inference, p -value

9 **1 Introduction**

10 Much and, one might think, enough has been said about the misuses and misinterpretations of the p -value,
11 both before and after the ASA-statement (WASSERSTEIN and LAZAR 2016). But the problems seem to be
12 there to stay. Two related features of the frequentist null hypothesis significance testing (NHST) frame-
13 work are at the origin of most errors: first, the dichotomization of results depending on whether the p -
14 value is below or above some arbitrary threshold (usually 0.05). Second, the associated terminology that
15 speaks of “hypothesis testing” and “statistically significant” as opposed to “statistically non-significant”
16 results. Dichotomization in conjunction with misleading terminology have propagated cognitive biases
17 that seduce even experienced researchers to make logically inconsistent and overconfident inferences,
18 both when p is below and when it is above the conventional “significance” threshold. The following errors
19 seem to be particularly widespread:

- 20 1) use of p -values when there is neither a random sample nor a treatment after random assignment
- 21 2) confusion of statistical and practical significance or complete neglect of effect size
- 22 3) unwarranted binary statements of there being an effect as opposed to no effect, coming along with
- 23 - misinterpretations of p -values below 0.05 as posterior probabilities of the null hypothesis
- 24 - misinterpretations of “significant” results as evidence in favor of the estimated coefficients/effects
- 25 - treatment of effects that are “statistically non-significant” as being zero (confirmation of the null)
- 26 4) inflation of evidence against the null caused by p -hacking or unconsidered multiple comparisons
- 27 5) inflation of effect sizes caused by considering “significant” results only

28 The ASA-statement has highlighted that the p -value does not provide a good measure of evidence regard-
29 ing a hypothesis. It is nonetheless a continuous measure of the strength of evidence against the null hy-
30 pothesis, but only in the sense that small p -values will occur more often if there is an effect compared to
31 no effect (HIRSCHAUER et al. 2017). Joining AMRHEIN et al. (2017), BERRY (2017), GELMAN and CARLIN

1 (2017), GREENLAND (2017), MCSHANE et al. (2017), TRAFIMOW et al. (2017), and many others, we be-
2 lieve that degrading the p -value's continuous message into binary "significance" declarations ("bright line
3 rules") is at the heart of the problem. Since the p -value is so deeply anchored in the minds of most applied
4 researchers, we believe that demanding drastic procedural changes, such as renouncing p -values or com-
5 pletely replacing frequentist approaches by Bayesian statistics, is not the most promising approach for
6 mitigating the serious inferential errors that we see today. Dispensing with significance testing but retain-
7 ing the p -value and adopting small and manageable but efficient steps towards improvement seems to be
8 more promising (AMRHEIN et al. 2017). Adequate steps will have to take account of the idiosyncrasies of
9 each scientific discipline. For example, requirements in the biomedical sciences, which often focus on risk
10 ratios or mean differences between treatments, will at least partly differ from those in the social sciences
11 including economics, which mainly resort to multiple regression analysis.

12 **2 Things to consider in general**

13 Even if one is fully aware of the fundamental pitfalls of NHST, it is difficult to escape the categorical rea-
14 soning that is so entrancingly suggested by its dichotomous "significance" declarations. It is an even more
15 difficult task to provide readers with an interpretative evaluation of the often large numbers of coefficient
16 estimates in multiple regressions that avoids inferential errors. Imagine a regression with several focal
17 variables (predictors), a set of controls, and possibly some secondary covariates (interaction terms, higher-
18 order polynomials, etc.) introduced in the process of model specification. How should we evaluate and
19 comment on the "evidence" as represented in the large number of regression coefficient estimates and
20 their associated p -values? We know that we should dispense with the extremely convenient but misleading
21 dichotomous interpretation. This creates a big problem, however. Despite hundreds or even thousands of
22 papers criticizing the misuses and misinterpretations of the p -value (the don'ts), we still lack appropriate
23 wordings that describe the informational content of a p -value correctly *and* meaningfully. This is because
24 the p -value does not provide a clear rationale or even calculus for statistical inference (GOODMAN 2008)
25 or, as BERRY (2017: 896) formulates more drastically, because "as such it has no inferential content."

26 Our wording difficulties are due to the obscure as well as uncertain and inconclusive informational content
27 of the p -value that as such even precludes making probability statements about hypotheses (GELMAN
28 2016). On the one hand, we know that for a two-sided test "any p -value less than 1 implies that the test
29 [null] hypothesis is not the hypothesis most compatible with the data, because any other hypothesis with a
30 larger p -value would be even more compatible with the data" (GREENLAND et al. 2016: 341). Along the
31 same lines but with a focus on experimental data, GOODMAN (2008: 136) notes that "the effect best sup-
32 ported by the data from a given experiment is always the observed effect, regardless of its significance."
33 On the other hand, while realizing that there is evidence in the data, we commonly interpret the p -value as

1 a first defense line against being fooled by the randomness of sampling (BENJAMINI 2016) when general-
2 izing from our findings to the population. We should meet this defense-line interpretation with caution,
3 however, because the p -value itself is but a summary statistic of data obtained from one-time random
4 sampling. In plausible constellations of noise and sample size, the p -value exhibits wide sample-to-sample
5 variability (HALSEY et al. 2015). This is paralleled by the variability of estimated coefficients over replica-
6 tions. We may easily find a large coefficient in one random sample (overestimation) and a small one in
7 another (underestimation). We must not forget that unbiased estimators estimate correctly on average
8 (HIRSCHAUER et al. 2017). We would thence need *all* estimates from frequent replications – *irrespective*
9 of their p -value and their being large or small – to obtain a good idea of the population effect size. Based
10 on a single sample, we have no way of identifying the p -value below (above) which the associated effect
11 size estimate is too large (too small), but we are very likely to overestimate effect sizes when taking “sig-
12 nificant” results at face value (HIRSCHAUER et al. 2017). Even when finding a highly “significant” result
13 (with, let’s say, a p -value of 0.001), which ironically would be a highly appreciated case in the conven-
14 tional NHST-approach, we cannot make a direct inference and assume the estimated effect to accurately
15 reflect the population effect size. Quite on the contrary! “Under reasonable sample sizes and reasonable
16 population effect sizes, it is the abnormally large sample effect sizes that result in p -values that meet the
17 .05 (or the .005) criterion” (TRAFIMOW et al. 2017: 10). Hence, even seemingly neutral, non-dichotomous
18 representations such as “the retail prices of product A exceed the retail prices of product B by 20% on
19 average ($p < 0.001$)” may be misleading because they insinuate that the evidence against the null can be
20 translated into evidence in favor of the concrete effect (AMRHEIN et al 2017) that we happened to find in a
21 sample.

22 The problem of interpreting p -values is further exacerbated by the fact that multiple comparisons, which
23 are inherent to multiple regression, inflate the strength of evidence against the null as indicated by the p -
24 value. Since the extent of multiple comparisons varies between studies, p -values cannot be compared
25 across different studies. A p -value is a summary statistic that tells us how incompatible the data are with
26 the specified statistical model including the null hypothesis. In a *single* regression, a p -value (for example
27 0.05) represents the conditional probability of finding the observed effect (or even a larger one) in random
28 replications *if* the null hypothesis were true. In contrast, in a *multiple* regression with, let’s say, ten focal
29 predictor variables, we unavoidably make ten comparisons in that we assess the strength of evidence
30 against the null as many times as there are variables of interest. Even if all ten null hypotheses were true,
31 we would have a 40.1% ($1-0.95^{10}$) probability of finding at least one coefficient with $p \leq 0.05$. Finding a
32 low p -value for a coefficient in a multiple regression represents much weaker evidence against the null
33 than finding the same p -value in a single regression. We must furthermore not forget that it is common
34 practice amongst applied economists to retain one model as the final (“best”) model after an often large
35 number of different models have been tried out and evaluated by using some measure of model fit such as

1 the likelihood ratio or the Akaike Information Criterion. We necessarily produce inflated effects and arrive
2 at overconfident conclusions if we assess the strength of evidence in only the “best” model even though
3 multiple alternatives had been tried out before (FORSTMEIER et al. 2016).

4 Remembering that the p -value is but a summary statistic of the data at hand is important because we must
5 avoid all wordings that invite confusion with the posterior (or Bayesian) probability, i.e., the epistemic
6 probability that a hypothesis or scientific proposition about the world is true *given* the evidence from the
7 data. There is a big tension between the correctness and the intuitive meaningfulness of the p -value inter-
8 pretation as people, especially when confronted with “significance” language, seem to be prone to the
9 “inverse probability error.” That is, they often confuse the “conditional probability of data given a hypoth-
10 esis” (p -value) with the “conditional probability of a hypothesis given the data” (posterior probability).
11 “Inverse probability error” is a term coined by COHEN (1994: 997) to emphasize that the p -value “does not
12 tell us what we want to know [i.e., the posterior probability], and [that] we so much want to know what we
13 want to know that, out of desperation, we nevertheless believe that it does.”

14 **3 Suggestions for the use and interpretation of p -values in multiple regressions**

15 Widely scattered over time, disciplines, and journals, a huge amount of criticism regarding the use of p -
16 values (the don'ts) as well as a large number of suggestions for reform (the do's) have accumulated. How-
17 ever, neither abundant criticisms of misuses nor grand visions of how to replace the p -value through other
18 tools of statistical inference such as Bayesian statistics have been of much avail. On the contrary! Not
19 even the ASA-statement seems to have produced much change so far (MATTHEWS 2017). We believe that
20 this is not so much due to researchers not recognizing the problems associated with conventional signifi-
21 cance testing, but rather due to their not knowing what to do instead. With a view to the apparent need for
22 both guidance and (at least some degree of) consensus among scientists, this commentary discusses and
23 suggests immediate reforms that seem to be realistic in the light of the p -value's deep entrenchment in
24 current research practice.

25 We systematically compile suggestions (do's) – none of them new and none of them our own – that jointly
26 seem to represent the most promising set of concrete and immediately actionable steps. Being economists,
27 we focus on suggestions that are relevant for correctly interpreting the results of multiple regression analy-
28 sis, which is the working horse in econometric research. Being pragmatic, we focus on suggestions that
29 are concerned with the analysis of single-sample data, even though we are aware of the advantages of
30 multiple-study designs, meta-analysis, and Bayesian approaches for making valid inferences. In brief, the
31 criteria for the selection of suggestions were as follows: (i) their suitability for furthering a *correct* and
32 *meaningful* interpretation of p -values associated with regression coefficient estimates in single studies;
33 (ii) their capacity to provide small and efficient changes that are *manageable* for all those who are so

1 much accustomed to using p -values that they are probably not ready (yet) to meet the huge challenges of
2 fully Bayesian analyses within a multiple regression framework.

3 Contenting ourselves for the time being with compiling small incremental steps for single-study designs
4 must not be understood as opposition towards more substantial change in the future. Quite on the contrary!
5 We hope that our suggestions will help prepare the field for better study designs and inferential tools and
6 especially more meta-analytical thinking in the long run. More immediately, however, we hope that they
7 will serve as a discussion base or even tool kit that is directly helpful, for example, to editors of economic
8 journals who reflect on best practices and try to revise their editorial policies and guidelines in order to
9 increase the quality of published research. In brief, we address the question of how a typical econometric
10 study, which for the time being refrains from Bayesian statistics and continues to use p -values, should
11 proceed and which wordings it should use to avoid the many inferential errors that are so pervasive at
12 present. It is important to note that some suggestions such as displaying random errors could be criticized
13 as asking for redundant information. Readers of a research paper could in principle compute standard er-
14 rors when effect sizes and p -values are provided. Mathematical redundancy is not a good argument, how-
15 ever. Instead, the question is how we should present information to avoid cognitive biases and foster the
16 logical consistency of inferential arguments through good intuition.

17 Our suggestions are best preceded by a quote by VOGT et al. (2014: 242; 244) who emphasize that the
18 classical tools for statistical inference (including p -values) are inherently based on probability theory: “in
19 research not employing random assignment or random sampling, the classical approach to inferential sta-
20 tistics is inappropriate. [...] In the case of a random sample, the p -value addresses the following question:
21 ‘If the null hypothesis were true of the population, how likely would we have been to obtain a sample
22 statistic this large or larger in a sample of this size?’ [...] In the case of random assignment, the p -value
23 targets the following question: ‘If the null hypothesis were true about the difference between treated and
24 untreated groups, how likely is it that we would have obtained a difference between them this big (or big-
25 ger) when studying treatment and comparison groups of this size?’ [...] If the experimental and control
26 groups have not been assigned using probability techniques, or if the cases have not been sampled from a
27 population using probability methods, inferential statistics are not applicable. They are routinely applied in
28 inapplicable situations, but an error is no less erroneous for being widespread.”

29 *(a) Fundamental prerequisites for using the p -value*

30 **Suggestion 1:** Do use neither p -values nor other inferential tools such as random errors or confidence
31 intervals if you have (a 100% sample of) the population of interest. In this case, no generalization from the
32 sample to the population (statistical inference) is necessary and you can directly describe the population
33 properties. Do not use p -values either if you simply provide descriptive statistics or if your findings hold
34 only for the sample under study – for example because it is a non-random sample that you have chosen for

1 convenience reasons instead of using probability methods. While not obvious at first view, studying a non-
2 random sample is a special case of studying a population because the sample already represents the entity
3 beyond which no statistical inference based on probability theory can be made.

4 **Suggestion 2:** Provide p -values if you deal with a random sample or a random assignment. But be clear
5 that the function of the p -value is different in the two cases. In the random sample case, you are concerned
6 with generalizing from the sample to the population. In the random assignment case, you are concerned
7 with the internal validity of an experiment in which you randomly assign experimental subjects to groups
8 that you subject to different treatments. For random assignments, the p -value is a continuous measure of
9 the strength of evidence against the null hypothesis of there being no treatment effect in the experiment. It
10 is *no* help whatsoever to assess the generalizability of results towards the population from which the ex-
11 perimental subjects themselves have been recruited. They may, or may not, be a random sample of a cer-
12 tain population.

13 **Suggestion 3:** Random samples from a population are often costly to come by and thence frequently not
14 available. When using p -values as a tool that is to help generalize from a sample to a population, provide
15 convincing arguments that your sample represents at least approximately a random sample. To avoid mis-
16 understandings, transparently state how and from which population the random sample was drawn and to
17 which population you want to generalize.¹

18 *(b) Wording guidelines for avoiding misunderstandings*

19 **Suggestion 4:** Use wordings that ensure that the p -value is understood as a *continuous* measure of the
20 strength of evidence against the null but *not* as the probability of the null (or the probability of being
21 wrong when rejecting the null). Avoid the term “error probability” because it suggests making this error.

22 **Suggestion 5:** Avoid wordings that insinuate that the p -value denotes an epistemic (posterior) probability
23 that you can attach to a scientific hypothesis (the null) *given* the evidence you found in your data.

24 **Suggestion 6:** Avoid wordings that insinuate that a low p -value indicates a large or even practically or
25 economically relevant size of the estimate, and vice versa.

26 **Suggestion 7:** Do not suggest that high p -values can be interpreted as an indication of no effect (“evidence
27 of absence”) even though in the NHST-approach “non-significance” leads to non-rejection of the null

¹ Emphasizing the p -value’s probabilistic foundation, DENTON (1988: 166f.) points out that “where there is a sample there must be a population.” He notes that conceiving of the population can be difficult. The easiest case is a sample drawn from a finite population such as a country’s citizens. A less intuitive sample-population relationship arises when we generate a sample by conducting an experiment such as flipping a coin n -times. Here, the population is an imaginary set of infinitely repeated coin flips. When studying observational macro-data, maintaining the p -value’s probabilistic foundation poses serious conceptual challenges. One would have to imagine an “unseen parent population” that is subjected to a noise producing process from which we observe a random realization.

1 hypothesis of no effect. Do not even suggest that high p -values can be interpreted as “absence of evi-
2 dence.” Doing so would negate the evident effects that you found in the data.

3 **Suggestion 8:** Avoid formulations and representations that could suggest that p -values below 0.05 (“sig-
4 nificant” results) can be interpreted as evidence in favor of the concrete coefficient estimates that you
5 happened to find in your study.

6 **Suggestion 9:** Do use neither the term “hypothesis *testing*” nor the term “*confirmatory* analysis.” It is
7 logically impossible to infer from the p -value whether the null hypothesis or an alternative hypothesis is
8 true. We can even not derive probabilities for hypotheses based on what has delusively become known as
9 “hypothesis *testing*.” p -values cannot “test” or “confirm” any hypothesis at all, but only describe data fre-
10 quencies under a certain statistical model including the null hypothesis.

11 **Suggestion 10:** Restrict the use of the word “evidence” to the findings in your data. Do not use “evidence”
12 for your inferential conclusions.

13 *(c) Things to do and discuss explicitly*

14 **Suggestion 11:** Do explicitly state whether your study is *exploratory* and thus aimed at generating new
15 research questions/hypotheses, or whether it is aimed at producing new evidence with regard to *pre-*
16 *specified* research questions/hypotheses. While the latter is conventionally termed “confirmatory analy-
17 sis,” this term should be avoided. It might mislead people to expect categorical yes/no answers that we
18 cannot give (see suggestion 9). Your paper may also contain both types of analysis. If so, explicitly com-
19 municate *where* you change from the study of pre-specified issues to exploratory search.

20 **Suggestion 12:** In the *exploratory* search for potentially interesting associations (e.g., in the control varia-
21 bles), p -values can be used as a flagging device to identify what might be worth investigating with new
22 data in the future. To prevent overhasty generalizations in this case, it might be worthwhile considering
23 BERRY’S (2017: 897) recommendation to use the following warning: “Our study is exploratory and we
24 make no claims for generalizability. Statistical calculations such as p -values and confidence intervals are
25 descriptive only and have no inferential content.”

26 **Suggestion 13:** If your study is aimed at producing evidence regarding *pre-specified* research ques-
27 tions/hypotheses, exactly report in your paper the list of questions/hypotheses that you drafted before run-
28 ning the analysis. In the results section, clearly relate findings to these initial questions or hypotheses.

29 **Suggestion 14:** When studying pre-specified questions or hypotheses, clearly distinguish two parts in your
30 analysis: (i) the description of the *evidence* (estimates) that you actually happened to find in your single
31 study (What is the evidence in this data?); (ii) the *inferential reasoning* that you base on this evidence
32 under consideration of p -values, confidence intervals, the study design, and all relevant external evidence
33 (What should one believe after seeing this data?). If applicable, a third part should outline the recommen-

1 dations or *decisions* that you would make all things considered (What should one do after seeing this da-
2 ta?).

3 **Suggestion 15:** Transparently report all analytical steps including data cleansing and the multiple models
4 that you tested in the process of model specification. When interpreting your findings, explicitly comment
5 on multiple comparisons that inflate the strength of the evidence against the null as indicated by the p -
6 value. Doing so, distinguish between (i) the multiple comparisons that you make because you study multi-
7 ple variables in your final regression model and (ii) the multiple comparisons that you make because you
8 tried multiple models before retaining one model as the “best” model. If appropriate, use robustness
9 checks to show how substantially stable (“robust”) your findings are over a reasonable range of analytical
10 variants including measurement and modeling alternatives.

11 **Suggestion 16:** In inferential reasoning, explicitly distinguish between statistical and scientific inference.
12 *Statistical inference* and the p -value are concerned with the random sampling error, i.e., the fact that even
13 a random sample may not exactly reflect the properties of the population. Generalizing from a random
14 sample to its population is only the first step of *scientific inference*, which is the totality of reasoned judg-
15 ments (inductive generalizations) that we make in the light of our own study and the available body of
16 external evidence. We might want to know, for example, what we can learn from a random sample of a
17 country’s agricultural students for its student population, its citizens, or even human beings in general. Be
18 clear in your inferential reasoning that a p -value, being a probabilistic concept, can do *nothing* to assess
19 the generalizability of results beyond the parent population (here: the country’s agricultural students) from
20 which the random sample has been drawn.

21 *(d) Operative rules*

22 **Suggestion 17:** Provide information regarding the size of your estimate (point estimate). In many regres-
23 sion models, a meaningful representation of magnitudes will require going beyond coefficient estimates
24 and displaying marginal effects.

25 **Suggestion 18:** Do not use asterisks (or the like) to denote different levels of “statistical significance.”
26 Doing so could instigate erroneous categorical reasoning.

27 **Suggestion 19:** Provide exact p -values for all coefficient estimates or marginal effects and avoid the clas-
28 sification of results in being “statistically significant” as opposed to “statistically non-significant.” That
29 said, avoid using the terms “statistically significant” and “statistically non-significant” altogether. Dis-
30 pensing with these two categorical labels enables you for the first time to use “relevant” and “significant”
31 as interchangeable terms without causing confusion.

1 **Suggestion 20:** Provide standard errors for all coefficient estimates or marginal effects. Additionally pro-
2 vide confidence intervals for the focal variables of interest associated with your pre-specified research
3 questions/hypotheses.

4 We hope that we give an impulse to applied researchers to comply with these single-study suggestions in
5 all cases in which coordinated multiple-study designs are not feasible. We furthermore believe that the
6 quality of published research could be improved by incorporating these suggestions as best practice rules
7 into journal guidelines. With a view to the inflation of the strength of evidence through covert multiple
8 comparisons (*p*-hacking), SIMMONS et al. (2012) propose a formalization outside of the paper that seems
9 worthwhile considering. Similar to the standard “no-competing-interests” statements, they suggest to
10 oblige authors to make a *formal* “no-*p*-hacking” declaration. The problem is that it is difficult to unambig-
11 uously define the practices that are outlawed as *p*-hacking. A substantiated selection of an analytical ap-
12 proach is not *p*-hacking. But results will be biased if researchers covertly engage in multiple comparisons
13 and selectively publish those analytical variants that “work” in that they produce lower *p*-values than other
14 variants (HIRSCHAUER et al. 2016). For a formal declaration to make sense, journals must clearly specify
15 outlawed practices. Some people may think that, given the perverse publish-or-perish conditions that
16 many researchers face today, a formal no-*p*-hacking declaration is just an empty phrase. However, we
17 believe that it could produce a practically significant reinvigoration of science ethics’ call for transparency
18 and integrity that leads to published research better reflecting reality than what we have seen in the past. In
19 this sense, we agree with SIMMONS et al. (2012: 6) that “changes need not to be judged in terms of their
20 perfection, but merely in terms of their improvement.”

21 While we believe that our suggestions represent practically significant steps towards improvement, we do
22 not expect that all researchers will endorse all of them at once. With a view to their acceptance and imme-
23 diate viability, there seem to be three categories: some suggestions, such as the eschewal of asterisks and
24 the requirement to display random errors, are likely to cause little controversy. Others, such as renouncing
25 dichotomous significance declarations and giving up the term “statistical significance” altogether, will
26 possibly be questioned – at least by some scientists. And two suggestions, both concerned with leaving
27 behind categorical yes/no declarations, require more than just debate before they can be implemented.
28 They require solving problems that linger on because the existing statistical literature does not yet meet
29 applied researchers’ needs for guidance on how to do things in a post $p < 0.05$ era.

30 First, we do not yet have formulations at our disposition that ensure that the *p*-value is understood as a
31 *continuous* measure of the strength of evidence against the null. Which wording is appropriate to convey
32 the information of a *p*-value of, let’s say, 0.37 as opposed to 0.12 or 0.06 or 0.005 – for large and small
33 effects, respectively? Our troubles do not come as a surprise since the difficulties of translating the *p*-value
34 concept adequately into natural language are at the heart of the problem. BERRY (2017: 896) puts it in a

1 nutshell: “I forgive nonstatisticians who cannot provide a correct interpretation of $p < 0.05$. p -Values are
2 fundamentally un-understandable. I cannot forgive statisticians who give understandable—and therefore
3 wrong—definitions of p -values to their nonstatistician colleagues. But I have some sympathy for their
4 tack. If they provide a correct definition then they will end up having to disagree with an unending se-
5 quence of ‘in other words’. And the colleague will come away confused [...]” While this statement may
6 seem overly pessimistic, we agree with the problem description. The only way out is to *find* and *agree* on
7 formulations that convey the limited but existing informational content of the p -value in both a *correct*
8 and *meaningful* way, lest we better abandon its use altogether. This is what BERRY (2017): demands: “We
9 created a monster [the p -value]. And we keep feeding it, hoping that it will stop doing bad things. It is a
10 forlorn hope. No cage can confine this monster. The only reasonable route forward is to kill it.” Contrary
11 to Berry, we believe that we should only dispense with significance testing and categorical reasoning but
12 retain the p -value itself because of its familiarity and potential usefulness. However, if retaining the p -
13 value is to make sense, we need an organized approach – maybe under the aegis of the ASA – that gets
14 some work done and comes up with practical guidance for applied researchers who, rightly leaving behind
15 dichotomous significance declarations, are in need of correct and understandable formulations of what a p -
16 value means.

17 A second and analogous problem arises because it is equally hard to provide a correct wording and good
18 intuition regarding the meaning of the confidence interval (CI). Imagine observing a mean difference of
19 10 g in daily weight gains between two randomly assigned animal groups that were subjected to different
20 dietary treatments. Finding a 95% CI of [8, 12], it would not be correct to say that the difference is be-
21 tween 8 g and 12 g with 95% probability. Not much change for the better is obtained by replacing “proba-
22 ble” by words such as “likely/plausible” or “confident.” Stating, for example, that the CI indicates the
23 precision of the estimation and that “in other words, we can be 95% confident that the difference is be-
24 tween 8 g and 12 g” is extremely deceptive. This holds even though researchers could argue that they use
25 the word “confidence” as a technical term that by convention is attached to the said interval. Even though
26 such statements sound like uncertainty statements, they promise too much certainty. In the words of
27 GELMAN (2016), they are to be qualified as “uncertainty laundering” because they neglect the inherent
28 uncertainty of the CI itself. A correct interpretation requires realizing that CI (analogous to p -values) vary
29 from one random sample to the other. A 95% CI only means that 95% of CI computed for repeatedly
30 drawn random samples will capture the “true” value (GREENLAND et al. 2016). As in the case of p -values,
31 we must realize that providing a correct technical definition of a difficult-to-understand concept is not
32 enough. It is likely to provoke an unending sequence of false “in-other-words statements.” In fact, we
33 rarely see a proper and understandable interpretation of CI in empirical papers and even textbooks. Most
34 formulations seem to communicate in one or the other way that a CI describes the probability that the
35 specified range contains the “true” value with the probability of 95% or 99%. They thus insinuate that we

1 could make an epistemic probability statement regarding the population effect size based on the results of
2 a single study. Such a statement must be reserved to Bayesian analysis (here: the Bayesian credibility in-
3 terval), however. The lack of appropriate wordings is especially serious since CI have been recommended
4 as being a part of the solution for the p -value problem (e.g., CUMMING 2014). But to be a part of the solu-
5 tion, we first need guidance and consensus regarding the wordings that are able to communicate the in-
6 formational content of a CI not only *correctly* but also *meaningfully*. Being again a problem that is not
7 restricted to a particular applied science discipline, the process of finding a solution could possibly be
8 initiated and organized by leading statistical associations.

9 **4 Reforms under way and reforms still to be undertaken in publishing economic research**

10 *Measures addressed at the author(s) of a single study*

11 Both the accumulation of knowledge and technological developments in computing continuously shift
12 what constitutes best methodological practice in statistical analysis. Given these dynamics, sticking tradi-
13 tions as well as rigid formalizations such as tight but rarely scrutinized journal guidelines may slow down
14 or even prevent necessary change. A particular challenge arises for interdisciplinary journals which need
15 sufficiently flexible rules to cater for the fact that they receive manuscripts from collaborative research
16 teams and authors with different traditions in statistical analysis. While overly rigid and sticking rules are
17 detrimental with respect to interdisciplinary research and dynamic adjustments, guidelines can also be a
18 pertinent means to communicate new best practice procedures and induce overdue change in disciplinary
19 traditions and researchers' inert habits.

20 Trying to get an impression of “what is going on” in the practice of publishing research, we asked the
21 editors of 100 leading economic journals (cf., <http://www.scimagojr.com/journalrank.php?category=2002>)
22 about policy changes with respect to the use of p -values. Overall, journals seem still a long way from
23 translating a significant portion of recent reform suggestions into concrete journal policies. Despite the
24 prominence of the current p -value debate, a significant share of journal editors believe that their reviewing
25 system is sufficiently effective to prevent inferential errors. Consequently, they do not see a need to bring
26 about formal change. The editors of some journals, however, seem to be seriously worried about the mis-
27 uses and misinterpretations of the p -value. What is more, some editorial boards are deliberating concrete
28 future steps or have already started using their guidelines to prevent misleading practices and false infer-
29 ential conclusions. For example, leading journals, such as the American Economic Review, Econometrica,
30 and the four AEJs (Applied Economics, Economic Policy, Macroeconomics, Microeconomics), now re-
31 quest authors not to use asterisks or other symbols to denote “statistical significance.” While this seems a
32 small change, it represents a distinct breach of a convention that many researchers considered to be set in
33 stone for a long time. The basic idea behind banning asterisks is to prevent overconfident categorical con-

1 clusions induced by arbitrary thresholds. Other noteworthy editorial policy changes in leading economic
2 journals (e.g., American Economic Review, Econometrica) include the request to explicitly report effect
3 sizes (marginal effects) and display standard errors (or even confidence intervals) in results tables. With
4 respect to the p -value, they call upon authors to display standard errors *instead of* p -values in results tables
5 but do not ban the use of p -values when interpreting results in the text. This is consistent with a suggestion
6 put forward by many critical voices in the recent debate, namely to demote p -values from their pedestal
7 and consider them as a tool amongst many that may help researchers make appropriate inferences (cf.,
8 e.g., AMRHEIN et al. 2017; MACSHANE et al. 2017; TRAFIMOW et al. 2017).

9 The fact that some of the leading journals have initiated modest but sensible changes with regard to the
10 use of the p -value is a promising signal. GOODMAN (2017: 559) notes that “norms are established within
11 communities partly through methodological mimicry.” If a field’s flagship journals, opinion leaders, and
12 professional associations take up the lead, they may be able to set a trend. “Once the process starts, it
13 could be self-enforcing. Scientists will follow practices they see in publications; peer reviewers will de-
14 mand what other reviewers demand of them.” Besides their general function as beacons for best practice,
15 the guidelines of flagship journals may become more direct agents of change due to prevalent submission
16 practices. Researchers often submit their papers to leading journals first. If declined, they regularly try
17 alternative publication outlets and submit their papers, written according to the guidelines of the flagship
18 journal, more or less unchanged to less prominent journals. This might cause a trickle-down effect that
19 generates new best practice standards for the less prominent journal.

20 *Measures beyond the single study*

21 Avoiding mistakes within the single study is a necessary but not sufficient condition for making correct
22 inferences. Instead, we need to embed each study in its wider context and consider the body of evidence
23 including all external (“prior”) knowledge in the field under research. Several propositions towards im-
24 provement beyond the realms of the single study have been made. Meta-analysis that systematically con-
25 solidates the body of evidence and the formal consideration of prior knowledge through Bayesian analysis
26 are prominent examples. Furthermore, two reforms on the level of scientific institutions (journals, scien-
27 tific associations) are practically important in some disciplines but only nascent or non-existing in others.

28 First, many leading journals now oblige authors to provide their raw data and analytical protocols in the
29 appendix of their paper to facilitate *replication* studies aimed at scrutinizing a study’s findings. While
30 compulsory sharing of raw data and analytical protocols seems to slowly trickle down to more and more
31 journals, institutionalized efforts to strengthen replication are weak in economics compared to other fields.
32 According to DUVENDACK et al. (2015), most of the 333 economic Web-of-Science journals still give low
33 priority to replication. The same holds for initiatives targeted at counteracting publication bias. While a
34 global initiative [All Trials Registered/All Results Reported](#) was launched in 2013 in the medical sciences,

1 for example, similar efforts are rare in economics. Among the few exceptions are [The Replication Net-](#)
2 [work](#), and [Replication in Economics](#). Both platforms are aimed at fostering the scrutiny of scientific
3 claims and at counteracting publication bias by providing not only databases for replications but also equal
4 opportunities for publishing positive and negative results.

5 A second important reform on the institutional level is *pre-registration*. Pre-registration goes not only
6 beyond the post-study provision of raw data and analytical protocols but also beyond the mere appeal to
7 honestly report *pre-specified* analysis plans in the paper. Instead, it obliges researchers to disclose their
8 hypotheses, data, and analytical approach *before* running the analysis. Eventually necessary deviations
9 from the pre-analysis plan must be justified in the paper. Pre-registration is aimed at preventing *p*-hacking
10 and providing equal chances of being published independent of which results are eventually found. In
11 other words, it is to prevent not only multiple comparisons and selective reporting but also selective pub-
12 lishing and thus the bias towards positive (“statistically significant”) findings (cf., ROSENTHAL 1979),
13 which seems to be widespread even in economic flagship journals (BRODEUR et al. 2016). Contrary to
14 clinical drug trials for which pre-registration is standard (<http://www.who.int/ictrp/network/primary/en/>),
15 it is still rare in the social sciences. There are, however, two new initiatives. Within the “\$1 Million Pre-
16 registration Challenge,” the Center for Open Science (<https://cos.io/prereg/>) provides \$1,000 to 1,000 re-
17 searchers who pre-register their research projects. Because existing registries were not considered a good
18 fit for the needs of the social sciences, the American Economic Association launched an initiative in 2017
19 to register randomized controlled trials on its AEA RCT platform (<https://www.socialscienceregistry.org/>).
20 After peer-approval of the study design and analysis plan, research projects are accepted and published
21 before they are implemented.

22 The poor development of replication, pre-registration, and meta-analysis in economics stems from the
23 discipline’s culture and its preoccupation with observational study, data-driven model specification, and
24 multiple regression analysis. In other words, there are questions to be answered before approaches from
25 other fields can be transplanted to (non-experimental) economic research. It is not clear, for example, how
26 replication or preregistration should work within a disciplinary culture in which it is not only common but
27 also highly appreciated practice to specify regression models *after* seeing the data (“model fitting”). Relat-
28 ed to that, the question arises of how to carry out quantitative meta-analysis and consolidate the body of
29 evidence when even within a narrow field of research there are often as many data-dependent model spec-
30 ifications as studies. The fact that economic research is mainly a bottom-up research exercise is responsi-
31 ble for the lacking comparability across studies. Non-programmed bottom-up research produces a large
32 quantity of empirical results on topical issues, but is plagued by an enormous heterogeneity of empirical
33 measures and model specifications. Besides differing measures for the focal variables of interest,
34 databased models are regularly populated by differing interaction terms, transformed variables, lagged

1 variables, higher-order polynomials, and control variables. Given the heterogeneity of econometric mod-
2 els, applied economists need guidance and consensus regarding best practices for specification search,
3 replication, and meta-analysis. We hope that professional associations in the field of economics take up
4 the cause and organize a debate on the urging question of how to systematically build up knowledge and
5 scrutinize scientific claims derived from nearly limitless variants of databased models.

6 References

- 7 Amrhein, V., Korner-Nievergelt, F., Roth, T. (2017): The earth is flat ($p > 0.05$): significance thresholds
8 and the crisis of unreplicable research. PeerJ, DOI 10.7717/peerj.3544.
- 9 Benjamini, Y. (2016): It's not the p -values' fault. The American Statistician 70 (2): Supplemental Material
10 to the ASA Statement on P-Values and Statistical Significance
11 ([http://amstat.tandfonline.com/doi/suppl/10.1080/00031305.2016.1154108/suppl_file/utas_a_1154108_sm](http://amstat.tandfonline.com/doi/suppl/10.1080/00031305.2016.1154108/suppl_file/utas_a_1154108_sm5354.pdf)
12 [5354.pdf](http://amstat.tandfonline.com/doi/suppl/10.1080/00031305.2016.1154108/suppl_file/utas_a_1154108_sm5354.pdf)).
- 13 Berry, D. (2017): A p-Value to Die For. Journal of the American Statistical Association 112(519): 895-
14 897 (DOI: 10.1080/01621459.2017.1316279).
- 15 Brodeur, A., Lé, M., Sangnier, M., Zylberberg, Y. (2016): Star Wars: The Empirics Strike Back. Ameri-
16 can Economic Journal: Applied Economics 8(1): 1-32.
- 17 Cumming, G. (2014): The New Statistics: Why and How. Psychological Science 25(1): 7-29.
- 18 Denton, F.T. (1988): The significance of significance: Rhetorical aspects of statistical hypothesis testing
19 in economics. In: Klamer, A., McCloskey, D.N., Solow, R.M. (eds.): The consequences of economic rhet-
20 oric. Cambridge: Cambridge University Press: 163-193.
- 21 Duvendack, M., Palmer-Jones, R., Reed, W.R. (2015): Replications in Economics: A Progress Report.
22 Econ Journal Watch 12(2): 164-191.
- 23 Goodman, S. (2008): A dirty dozen: Twelve p-value Misconceptions. Seminars in Hematology 45: 135-
24 140.
- 25 Cohen, J. (1994): The earth is round ($p < 0.05$). American Psychologist 49(12): 997-1003.
- 26 Forstmeier, W., Wagenmakers, E.-J., Parker, T.H. (2016): Detecting and avoiding likely false-positive
27 findings – A practical guide. Biological Reviews of the Cambridge Philosophical Society, 92(4): 1941-
28 1968 (doi:10.1111/brv.12315).
- 29 Gelman, A. (2016): The problems with p-values are not just with p-values. American Statistician, supple-
30 mental material to the ASA statement on p-values and statistical significance, 10, 2016.
- 31 Gelman, A., Carlin, J. (2017): Some natural solutions to the p-value communication problem-and why
32 they won't work. Blogsite: Statistical Modeling, Causal Inference, and Social Science.
- 33 Greenland, S. (2017): Invited Commentary: the Need for Cognitive Science in Methodology. American
34 Journal of Epidemiology 186(6): 639-645.
- 35 Greenland, S., Senn, S.J., Rothman, K.J., Carlin, J.B., Poole, C., Goodman, S.N., Altman, D.G. (2016):
36 Statistical tests, P values, confidence intervals, and power: a guide to misinterpretations. European Journal
37 of Epidemiology 31(4): 337-350.
- 38 Goodman, S.N. (2017): Change norms from within. Nature 551: 559 ([https://www.nature.com/magazine-](https://www.nature.com/magazine-assets/d41586-017-07522-z/d41586-017-07522-z.pdf)
39 [assets/d41586-017-07522-z/d41586-017-07522-z.pdf](https://www.nature.com/magazine-assets/d41586-017-07522-z/d41586-017-07522-z.pdf)).
- 40 Halsey, L.G., Curran-Everett, D., Vowler, S.L., Drummond, B. (2015): The fickle P value generates irre-
41 producible results. Nature Methods 12(3): 179-185.

- 1 Hirschauer, N., Grüner, S., Mußhoff, O., Becker, C. (2017): Pitfalls of significance testing and p-value
2 variability – implications for statistical inference. The Centre for Statistics working paper series 07/2017,
3 Georg August University Goettingen (<https://www.uni-goettingen.de/de//511092.html>).
- 4 Matthews, R. (2017): The ASA's *p*-value statement, one year on. The Royal Statistical Society. Signifi-
5 cance 14(2): 38-40.
- 6 McShane, B., Gal, D. (2017): Statistical Significance and the Dichotomization of Evidence. Journal of the
7 American Statistical Association 112(519): 885-895 (DOI: 10.1080/01621459.2017.1289846).
- 8 McShane, B., Gal, D., Gelman, A., Robert, C., Tackett, J.L. (2017): Abandon Statistical Significance.
9 <http://www.stat.columbia.edu/~gelman/research/unpublished/abandon.pdf>
- 10 Rosenthal, R. (1979): The file drawer problem and tolerance for null results. Psychological Bulletin 86(3):
11 638-641.
- 12 Simmons J.P., Nelson L.D., Simonsohn U. (2012): A 21 word solution. Dialogue. The Official Newsletter
13 of the Society for Personality and Social Psychology 26(2):4-7.
- 14 Trafimow, D. et al. (2017): Manipulating the alpha level cannot cure significance testing – comments on
15 “Redefine statistical significance”. PeerJ Preprints 5:e3411v1
16 (<https://doi.org/10.7287/peerj.preprints.3411v1>).
- 17 Wasserstein, R.L., Lazar N.A. (2016): The ASA's statement on p-values: context, process, and purpose,
18 The American Statistician 70(2): 129-133.