# Pitfalls of significance testing and p-value variability – implications for statistical inference

**Norbert Hirschauer and Sven Grüner**
Institute for Agricultural and Nutritional Sciences
Agribusiness Management
Martin Luther University Halle-Wittenberg
Halle (Saale)


**Oliver Mußhoff**
Department of Agricultural Economics and Rural Development
Farm Management Group
Georg-August University Goettingen
Göttingen


**Claudia Becker**
Faculty of Law, Economics, and Business
School of Economics and Business
Statistics
Martin Luther University Halle-Wittenberg
Halle (Saale)

This version: July 10, 2017

# Pitfalls of significance testing and $p$-value variability – implications for statistical inference

*Norbert Hirschauer (corresponding author)*

Martin Luther University Halle-Wittenberg, norbert.hirschauer@landw.uni-halle.de


*Sven Grüner*

Martin Luther University Halle-Wittenberg


*Oliver Mußhoff*

Georg-August-University Goettingen


*Claudia Becker*

Martin Luther University Halle-Wittenberg

**Abstract:** Data on how many scientific findings are reproducible are generally bleak and a wealth of papers have warned against false findings in recent years. This paper discusses the question of what we can(not) learn from the $p$-value, which is still widely considered as the gold standard of statistical validity. We describe classical pitfalls of significance testing and the disregard of $p$-value sample-to-sample variability.

**JEL:** B41; C18


**Key Words:** meta-analysis; multiple testing; $p$-hacking; publication bias; $p$-value misinterpretations; $p$-value sample-to-sample variability; statistical inference; statistical significance

# 1   Introduction

Data on how many scientific findings are reproducible are generally bleak and a wealth of papers have warned against false findings in recent years. The reasons for false discoveries are manifold, but misuses and misinterpretations of statistical significance testing based on *p*-values are the most prominently decried ones. In the light of prevalent and persistent misunderstandings, even the American Statistical Association (ASA) felt compelled to issue a warning that the *p*-value can neither be used to determine whether a scientific hypothesis is true nor whether a finding is important (Wasserstein and Lazar 2016). In a *Nature* paper, Baker (2016: 151) comments: "This is the first time that the 177-year-old ASA has made explicit recommendations on such a foundational matter in statistics, says executive director Ron Wasserstein. The society's members had become increasingly concerned that the *P* value was being misapplied in ways that cast doubt on statistics generally, he adds." The ASA apparently considers inappropriate interpretations and uses of significance tests so serious a threat to statistics and science in general that, as a follow-up to its statement, it organizes a symposium under the heading *Scientific Method for the 21st Century: A World beyond p < 0.05* in fall 2017.

The literature concerned with the "reproducibility crisis" can be largely divided into four classes. One group of papers focus on multiple testing that is left uncorrected for in many studies. This leads to inflated claims of statistical significance. Even though multiple testing is often *evident*, researchers in the economic and social sciences usually disregard adjustment requirements. An example are multiple regressions where many scientists ignore the problem due to the "conventionality" of the model. Imagine a set of 10 regressors all of which are independent and none of which is predictive. Despite the completely random probabilistic structure, there is a 40% chance ($= 1 - 0{,}95^{10}$) of finding at least one statistically significant coefficient at the conventional threshold of 0.05 (Altman and Krzywinski 2017). Even more disastrous is *covert* multiple testing in conjunction with selective reporting. Testing alternative data, hypotheses, or analytical variants, and reporting only what has produced low *p*-values has been coined "*p*-hacking." Simmons et al. (2011: 1359) note that it is common "to explore various analytic alternatives, to search for a combination that yields 'statistical significance', and to then report only what 'worked'. The problem, of course, is that the likelihood of at least one (of many) analyses producing a falsely positive finding at the 5% level is necessarily greater than 5%."

A second class of papers tackle the semantically induced misunderstandings (cognitive biases) that the delusive language of frequentist statistics causes even among non-statistician scientists. This group of papers stress the limitations of statistical significance testing that are present even if there are no multiple testing problems. They emphasize that the frequentist *p*-value concept can do much less to inform us about the reliability of scientific findings than what the persistently recurring colloquial associations with statistical terms such as significance, error probability, confidence interval, and inference suggest. In the words of Greenland et al. (2016: 337), the departure point of these papers can be characterized as follows: "Misinterpretation and abuse of statistical tests, […] have been decried for decades, yet remain rampant. A key problem is that there are no interpretations of these concepts that are at once simple, intuitive, correct, and foolproof. Instead, correct use and interpretation of these statistics requires an attention to detail which seems to tax the patience of working scientists. This high cognitive demand has led to an epidemic of shortcut definitions and interpretations that are simply wrong, sometimes disastrously so—and yet these misinterpretations dominate much of the scientific literature."

A third class of papers is concerned with the question of how to summarize the findings of individual studies to obtain an appropriate picture of the state of knowledge in a given field (meta-analysis). Borenstein et al. (2009: xxi) state that "rather than looking at any study in isolation, we need to look at the body of evidence." The essence of meta-analysis is best explained by comparing it to narrative reviews that simply count the studies that were declared

statistically significant, or not, at the arbitrary threshold of 0.05. Contrasting the tallies ("vote counting"), together with the mistaken belief that studies with *p*-values on opposite sides of the conventional level of 0.05 are conflicting (Goodman 2008), leads to a wrong picture of the body of evidence. Borenstein et al. (2009: 14) claim that doing narrative reviews boils down to "*doing arithmetic with words*" and that "when the words are based on *p*-values *the words are the wrong words*." They contend that this problem practically "gallops" through many research fields. In contrast, meta-analysis addresses the question of how to arithmetically synthetize meta summary statistics (usually the meta effect size and its *p*-value and confidence interval) based on the statistics provided in individual studies.

Looking at the publication bias (Smith 1980) or file drawer problem, a fourth class of papers focus on researchers' incentives and the perverse effects of the scientific publishing system with its pressure to "publish (statistically significant results) or perish." Even if there were no inflation from multiple testing and no misinterpretation of statistical significance tests, the file drawer effect would distort the body of evidence towards results that can be declared "statistically significant." Starting with Sterling (1959) as early precursor, researchers have increasingly realized this bias in recent years. Even back in the 70s, Rosenthal (1979: 638) vividly described the harmful consequences that result from the preferences of researchers, reviewers and publishers for "significant" novelties: "The extreme view of this problem, the 'file drawer problem,' is that the journals are filled with the 5% of the studies that show Type I errors, while the file drawers back at the lab are filled with the 95% of the studies that show nonsignificant (e.g., p > .05) results." The file drawer effect does not only lead to the publication of unsubstantiated claims that are too rarely subjected to independent scrutiny. It may also lead to wrong medical treatment or policy recommendations with dire practical consequences.

Table 1: Classical pitfalls of significance testing

|  | Flaws when performing significance tests | Mistakes when interpreting significance tests |
|---|---|---|
| **Within study** | **(1) Uncorrected multiple testing** (= inflated claims of statistical significance) <br> - Unintentional disregard of evident multiple testing <br> - *p-hacking*: covert testing of multiple analytical variants and selective reporting of those that yielded "statistical significance" | **(2) Semantically induced misinterpretations** <br> - *Inverse probability error* (interpreting the *p*-value as the probability of the null) <br> - *Sizeless stare* or even equation of significant effects with large or important effects <br> - *False dichotomy* (interpreting not statistically significant results as confirmation of the null) |
| **Across studies** | **(3) Exaggerated focus on one-shot studies** (= disregard of prior knowledge) <br> - Lacking meta-analysis <br> - Improper meta-analysis (*vote counting*) <br> - Lacking Bayesian analysis | **(4) Publication bias/file drawer effect** (= distortion towards positive results) <br> - Selective preparation and submission <br> - Selective reporting (*p-hacking*) <br> - Selective publishing |

Table 1 lists the pitfalls marked up above. While we will describe them one by one, it should be noted that they are intimately linked: first, misuses and misinterpretations often occur consecutively and may reinforce each other and accumulate. Second, mistakes and distortions in the research process may render subsequent procedures useless even if they are appropriate as such. For example, even the most elaborate meta-analysis aimed at consolidating the available evidence from prior studies will yield nonsense results in the presence of serious publication bias since it simply consolidates the distortion (Kline 2013: chapter 9).

We have attached the label "classical" to the pitfalls in Table 1 because they have already been widely discussed in the past. Despite the large body of literature that has accumulated over the last decades on these issues (OAKES 1986; Cohen 1994; Nickerson 2000; Ioannidis 2005; Armstrong 2007; Simmons et al. 2011; Motulsky 2014; Hirschauer et al. 2016; Wasserstein and Lazar 2016, and many others), misapplications of statistical significance testing continue to be alarmingly "normal" practice for many scientists. With few exceptions (e.g., Krämer 2011; Ziliak and McCloskey 2008), the acknowledgement of the critical issues in

statistical significance testing seems to be particularly low in economics. This can be partly ascribed to the fact that papers concerned with the pitfalls of significance declarations are not only scattered over disciplines often remote from economics but also concerned with analytical measures such as mean differences or response ratios in many instances. Hence, an easily accessible informational resource for economists, who heavily rely on regression analysis, is lacking. Furthermore, the *p*-value's sample-to-sample variability, even though it is a fundamental feature that severely limits its suitability to indicate the strength of evidence, has been underexposed in the economic literature so far (Berry 2016; Halsey et al. 2015).

Against this background, we provide a systematic overview of the "classical" pitfalls of statistical significance testing in section 2. To complete the understanding that the *p*-value is a descriptive summary of a given data set (sample) but a poor tool for making inferences and generalizations, section 3 focuses in detail on the *p*-value's sample-to-sample variability. Section 4 concludes with a brief discussion of questions that need to be answered to arrive at an informed approach to statistical inference in multiple regression analysis.

## 2 Classical pitfalls in significance testing

### 2.1 Uncorrected multiple testing

*Unintentional disregard of evident multiple testing*

Quantitative research in the social sciences including economics is frequently based on multiple regression. While often ignored, multiple testing is inherent to multiple regression since we test as many null hypotheses as we have variables of interest. This inflates claims of statistical significance. Let us take a closer look why. Adopting the frequentist null hypothesis testing assumption of no association, significance testing of multiple regression coefficients represents a Bernoulli experiment: we independently repeat, for each regressor, a trial with two possible outcomes "significant" $S$ and "not significant" $\bar{S}$ that occur with probabilities $P(S) = 0.05 = p$, and $P(\bar{S}) = 0.95 = 1 - p$. The chance of finding "significant" coefficients despite a fully random data structure is provided by the binomial distribution $B_{m,p}(k)$, with $m$ indicating the number of regressors and thus null hypotheses, $p = 0.05$ the probability of falsely claiming significance, and $k \in \{0, \dots, m\}$ the number of mistaken significance declarations.

Table 2: Probability of $k$ false significance declarations at the 0.05 threshold in multiple regressions with completely random data structure

| | Number of multiple tests (null hypotheses) | | | | | |
|---|---|---|---|---|---|---|
| | $m = 1$ | $m = 2$ | $m = 3$ | $m = 5$ | $m = 10$ | $m = 20$ |
| $B_{m,0.05}(0) = P(k = 0)$ | 0.95 | 0.903 | 0.857 | 0.774 | 0.599 | 0.358 |
| $B_{m,0.05}(1) = P(k = 1)$ | 0.05 | 0.095 | 0.135 | 0.204 | 0.315 | 0.377 |
| $B_{m,0.05}(2) = P(k = 2)$ | - | 0.003 | 0.007 | 0.021 | 0.075 | 0.189 |
| $P(k \geq 1) = 1 - 0.95^m$ | **0.05** | **0.098** | **0.143** | **0.226** | **0.401** | **0.642** |
| Inflation factor: $P(k \geq 1)/0.05$ | 1.000 | 1.950 | 2.853 | 4.524 | 8.025 | 12.830 |

Table 2 illustrates how multiple tests inflate claims of statistical significance. The implications are woeful if adjustment requirements are ignored. In this case, a researcher might be convinced of having found several "significant" results at the conventional level of 0.05 without realizing that (s)he faces an inflated probability of 0.098 (0.143, 0.226, 0.401, 0.642) of finding at least one significant coefficient in a multiple regression with 2 (3, 5, 10, 20) regressors even though there is no association whatsoever in the data. Contrary to *multiple* testing where conventional *p*-values do not reflect the probability under the null, there is of course a 0.05 probability of falsely rejecting the null in a *single* test when it is valid.

Multiplicity problems and the subsequent profusion of faulty scientific claims are not restricted to regression analysis in observational studies. Experimental studies regularly analyze how various treatments (multiplicity of treatments) affect various outcomes (multiplicity of effects) across various subgroups (multiplicity of groups). While it is common to study considerable multiplicities of treatments, outcomes, and subgroups in economic experiments, it is still, "with a few exceptions, […] uncommon for the analyses of these data to account for the multiple hypothesis testing (List et al. 2016: 1).

*p-hacking: covert testing of multiple analytical variants and selective reporting*

The arbitrary dichotomization of results into "significant" and "not significant," in conjunction with researchers' self-interested desire to obtain findings that can be declared "significant" is considered a major cause of the reproducibility crisis. The term "*p*-hacking" has been coined to describe the behavior of researchers who try a multiplicity of analytical alternatives and then report only the one that produced the desired result (Simmons et al. 2011). *p*-hacking has several noteworthy features: first, the list of possible analytical alternatives is near endless in most research contexts. Second, problem awareness among researchers is frequently low.[1] Third, the multiple analytical variants that can be tried often seem to have little in common at first view. They share one noxious quality, however: selective reporting of covert multiple tests may disastrously inflate claims of statistical significance.

The literature concerned with the pitfalls of significance tests and declarations has extensively discussed the various forms of *p*-hacking (see Hirschauer et al. 2016 for an overview). We do not intend to summarize this literature once more. However, to facilitate the understanding in which ways *p*-hacking represents a covert and thence especially harmful type of multiple testing, we briefly describe the four main forms of *p*-hacking that can be distinguished.

a) **Testing multiple data sets:** Researchers might be tempted to explore whether *p*-values can be reduced when some data are removed – for example on the seemingly justifiable grounds of being outliers. They might also try which *p*-values they can obtain when analyzing multiple data subsets. Testing a treatment in 20 arbitrary subgroups, for example, is analogous to testing 20 null hypotheses and leads to false positive probabilities as shown in Table 2 (last column) even if the effect is nil in all subgroups (Kerr 1998; Motulsky 2014). Finally, researchers might ad hoc increase sample size after an original sample has yielded "disappointing" *p*-values. A general feeling that larger samples are better may contribute to a lacking awareness that this constitutes a test of multiple sample sizes.

b) **Testing multiple data transformations:** Researchers might also be tempted to check which (combination) of many conceivable data transformations produces lower *p*-values than the original data. Possibilities are nearly unlimited: downgrading of measurement scales (e.g., age classes instead of age in years), log-transformation, squaring, normalization (ratios), and the synthetization of various variables including interaction terms. Any of these manipulations may be statistically appropriate in the light of the theory and the research questions. They inflate claims of statistical significance, however, if they are driven by a researcher's significance-pursuing behavior.

c) **Testing multiple variable sets:** The choice of variables to be included in a regression is often ambiguous: Which theoretical (latent) constructs are to be included into the analysis? Which manifest variables (e.g., survey items) should be used to measure these constructs? Which control variables should be considered? Mining for a combination of variables that yields low *p*-values inflates statistical significance claims. A researcher who

---

[1] Simmons et al. (2011) contend that *p*-hacking, while violating the rules of good scientific practice, is often not caused by malicious intent but by analytical ambiguities that facilitate the unconscious self-justification of choices that mesh with researchers' desire to obtain low *p*-values.

studies, for example, how people's attitudes towards organic farming affect their willingness to pay for organic products produces a distortion if (s)he keeps on searching for a survey item for "attitude" until (s)he finds one that produces a "significant" result.

**d) Testing multiple estimation models:** Selecting statistical tests and models also offers ample scope for decisions that "improve" $p$-values. For example, when facing an ambiguous choice of whether to use a simple OLS estimator or a panel data model, it would be good scientific practice to transparently compare the results of both models. However, the rules of good scientific practice are occasionally broken and data analyses are not performed as planned in a prior study design but ad hoc adjusted according to the criterion of which analytical model yields low $p$-values. Scientific transparency is lost when the results of competing models are neither explicitly reported nor comparatively discussed.

Disregarding multiplicities and notably $p$-hacking are problems that arise when confirmatory and exploratory study are conceptually mixed up. Exploring potentially interesting associations within a data set is an adequate primary step of the research process. In exploratory study, $p$-values may help identify what might be worth investigating with new data in the future (Head et al. 2015; Motulsky 2014).[2] This exploratory exercise, however, must be clearly distinguished from confirmatory analysis; i.e., the exploratory search for new hypotheses must not be presented as **h**ypotheses testing **a**fter **r**esults are **k**nown (HARKing; cf., Kerr 1998). Attaching the label "significance test" to $p$-values in exploratory studies is misleading per se since no testable hypotheses exist. As reminder for researchers and readers alike, Berry (2016: 2) thence suggests to include a "black-box warning" into all exploratory studies: "Our study is exploratory and we make no claims for generalizability. Statistical calculations such as p-values and confidence intervals are descriptive only and have no inferential content." In contrast, not accounting for multiple testing in confirmatory analysis inflates significance claims. It is nonetheless common in economics – even though ever-increasing numbers of variables are included in regression models due to better data availability (Ioannidis and Doucouliagos 2013).

*Multiple testing adjustment requirements*

We engage in multiple comparisons every time we perform and interpret significance tests for more than one statistical hypothesis based on *one* data set. Covert multiple testing ($p$-hacking) is hard to detect and overcome because it withholds the information that is needed to correct for multiplicities. By contrast, adjustment tools are available for those who are ready to account for the multiple tests they have made in accordance with their research interests. If, in a multiple testing setting, we are interested to make claims of statistical significance that are comparable to the one we would make for a single significance test (i.e., at a threshold comparable to the conventional 0.05), we have to adjust the significance levels.

Both the Bonferroni and the Bonferroni-Holm adjustment are used to control the so-called *family wise error rate* (FWER or multiple type I error rate). The logic behind the FWER correction is to restrict the probability of rejecting even one null hypothesis when it is true, irrespective of how many of the other null hypotheses are valid (Hochberg and Tamhane 1987; Pigeot 2000). The computationally easiest way to arrive at a family wise error rate is the Bonferroni correction (see Fisher 1935). Assuming the conventional significance threshold of 0.05, the Bonferroni correction implies that the adjusted threshold of $0.05/m$ is used for each of a total number of $m$ tests. For large $m$, this leads to extremely small significance levels. This, in turn, usually results in only few null hypotheses being rejected. In our example for $m = 3$ tests (see Table 3), only the effect associated with the "raw" $p$-value of 0.01 would be declared significant since the adjusted threshold is 0.0167. Hence, alternative and less con-

---

[2] This seems to be what Fisher had in mind when noting that low $p$-values *signify* "worth a second look" if we have little to none prior knowledge (Gigerenzer et al. 2004; Lecoutre and Poitevineau 2014; Nuzzo 2014).

servative corrections have been suggested. The Bonferroni-Holm adjustment (Holm 1979) is an example. It calculates different thresholds for each one of the $m$ tests. The respective threshold is computed according to the formula $0.05/(m + 1 - r)$ and increases with the rank $r$ of the raw $p$-values. The rationale is to identify the largest raw $p$-value that is below or equal to the adjusted threshold and then declare this and all smaller raw $p$-values "significant." Following this rule, we would still only declare the raw $p$-value of 0.01 "significant" even though the Bonferroni-Holm thresholds are indeed less rigorous than the Bonferroni thresholds.[3]

Table 3: Multiple testing adjustments for $m = 3$ consistent with a single-test threshold of 0.05

| | | | |
|---|---|---|---|
| Conventionally computed ("raw") $p$-values | 0.01 | 0.03 | 0.04 |
| Rank $r$ of raw $p$-values | 1 | 2 | 3 |
| **Family wise error rate (FWER) corrections** | | | |
| Significance threshold after Bonferroni adjustment: $0.05/m$ | **0.0167** | 0.0167 | 0.0167 |
| Significance threshold after Bonferroni-Holm adjustment: $0.05/(m + 1 - r)$ | **0.0167** | 0.0250 | 0.0500 |
| **False discovery rate (FDR) correction** | | | |
| Significance threshold after Benjamini-Hochberg adjustment: $0.05 \cdot r/m$ | **0.0167** | **0.0333** | **0.0500** |

A different perspective is adopted when using the concept of the *false discovery rate* (Benjamini and Hochberg 1995) according to which different significance levels for each of the $m$ tests are computed based on the formula $0.05 \cdot r/m$. Again, the rationale is to identify the largest raw $p$-value that is below or equal to the adjusted threshold and then declare this and all smaller raw $p$-values "significant." The FDR correction produces less rigorous thresholds than the FWER correction. With the FDR rule, we would declare all three raw $p$-values "significant." To provide intuition, one could say that the FDR is aimed at restricting the rate of valid null hypotheses being rejected relative to the total number of rejected hypotheses. One should keep in mind that, despite its delusive naming, the FDR (along with the other non-Bayesian multiple testing adjustments) shares the crucial limitation of the frequentist $p$-value approach: it does not estimate the post-study probabilities of hypotheses given the data. Analogous to the $p$-value, the frequentist FDR cannot work backwards and inform us about the probability of real-world phenomena (see section 2.2) – even though we as researchers would want to have that kind of (inherently Bayesian) information to assess the trustworthiness of our inductive conclusions.[4]

It was easy to determine the number of tests $m$ that must be corrected for in our multiple regression example. However, in many research settings, the need for multiple testing corrections may be less obvious and the determination of $m$ may often be difficult. Many questions may arise: what should we do when separate regressions are used for multiple variables of interest? How should we deal with interaction terms, higher-order polynomials etc. that are included in the regression besides the original variables? How should we consider tests of multiple model specifications? The answer is straightforward: whenever we perform multiple tests on *one* data set, they need to be corrected for to prevent the inflation of statistical significance claims.[5] Control variables, in contrast, which often populate regression models in large

---

[3] Another concept, the *global level of simultaneous tests,* can be directly linked to the binomial distributions as shown in Table 2. The global level of simultaneous tests is aimed at restricting the probability of rejecting one or more of the tested null hypotheses when *all of them* are valid?

[4] To avoid confusion, the *frequentist* false discovery rate must be clearly distinguished from the term's Bayesian interpretation as used, for example, by Colquhoun 2014, Hirschauer et al. (2016), and Motulsky (2014). The *Bayesian* false discovery rate (or: Bayesian error rate) is the post-study probability of "no effect" and therefore the probability of a faulty scientific claim when rejecting the null.

[5] Although adjusting the significance level is necessary in most cases of multiple testing, there are special constellations where an inherent adjustment takes place within the very system of tested hypotheses. This is the case for so-called *closed families of hypotheses* in the special case of a coherent test procedure. The least-significant-

numbers, need to be considered in multiple testing adjustment. Instead, they have to be clearly identified as control variables and separated from the $m$ variables of interest. Disregarding control variables is only adequate, however, if we explicitly distinguish between confirmatory and exploratory analysis. This may require to clearly divide studies in two parts: a confirmatory part where we perform a pre-defined number of multiple tests $m$, and an exploratory part where we look at potentially interesting correlations in the control variables. As has been said before, conventional $p$-values are an adequate focusing aid to identify what might be worth investigating with new data. If this exploratory search is not presented as HARKing (cf., Kerr 1998), which itself would be illegitimate hidden testing, no multiple testing corrections are needed.

We may summarize that the appropriate multiple testing adjustment depends on the research setting. It is an often ambiguous choice on which we cannot further elaborate in this paper. For a general overview and description of the great variety of multiple comparison adjustments and their eligibility criteria the reader is referred to Bretz et al. (2010) or Westfall et al. (2011).

## 2.2 Semantically induced misinterpretations of the $p$-value

Even if significance claims are not inflated, serious misunderstandings lurk due to the delusive technical terms of frequentist statistics that contradict everyday language. The fact that a vast body of literature has decried these misinterpretations over the last four decades (see Hirschauer et al. 2016 for an overview), has apparently been of little avail. Haller and Krauss (2002), Gigerenzer et al. (2004), Lecoutre and Poitevineau (2014), Krämer (2011), and Nickerson (2000) point out that the ubiquity and persistence of faulty interpretations are, not least, caused by the fact that they have been perpetuated over decades through academic teaching and even through best-selling statistics textbooks.

In econometric applications, the most common misinterpretations of the $p$-value are best understood when realizing that the majority of analyses are performed as if using the following misleading guidelines: (1) Run a multiple regression. (2) Compute coefficients and $p$-values. (3) Declare coefficients with $p$-values below a threshold (usually 0.05) "statistically significant." (4) Do not reflect on the arbitrary dichotomization of results into "statistically significant" and "not statistically significant." (5) Implicitly attach a high trustworthiness or probability (possibly even $1 - p$) to statistically significant coefficients. (6a) Do not discuss the effect size or (6b) even suggest that a "significant effect" is also large or important or (6c) simply claim the just-estimated effect to be real. (7) Attach the label "statistically nonsignificant" to coefficients with $p$-values above 0.05 instead of noting that they are "not statistically significant" at the defined threshold. (8a) Interpret your "statistically nonsignificant" results as minor or not noteworthy or (8b) even suggest that they can be interpreted as a proof of no effect.

As a consequence of such misleading practices, researchers and readers alike will not only overrate the inferential content of findings that have been declared "statistically significant" but also of finding that have been declared "statistically nonsignificant." These dichotomous declarations will make them draw similarly dichotomous existence/relevance conclusions that – to borrow a quote from H.L. Mencken in the New York Evening Mail from November 16, 1917 – are "neat, plausible, and wrong." The first neat but wrong conclusion is that effects with the label "statistically significant" can be considered to be real (and may be large) with a high probability. The second neat but wrong conclusion is that the label "statistically nonsignificant" is an indication or even proof of no or little effect. Interpretations along this erroneous dichotomy is entrenched practice for many researchers, a practice that is based on deeply

---

differences test of Fisher (1935) in case of one-way ANOVA falls into this class when three groups are compared, but not for more than three groups. For details, see Pigeot (2000) or Didelez et al. (2006).

engrained beliefs that result both from wishful thinking (desire for a "neat" interpretation) and the inevitably wrong connotation of the word "significance" in everyday language.

Unfortunately, the *p*-value, if correctly interpreted, has much less inferential content than what colloquial associations with the terms "statistically significant" and "statistically nonsignificant" suggest (Berry 2016). Following the misleading guidelines from above mirrors a serious misunderstanding of what the *p*-value, which is merely a summary statistic of a given data set, can tell about reality. It means falling prey to three fallacies that have been given distinct names: Cohen (1994) used the term *inverse probability error* to describe the belief that the *p*-value is the conditional probability of (falsely rejecting) the null hypothesis given the data under study. Instead, the *p*-value is the conditional probability of finding the observed effect (or even a larger one) in random replications *if*, as a thought experiment, we assumed the null hypothesis to be true. Per definition, it cannot work inversely and inform us on the underlying reality. But looking for answers to their scientific questions about reality, statistical practitioners often confuse frequentist and Bayesian probabilities and adopt a Bayesian interpretation of frequentist measures such as the *p*-value.[6] McCloskey and Ziliak (1996) coined the expression *sizeless stare* for the disregard of effect size or the implicit equation of statistical significance with relevance. In the light of the increasing availability of large samples, the naïve equation of significance with relevance becomes more and more misleading because any effect, even if very small and irrelevant, eventually becomes statistically significant in large samples.[7] Hirschauer et al. (2016) used the term *false dichotomy* to describe the logical fallacy that misleads people to first adopt an ill-founded either-or perspective and then use the ensnaring label "statistically nonsignificant" that finally makes them interpret *p*-values above 0.05 as an indication or even confirmation of the null.[8]

Colloquial associations that rashly equate "statistically significant" with "scientifically trustworthy" may furthermore prevent researchers from realizing that a meaningful interpretation of the *p*-value requires that the data at hand represent (at least approximately) a *random* sample of a defined parent population.[9] In many contexts, the fact that even a truly random sample may not exactly reflect the properties of the population (random sampling error) is the least of worries. Data are frequently not obtained through random sampling but affected by various types of selection bias or measurement errors. These problems are often more serious

---

[6] Cohen (1994: 997) succinctly described the harmful mixture of wishful thinking and semantic confusion that causes the inverse probability error: "[the *p*-value] does not tell us what we want to know, and we so much want to know what we want to know that, out of desperation, we nevertheless believe that it does! What we want to know is 'given these data, what is the probability that H0 is true?' But [...], what it tells us is 'given that H0 is true, what is the probability of these (or more extreme) data?' "

[7] It is unlikely that any two real-world variables exhibit zero correlation. This is why Lecoutre and Poitevineau (2014: 50) call the null hypothesis a "straw man" that significance testing tries to knock down. Similarly, Leamer (1978: 89) notes that since "a large sample is presumably more informative than a small one, and since it is apparently the case that we will reject the null hypothesis in a sufficiently large sample, we might as well begin by rejecting the hypothesis and not sample at all."

[8] The inverse probability error and the false dichotomy fallacy, even though jointly found in many instances, are inconsistent in themselves. Committing the inverse probability error, one would believe that $p = 0.01$ indicates a 1%-probability of the null being true and thus a 1%-probability of making a false claim when rejecting it. Similarly, one would have to believe that a "nonsignificant" result with $p = 0.15$, for example, indicates a 15%-probability of the null. After this error, one cannot logically interpret a "nonsignificant" result as a confirmation of the null.

[9] Noting that "where there is a sample there must be a population," Denton (1988: 166f.) points out that conceiving of the population can be difficult. The easiest case is a sample drawn from a finite population such as the entirety of a country's citizens. Slightly less intuitive is an experiment such as flipping a coin. Here, the population is an imaginary set of coin flip experiments that are infinitely repeated under constant conditions. More conceptual challenges arise when working with observational data. Here, the frequentist statistician must introduce an "unseen parent population" that is subjected to a noise producing process from which we observe one random realization.

than the random sampling error. Not having a random sample has fundamental implications for the utilization of the *p*-value that, arithmetically, can be computed for any data: (1) *p*-values computed from non-random samples have no meaningful interpretation. (2) The same holds when the sample itself represents the population of interest. In this case, the random sampling error is completely eliminated because we have already assessed the properties of the population. (3) The *p*-value can only be meaningfully interpreted if researchers explicitly define *from which* population the random sample has been drawn and thence *to which* population a statistical inference is to be made. (4) *Statistical* inference is concerned with the relationship between a random sample and its population. This is only the first step of *scientific* inference, which is the totality of reasoned judgments (inductive generalizations) that we (can) make in the light of our own observations and the available body of evidence found elsewhere. There is no easy-to-apply recipe of making such judgements. We instead have to critically rethink each situation. When trying to answer the question, for example, of what we can learn from an experiment with agricultural students for a country's student population or even its citizens, or even human beings in general, we must keep in mind that a *p*-value can do *nothing* to assess the generalizability of a result beyond the parent population from which the random sample has been drawn.

## 2.3 Exaggerated focus on one-shot studies

Economic analyses are frequently based on multiple regression in which a "dependent" or response variable (usually denoted as $y$) is modeled as a function of several "independent" variables (usually denoted as $x_j$). The independent variables are often divided into focal variables of interest (*focal predictors*) and variables that are used to control for confounding influences (*control variables*). The direct outcomes of multiple regression analysis are the estimated coefficients $\hat{\beta}_j$ (regression slopes) that relate the predictors $x_j$ to response $y$. If several studies have tackled the same *x-y* relation (e.g., between education $x_1$ and wage earnings $y$), we are naturally interested in summarizing these findings. Contrary to narrative reviews, quantitative surveys (meta-analyses) are aimed at computationally synthesizing the results from primary studies.

Arguably going back to the paper by Stanley and Jarrell (1989), some economic meta-analyses have been carried out on selected issues over the last decades. Zelmer (2003), Cooper and Dutcher (2011), Engel (2011), and Lange (2016) addressed economic experiments. Non-experimental meta-analyses were conducted by Card and Krueger 1995, Crouch (1995), Loomis and White (1996), Fitzpatrick et al. (2017), and Van Houtven et al. (2017). More than ten years ago, even a whole issue of the *Journal of Economic Surveys* (cf., Roberts 2005) was dedicated to the meta-analysis of regression coefficients. In the meanwhile, there is also a *Meta-Analysis of Economics Research Network* that provides a platform for economic meta-analyses. Nonetheless, the practice of meta-analysis is less common in economics than other fields.[10] What is more, even the statistical literature has mainly dealt with the synthetization of measures such as (standardized) mean differences or (risk or response) ratios that are primarily used in non-economic fields (e.g., medical sciences). In contrast, summarizing coefficients from multiple regressions, which are the working horse of economists, has attracted less attention (Becker and Wu 2007). The limited use of meta-analysis in economics can be attributed to the fact that economic research is mainly a non-programmed bottom-up research exercise. As such, it produces an enormous quantity of empirical results on topical issues, but is also plagued by an enormous heterogeneity of regression model specifications (Bruns

---

[10] Even a cursory look at economic publications will show that the critique by Stanley and Jarrell (1989: 162) still applies: "The reviewer often impressionistically chooses which studies to include in his review, what weights to attach to the results of these studies, how to interpret the results, […]. Traditionally, economists have not formally adopted any systematic or objective policy for dealing with the critical issues which surround literature surveys. As a result, reviews are rarely persuasive to those who do not already number among the converted."

2017). Attempts to summarize findings that deal with the same *x-y* relation (e.g., education-wage) are thus hampered by a deficient or even lacking comparability of the regression coefficients across primary studies.

The comparability of regression slopes across studies is severely constrained by some basic features found in most economic research fields: first, the metrics (measurement scales and units of measurement) of both independent and dependent variables usually differ across studies. Second, even if all variables are identically measured, using structurally different econometric models involves that the estimated coefficients are usually beyond comparison. Third, even if identical metrics and econometric model structures are used, comparability is jeopardized when models with different sets of independent (control) variables are estimated. Fourth, even if metrics and estimation models were similar, different studies may have drawn their samples from different populations. Consequently, there may simply be no data base to do a meta-analysis because each single study covers a different parent population. Regression coefficients for the education-wage relation in France, Ghana, and the US, for example, cannot be meaningfully synthesized into one summary coefficient.

We use the univariate *weighted least squares* approach to demonstrate what meta-analysis is about in principle. We assume that we are to summarize 20 individual studies based on different sample sizes ($n = 20, 30, 40, 50$). For the sake of easy intuition, we avoid all complications, such as different metrics between studies, by simulating 20 samples (random realizations) from a "reality" characterized by the linear relationship $y = \beta_0 + \beta_1 x_1 + e$, with $\beta_0 = 1, \beta_1 = 0.2, x_1 \in \{0.5, 1.0, 1.5, \dots, n/2\}$, and $e \sim N(\mu; \sigma)$, with $\mu = 0$ and $\sigma = 5$. In each single study, an OLS regression is used to estimate $\hat{\beta}_0$ and the focal coefficient $\hat{\beta}_1$.

Using the weighted least squares method, the summary coefficient $\hat{\beta}_1^{sum}$ that synthetizes the coefficients of the single studies is computed as follows (cf., Becker and Wu 2007: 7):

$$\hat{\beta}_1^{sum} = \sum_{i=1}^{I} \hat{\beta}_{1i} \cdot w_{1i} \bigg/ \sum_{i=1}^{I} w_{1i}, \quad \text{with } w_{1i} = 1/SE_{1i}^2 \tag{1}$$

where $I$ is the number of single studies and $\hat{\beta}_{1i}$ is the coefficient estimated in the *i*th study. The weight $w_{1i}$ that is attributed to the coefficient from each study $i$ is the reciprocal of its squared standard error $SE_{1i}^2$; and the ratio $w_{1i}/\sum_{i=1}^{I} w_{1i}$ denotes the percentage weight of each study. The standard error of the summary coefficient $\hat{\beta}_1^{sum}$ is:

$$SE_1^{sum} = \left( 1 \bigg/ \sum_{i=1}^{I} w_{1i} \right)^{0.5} \tag{2}$$

Table 4 describes the results of the meta-analysis. The calculation of *p*-values is based on test scores derived from a one-sided test. This reflects the assumption that the researchers who presumably had carried out the 20 primary studies had qualitative prior knowledge indicating a non-negative relation between $x_1$ and $y$. We compute the *p*-values assuming that all test scores follow a standard normal distribution. Several noteworthy findings and conclusions can be derived from Table 4.

1. A large majority of studies (14 in 20) have not found a statistically significant result. This might mislead narrative reviewers to contrast tallies and conclude that the results of these 20 studies represent contradictory evidence or even overall a confirmation of no effect.

2. Meta-analysis is capable of leaving behind the arbitrary either-or interpretation within each study. Instead, it synthesizes – with an adequate weight – the informational content of all studies given the fact that "the effect best supported by the data from a given experiment is always the observed effect, regardless of its significance" Goodman (2008: 136).

3. The meta effect size $\hat{\beta}_1^{sum} = 0.195$ approximates the true $\beta_1 = 0.2$ quite well. The meta $p$-value of 0,000000003 shows that even a great majority of studies that are not statistically significant can together represent a "highly statistically significant" effect. This is due to the fact that meta-analysis is capable of including the informational content of studies even if they are too small to produce statistical significance.

4. Given the meta effect size and its very low $p$-value, we would be confident that the real-world level of $\beta_1$ lies above 0 (and we would commonly refer to its estimated value of approximately 0.2), even though we realize that the $p$-value does not provide a clear rationale or even calculus for statistical inference (Goodman 2008).

5. Low $p$-values do not indicate results that are "more trustworthy" to consider than others. Considering only significant studies, for example, would introduce a distortion (see section 3) and we would find a summary coefficient of 0.3109. That is, the results of *all* studies jointly represent the body of evidence and are valuable and *necessary*, irrespective of their $p$-value, to provide an approximately correct picture of the real-world regularity.

6. The 20 single studies in our illustrative example were *not* distorted but based on 20 random realizations (simulations). If the single studies were distorted due to publication bias (see section 2.4), the basic *weighted least squares* method of meta-analysis, which is unable to control for such biases, would simply summarize the distortion.[11]

7. Being in the comfortable position to know all raw data, we also carried out a single large regression over all $n = 700$ observations that serves as benchmark for the meta-analytical calculus. The estimated $\hat{\beta}_1^{700} = 0.2123$ from the large regression is slightly above the true effect size $\beta_1 = 0.2$ and the computed meta effect size of $\hat{\beta}_1^{sum} = 0.195$.

---

[11] While the *weighted least squares* approach is not able to control for publication bias, meta-regression (cf., Stanley and Doucouliagos 2012; Stanley and Jarrell 1989) has been suggested to control for the idiosyncrasies of model specifications in primary studies and notably for publication bias.

Table 4: Meta-analysis for 20 single studies, each based on a simulated random sample from a reality characterized by the $x$-$y$ relation: $y = 1 + 0.2x_1 + e$, with $e \sim N(0; 5)$

| Study No. $i$ | Observations $n$ per study | Estimated coefficient $\hat{\beta}_1$ | Standard error | $p$-value | Statistically significant [a] | Weight (%) |
|---|---|---|---|---|---|---|
| 1 | 20 | 0.0084 | 0.2669 | 0.4875 | 0 | 1.56 |
| 2 | 20 | -0.0499 | 0.3663 | 0.5541 | 0 | 0.83 |
| 3 | 20 | 0.1377 | 0.4240 | 0.3726 | 0 | 0.62 |
| 4 | 20 | -0.0830 | 0.2945 | 0.6110 | 0 | 1.29 |
| 5 | 20 | 0.1658 | 0.4708 | 0.3624 | 0 | 0.50 |
| 6 | 30 | 0.1206 | 0.2008 | 0.2740 | 0 | 2.76 |
| 7 | 30 | 0.2397 | 0.2486 | 0.1675 | 0 | 1.80 |
| 8 | 30 | 0.5598 | 0.2058 | 0.0033 | 1 | 2.63 |
| 9 | 30 | 0.4409 | 0.2051 | 0.0158 | 1 | 2.65 |
| 10 | 30 | 0.3295 | 0.1783 | 0.0323 | 1 | 3.51 |
| 11 | 40 | 0.2442 | 0.1561 | 0.0588 | 0 | 4.58 |
| 12 | 40 | 0.1415 | 0.1229 | 0.1247 | 0 | 7.38 |
| 13 | 40 | 0.2533 | 0.1269 | 0.0230 | 1 | 6.92 |
| 14 | 40 | 0.2125 | 0.1432 | 0.0689 | 0 | 5.44 |
| 15 | 40 | 0.1770 | 0.1462 | 0.1131 | 0 | 5.21 |
| 16 | 50 | 0.2841 | 0.0946 | 0.0013 | 1 | 12.44 |
| 17 | 50 | 0.1648 | 0.1201 | 0.0850 | 0 | 7.72 |
| 18 | 50 | 0.0758 | 0.0848 | 0.1856 | 0 | 15.50 |
| 19 | 50 | 0.2675 | 0.1203 | 0.0131 | 1 | 7.70 |
| 20 | 50 | 0.0973 | 0.1116 | 0.1917 | 0 | 8.95 |
| **Summary** | **700** | **0.1950** | **0.0334** | **0.0000** | **1** | **100.00** |

[a] 1 = yes at the 0.05 level; 0 = no at the 0.05 level.

Finally, one should bear in mind that conventional meta-analysis stays within the narrow confines of the frequentist approach: it does not provide probabilities of scientific propositions given the data, or, as Kline (2013: 307) notes with regard to experimental data "a standard meta-analysis cannot answer the question, What is the probability that the treatment has an effect?" Only Bayesian methods can provide the post-study probabilities of scientific propositions that researchers and users of scientific results are ultimately interested in (see footnote 6). While it is difficult but possible to combine Bayesian methods with meta- analytical approaches (Howard et al. 2000; Kline 2013: 307), we can all the more do without Bayesian methods the more studies regarding a specific scientific issue are included into meta-analysis. This is due to the fact that we consider an increasing number of observations and thus increasing evidence by including more and more studies. Correspondingly, the amount of studies and thus (prior) knowledge that remains unconsidered is declining. In the extreme, we include *all* available studies and consequently have an uninformative (flat) prior beyond these studies. If so, the meta $p$-value approximates the post-study Bayesian error rate (Zyphur and Oswald 2015).

The complications in multiple regressions that restrict the feasibility of meta-analyses, as well as the eligibility of the meta-analytical approaches that are available to deal with these problems, are beyond this paper's scope. For further study of meta-analysis, the reader is referred to Becker and Wu (2007), Card (2012), Kline (2013: chapter 9), and Schmidt and Hunter (2014). For an introduction to Bayesian methods, the reader is referred to Hartung et al. (2008: chapter 12), Pitchforth and Mengersen (2013), and Zyphur and Oswald (2015).

## 2.4   Publication bias

The ill-directed incentives of the present publication system produce a bias towards statistical significance (Fanelli 2011). This bias is arguably more covert and higher in economics than

other fields (Fanelli 2010), not least because replication is not a popular exercise among non-experimental economists (Evanschitzky and Armstrong 2010). While some researchers may honestly but erroneously believe that "starless" results are not interesting enough to warrant publication (Sterne et al. 2008), the major problem is the one highlighted by the "publish or perish" witticism: in our competitive research system, most researchers are under pressure to produce journal papers with novel findings. If papers with statistically significant findings (positive results) are more likely to be published, researchers are likely to adopt one or several selection strategies: *selective preparation* means not to conduct (replication) studies that are likely to be a "waste of time" because they do not promise to produce statistically significant novelties. *Selective submission* implies not to submit papers with results that happen to be not statistically significant because they are unlikely to be published. *Selective reporting* means *p*-hacking and the selective presentation of an analytical variant that "worked best" in terms of producing significance. Beyond the single researcher's sphere of influence, there is finally the problem of *selective publishing*. Even if researchers are conscious and scrupulous enough to refrain from self-interested selection practices, reviewers and editors have ample discretion to promote papers for publication that contain seeming novelties. As result of all these selection processes, "significant" findings are overrepresented and studies with negative results tend to stay in researchers' "file drawers" (Rosenthal 1979) without ever being presented to the public. In other words, they are missing but not missing at random and thence cause a bias for significance. This leads to the definition by Kline (2013: 274) according to which publication bias implies that published studies have more "statistically significant" findings and larger effect sizes than unpublished studies (including the ones that have not been made in the first place).

Over the last decades, a variety of meta-analytical methods have been developed to gauge publication bias (see Table 5). For a description of these methods and their respective potential to identify selection procedures the reader is referred to Cooper et al. (2009), Rothstein et al. (2005), Song et al. (2000), or Weiß and Wagner (2011). In this paper, we will have to limit ourselves to briefly describing a few selected approaches.
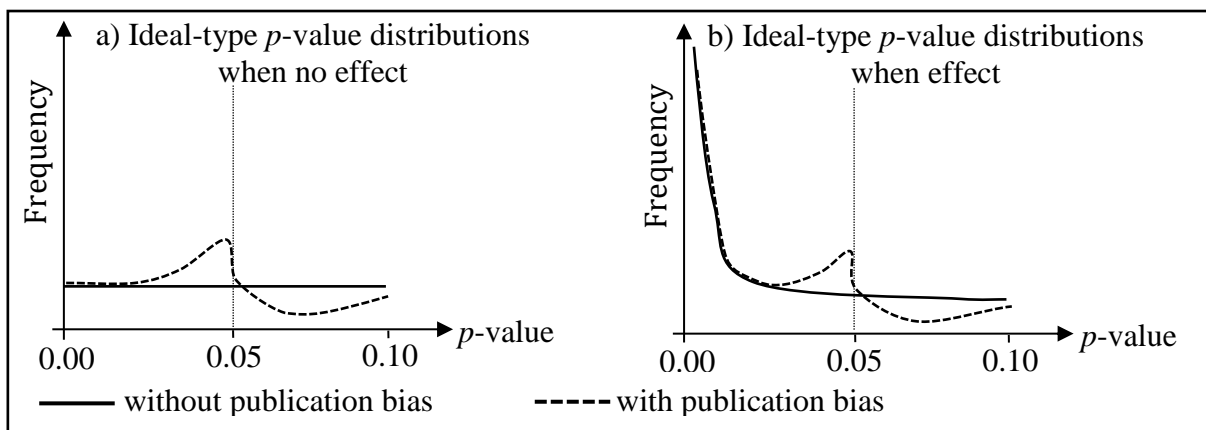
Table 5: Selected methods for detecting publication bias

| Method | Authors |
|---|---|
| Comparison of published and unpublished studies | Song et al. (2000); Sterne and Egger (2005) |
| Caliper test | Berning and Weiß (2016); Gerber and Malhotra (2008) |
| *p*-curve- analysis | Head et al. (2015); Simonsohn et al. (2014) |
| Funnel plots | Egger et al. (1997); Light and Pillemer (1984) |
| Capture-recapture method | Bennett et al. (2004); Poorolajal et al. (2010) |
| Fail-safe N | Rosenberg (2005); ROSENTHAL (1979) |
| Selection models | Kicinski et al. (2015); Silliman (1997) |

A first attempt to gauge publication bias is to compare gray literature (e.g., discussion and conference papers) with published studies (Song et al. 2000). Larger effect sizes (and smaller *p*-values) in published compared to non-published papers are an indication of publication bias caused by selective submission and/or selective editorial policies. Besides the problem that unpublished studies are less disseminated, neither selective preparation nor *p*-hacking can be identified through this comparison. This is why the assessment of publication bias is often based on "anomalies" in the structure of test statistics. A widely applicable method is the *caliper test* (Gerber and Malhotra 2008). Its idea is to compare the frequency of reported test scores within a small band above and below the usual significance thresholds. Gerber and Malhotra (2008: 6) claim that "there is no reason to expect that, in the narrow region just above and below the critical value, there will be substantially more cases above than below the critical value unless the .05 level was somehow affecting what is being published." In-

stead of test scores, *p-curve analysis* looks directly at the distribution of *p*-values (Head et al. 2015; Simonsohn et al. 2014). Figure 1 illustrates the idea: while we do not know whether there is an effect or not, we know that the *p*-value distribution is uniform if there is no effect. We also know that it is exponential with a right skew if there is an effect. In both cases, anomalies around the critical threshold should not occur and, if found, are taken as an indication of publication bias. Whereas an overrepresentation of *p*-values below 0.05 can be attributed to *p*-hacking, all selection procedures jointly contribute to the underrepresentation of values above the threshold.

Figure 1: illustration of the basic idea of *p-curve analysis*



In many research fields, publication bias itself has become an important object of study. Joober et al. (2012: 149), for example, report that in some medical areas almost no negative studies exist. They also find that publication bias has increased in many fields over the last years. The issue has also been taken up in the political and sociological sciences (cf., e.g., Auspurg and Hinz 2011; Gerber et al. 2010). In a recent study using the caliper test, Berning and Weiß (2016) find strong evidence for publication bias in papers published from 2001 to 2010 in three flag ship journals of the German social sciences (Kölner Zeitschrift für Soziologie und Sozialpsychologie, Zeitschrift für Soziologie, and Politische Vierteljahresschrift). With a few early exceptions (e.g., Denton 1985; Lovell 1983), the awareness and study of publication bias has been less pronounced in economics in the past. But recently, a large-scale study by Brodeur et al. (2016) analyzed the distribution of about 50,000 test statistics that were published from 2005 to 2011 in three of the most prestigious economic journals (American Economic Review, Journal of Political Economy, Quarterly Journal of Economics). They find a considerable overrepresentation of marginally significant test statistics as well as a sizeable underrepresentation of marginally "nonsignificant" statistics. Their "interpretation is that researchers inflate the value of just-rejected tests by choosing 'significant' specifications" (Brodeur et al. 2016: 1).

Many suggestions have been made to mitigate publication bias (see Munafò et al. 2017 for an overview). Song et al. (2013) and Weiß and Wagner (2011) propose to strengthen alternative publication outlets. They also call for a general change of editorial policies towards giving equal publishing chances to all scientific results, including replications and negative findings. To reduce the risk of selective reporting and increase the chance of being published independent of whether positive or negative results are eventually found, Munafò et al. (2017) suggest to go beyond the after-study provision of raw data, and peer-review and register complete study designs *before* they are carried out. Across a variety of disciplines, various initiatives try to institutionalize efforts counteracting distorting selection procedures. With a view to the dire consequences of publication bias in medical research, a global initiative *All Trials Registered/All Results Reported* was launched in 2013. Along the same lines, the *Journal of Negative Results in BioMedicine*, the *PLOS ONE Journal*, and the *All Results Journals* explicitly encourage replication studies and pursue policies of publishing positive and negative results.

Institutionalized efforts to strengthen the practice of replication and preregistration seem be weak in economics compared to other fields such as the medical sciences that spearhead the development. In a study of all 333 economic Web-of-Science journals, Duvendack et al. (2015) find that most of them still give very low priority to replication. Pre-registration of studies also seems to lag behind other fields. No economic journal, for example, is among the approximately 40 journals that, according to the listing of The *Center for Open Science*, have adopted a policy of peer reviewing and registering study designs before results are known (Duvendack et al. 2017). But things are changing. A topical initiative is the call of the *economics-ejournal* in which researchers are asked to select a published study as a candidate for replication and to discuss how they would carry out the replication. There are also some noteworthy replication platforms for economists such as *The Replication Network* and *Replication in Economics* that provide data bases of replications and the opportunity to publish replication studies. The issue has also attracted the attention of professional economic societies such as The American Economic Association that concerned itself with replication on its 2017 annual meeting and has launched a preregistration scheme for randomized controlled trials. In this scheme, a study's design is peer-reviewed based on its methodological quality and registered if accepted. Peer-review of a study's design and formal registration by a prestigious institution are meant to provide equal chances of being published independent of which results are eventually found.
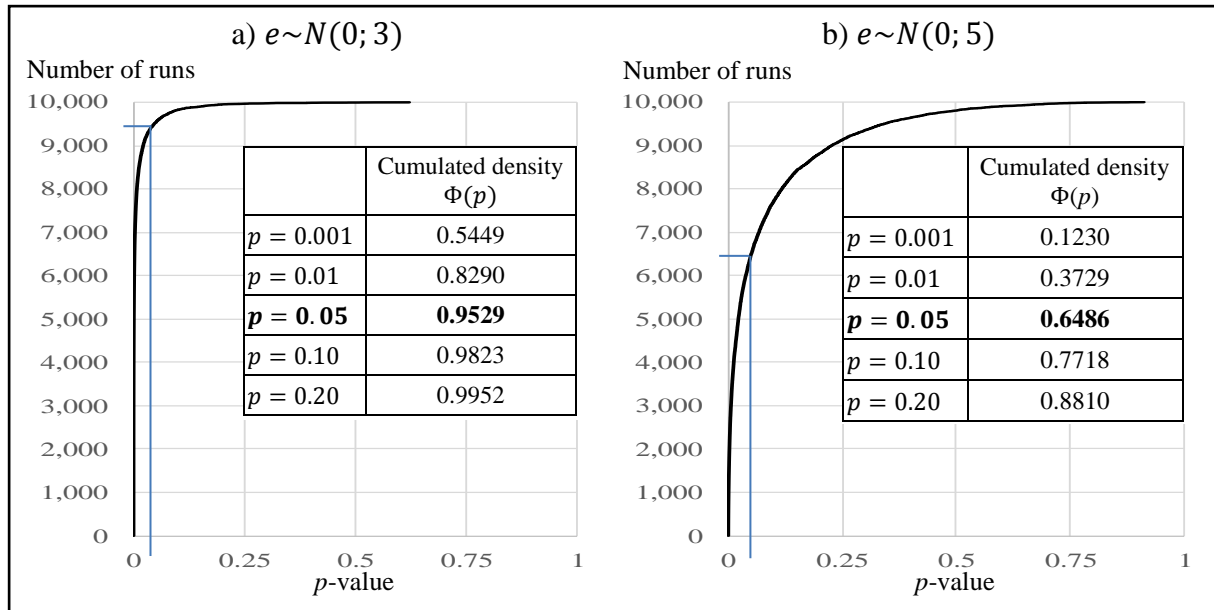
## 3  Disregard of *p*-value sample-to-sample variability

In econometric practice, results with small *p*-values are usually declared significant and claimed to represent substantial evidence for a real phenomenon or even simply assumed to be real. But even if researchers do not commit the inverse probability error of interpreting the *p*-value as the post-study probability of making a false claim, many believe that after having found a low *p*-value, repeated random sampling or repetition of an experiment will produce a similar statistical verdict (Goodman 1992; Halsey et al. 2015). The erroneous belief that the *p*-value indicates the likelihood that significant results can be replicated has been called "replication fallacy" (Gigerenzer et al. 2004). It is partly due to the fact that, contrary to other statistical estimates, no measure of the *p*-value's variability – even though it can be of considerable magnitude – is usually reported (Boos and Stefanski 2011). This omission is reflected in the conventional econometric practice to not reflect on, let alone, use the concept of statistical power.

Figure 2 illustrates why relying on the *p*-value as a measure of evidence without considering its sample-to-sample variability falls short of the mark even within an otherwise correct frequentist interpretation. We have simulated the data and thence know more than researchers who are limited to analyzing data from *one* random realization. We have generated samples for two "realities." Both are characterized by the linear relationship $y = 1 + 0.2x_1 + e$, but their normally distributed error terms are $e \sim N(0; 3)$ and $e \sim N(0; 5)$, respectively. Sample size is $n = 50$, with $x_1$ varying from 0.5 to 25 in equal steps of 0.5. For each reality, we run OLS-regressions for each of the 10,0000 simulation runs. That is, we use an estimator that perfectly fits the data. The *p*-values resulting from these regressions are summarized in Figure 2. The left-hand side shows the cumulated densities (number of runs) of the *p*-value for the normally distributed error term $e \sim N(0; 3)$. The right-hand side shows them for the error term $e \sim N(0; 5)$.

The cumulated densities that are tabulated for selected *p*-values reveal that the size of the error term has a considerable impact on the *p*-value's sample-to-sample variability. We find *p*-values above 0.1 in only 1.77% (= $1 - 0.9823$) of replications in the case of $e \sim N(0; 3)$. This rate increases to 22.82% (= $1 - 0.7718$) in the case of $e \sim N(0; 5)$. The tabulated figures illustrate why many statisticians call for complementing the *p*-value approach with statistical power considerations. Statistical power is defined as the conditional probability of a significant result over many replications if the effect is true. Having generated the data, we know the true effect $\beta_1 = 0.2$. We thence also know the true power to be the cumulated density $\Phi$ of

the $p$-value distribution at the 0.05 level. For $e \sim N(0; 3)$, we find a power of 95.29% (i.e., 9,529 runs, out of 10,000, show $p \leq 0.05$). In contrast, the power is only 64,86% in the case of $e \sim N(0; 5)$.

Figure 2: $p$-value distribution over 10,000 replications ($y = 1 + 0.2x_1 + e$; sample size $n = 50$) [a)]



a) One-sided test; normal test statistic.

Cumming (2008) notes that $p$-values, unless they are very low, provide a very unreliable signal of what is going to happen in replications. While this is a correct qualitative statement, the "unless very low" comment requires further attention. We need to keep in mind that the $p$-value is merely a property of the data (sample). This implies that "P values are only as reliable as the sample from which they have been calculated. A small sample taken from a population [with big noise] is unlikely to reliably reflect the features of that population" Halsey et al. (2015: 180). While most researchers realize that small samples and big noise increase the *level* of the $p$-value that is to be expected, they are less likely to be fully aware of the fact that they also considerably increase the sample-to-sample *variability* of the $p$-value.

Comparing the two $p$-value distributions in Figure 2 helps to understand that, besides a single study's $p$-value, its variability – and in dichotomous significance testing the statistical power – determines the informational value of a finding.[12] High-powered studies are more reliable in that they reduce not only the $p$-value's average level but also its variability over random replications. In many cases, power is below the 64.86% found in our $e \sim N(0; 5)$ simulation. But even here it is interesting that, for example, one has a 12% probability of finding a $p$-value below 0.001, but likewise a 12% probability of finding a $p$-value above 0.2. We should very cautiously interpret a $p$-value since we can easily find a significant result in one random sample and not find a significant result in another. This implies that $p$-values found in one study, albeit very low, need to be put into perspective by providing a measure of their variability.

Unfortunately, researchers usually face but one random realization and ignore the true effect and consequently the $p$-value's variability and the power (Halsey et al. 2015). In each of the 10,000 random realizations, we would estimate a different $\hat{\beta}_1$ and a different $p$-value. If we naïvely assumed that the coefficient $\hat{\beta}_1$ that we happened to estimate from one data set were

---

[12] Power is a zeroth order (lower) partial moment of the $p$-value distribution over replications. As such, it contains only a part of the distributional information. This part, however, is sufficient if one confines oneself to the dichotomous "statistically significant" vs. "not statistically significant" distinction; i.e., "statistical power quantifies the repeatability of the $P$ value, but only in terms of the either-or interpretation" (Halsey et al. 2015: 180).

true, we could estimate the power. Figure 3 shows that under this assumption power can be directly computed from a given $p$-value. We exemplarily illustrate the case for $p = 0.01$ (corresponding to a test score $Z_{Score}^{H0} = 2.33$). Assuming a one-sided significance test at the 0.05 level, we know the critical test value $Z^* = Z_{0.95}^{H0} = 1.64$. Presuming $\beta_1 = \hat{\beta}_1$ (alternative hypothesis $H1$), we can estimate $\widehat{Power} = 1 - F^{H1}(1.64)$ without computationally resorting to the estimated $\hat{\beta}_1$ and standard error ($SE_1$).

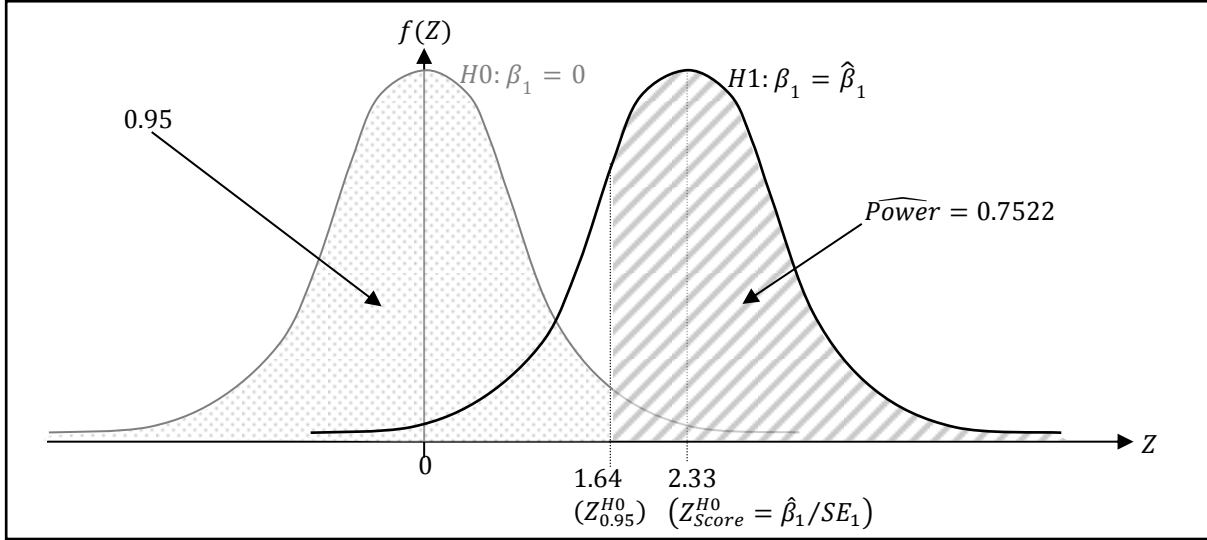Figure 3: Estimated power for $p = 0.01$ based on the assumption $\beta_1 = \hat{\beta}_1$



Table 6 tabulates selected $p$-values along with the estimated coefficients $\hat{\beta}_1$, the cumulated $p$-value densities $\Phi$, and the estimated power for the simulated $e \sim N(0; 3)$ and $e \sim N(0; 5)$ realities. While low $p$-values are often associated with a high reliability of estimated effects, there is an important message to be learned from Table 6: an unbiased estimator (in our case the OLS-estimator) estimates correctly *on average*. We indeed find in both realities an average coefficient across all 10,000 simulation runs that is *very* close to 0.2. We overestimate the effect size, however, in the case of highly significant results. For $e \sim N(0; 5)$, for example, $p = 0.00001$ corresponds to an estimated coefficient $\hat{\beta}_1 = 0.491$. Considering only the results that are significant at the 0.05 level, we find an average coefficient estimate of 0.2565. This is not surprising: by averaging over "significant" results only, we have right-truncated the distribution of the $p$-value which, in turn, involves a left-truncation of the distribution of the coefficient.

Table 6: Estimated coefficients and power for selected *p*-values

| | $e\sim N(0;3)$ | | $e\sim N(0;5)$ | | $Z_{Score}^{H0}$ $=\hat{\beta}_1/SE_1$ | $\widehat{Power}$ $=1-F_{H1}(Z^*)$ |
|---|---|---|---|---|---|---|
| | $\hat{\beta}_1$ [a] | $\Phi(p)$ | $\hat{\beta}_1$ | $\Phi(p)$ | | |
| $p=0.00001$ | 0.224 | 0.1209 | 0.491 | 0.0085 | 4.2649 | 0.9956 |
| $p=0.0001$ | 0.248 | 0.2831 | 0.317 | 0.0330 | 3.7190 | 0.9810 |
| $p=0.001$ | 0.174 | 0.5449 | 0.304 | 0.1230 | 3.0902 | 0.9258 |
| $p=0.01$ | 0.162 | 0.8290 | 0.249 | 0.3729 | 2.3263 | 0.7522 |
| $p=0.02$ | 0.138 | 0.8928 | 0.198 | 0.4800 | 2.0537 | 0.6587 |
| $p=0.03$ | 0.126 | 0.9240 | 0.192 | 0.5572 | 1.8808 | 0.5933 |
| $p=0.04$ | 0.106 | 0.9423 | 0.157 | 0.6057 | 1.7507 | 0.5421 |
| $p=0.05$ | 0.112 | 0.9529 | 0.170 | 0.6486 | 1.6449 | 0.5000 |
| $p=0.1$ | 0.083 | 0.9823 | 0.127 | 0.7718 | 1.2816 | 0.3582 |
| $p=0.2$ | 0.054 | 0.9952 | 0.088 | 0.8810 | 0.8416 | 0.2109 |

[a] Lower *p*-values go hand in hand with larger coefficient estimates in general. But $\hat{\beta}_1$ does not exactly follow the inverse order of the *p*-value because both the mean and the variance may vary from one random sample to the next; i.e., we may have different combinations of $\hat{\beta}_1$ and $SE_1$ resulting in identical *p*-values.

Besides the fact that the *p*-value does not provide a clear rationale or even formal calculus for statistical inference (Goodman 2008), its variability over replications undermines its already weak informative value. Even if there were *no* uncorrected multiple testing, *no* misinterpretation of the *p*-value, and *no* distorting selection procedures, researchers might still be overconfident and believe that at least a very low *p*-value of, let's say, 0.001 provides a trustworthy indication of the true effect. Alas, we cannot deduce even that much by just looking at the *p*-value. The consequences for conventional significance testing are quite sobering: if we are oblivious to the *p*-value's unknown sample-to-sample variability, we will grossly overestimate its limited informational content in a single study; and if we account for the *p*-value's sample-to-sample variability, we must concede that its suitability to indicate the strength of evidence is *very* limited. Let us briefly summarize the main reasons behind this sobering insights:

1. The variability of the *p*-value over random replications may be high or low. Being reduced to having to analyze data from one random realization, we do not know the degree of variation. We would need some prior assumption regarding the true effect size.

2. In plausible constellations of noise and sample size, we can very easily find a significant result in one random sample and not find a significant result in another.

3. The variability of the *p*-value is paralleled by the variability of the estimated coefficient. We may thence find a large coefficient in in one random sample and a small one in another.

4. Unbiased estimators estimate correctly *on average*, but we have no way of identifying the *p*-value below which (above which) we overestimate (underestimate) the effect size. In our example, a *p*-value of 0.001 may be associated with a coefficient estimate of 0.174 (= underestimation of the effect size) or a with a coefficient estimate of 0.304 (= overestimation).

5. We may grossly overestimate the effect size in the case of a highly significant result. That is, even in such a case, we cannot make a direct inference regarding the effect.

6. If we rashly claimed a just-estimated coefficient to be true, we would not have to be worried if it cannot be replicated. For example, if an effect associated with a *p*-value of 0.05

were real, we would *necessarily* have a mere 50% probability of finding a statistically significant effect in replications (one-sided test). Things improve with lower *p*-values. But even at the 0.01 level, we have only a 75% probability of re-finding significance (see Table 6).

The most important message to be learned is that the *p*-value is not the probability of a hypothesis but a measure that indicates how (in)compatible the particular data at hand (sample) is with a specified statistical model including the null hypothesis (Wasserstein and Lazar 2016). While the *p*-value is only a statement about a data set *conditional* on the null hypothesis of no effect, small *p*-values nonetheless give a hint that there might rather be an effect than none, in the sense that small *p*-values will occur more often if there is an effect compared to no effect (cf., the ideal-type *p*-value distributions without publication bias in Figure 1). However, due to the *p*-value's data dependence and the imponderables of random sampling (particularly in the case of small samples and big noise), large *p*-values may occur in a considerable fraction of random sampling replications even if there is an effect (cf., Figure 2), and vice versa. Unfortunately, in real-world applications we are usually not able to draw conclusions based on multiple random samples of the same population. Hence, we must be cautious when interpreting *p*-values and realize that "by itself, a *p*-value does not provide a good measure of evidence regarding a model or hypothesis" (Wasserstein and Lazar 2016: 132). When trying to assess the evidence from a particular data set regarding a hypothesis, statistical power calculations (or more generally, the consideration of the *p*-value's sample-to-sample variability) would be a helpful complement to the conventional *p*-value approach. However, since the size of the true effect is unknown, we are limited to variant power calculations for varying but plausible effect sizes. Assuming such plausible effect sizes, in turn, requires some degree of prior knowledge.

Furthermore, mixing up the *p*-value concept ("null hypothesis significance testing") and the estimation of effect size in the same step is problematic since many of the best estimation procedures are based on the concept of unbiasedness. Claiming the effect size to be real that we happened to estimate out of a sample where the effect showed up as "significant" bears the risk of overestimating the effect. Unfortunately, with decreasing *p*-values this risk seems to increase. When samples are large enough, one might thence think about splitting the data at hand and using one part for the testing part of the analysis and the other one for the estimation part. This might in further perspective lead to cross-validation-like approaches or to resampling procedures. But at this point we can only postpone such thoughts to further research.

## 4    Conclusion and outlook

The goal of this paper was to provide a non-technical and easily accessible resource for statistical practitioners who wish to spot and avoid misinterpretations and misuses of statistical significance tests. Our main insights are as follows: *first*, as empirical economists, we are naturally interested in statistical inference, i.e., we want to draw inductive conclusions and make general propositions about regularities in economic life given the evidence from the data. *Second*, we are also naturally interested in assessing the trustworthiness of these inductive conclusions by assigning epistemic probabilities to our propositions (hypotheses). *Third*, even though the label "hypotheses testing" is commonly attached to statistical significance testing, *p*-values cannot be used to test (the trustworthiness or probability of) hypotheses because they are not the probabilities that can be assigned to hypotheses given the analyzed data. Instead, *p*-values describe the conditional probability of the data given an assumed (null) hypothesis. *Fourth*, besides being per se a poor tool for assessing the trustworthiness of inductive generalizations from the particular (sample) to the general (population), *p*-values and associated effect size estimates may exhibit a wide sample-to-sample variability. We may easily find a small *p*-value and a large effect in one random sample and a large *p*-value and a small effect in another. *Fifth*, we rarely start from scratch. Besides the evidence provided by our own data,

we would need to summarize the knowledge from prior studies and use Bayesian statistics if we wanted to assign probabilities to our scientific propositions.

We must realize that statistical inference is not as trivial as the dichotomous "significance" vs. "nonsignificance" declarations suggest at first glance. The corresponding either-or interpretations regarding the trustworthiness and importance of estimated effects are neat and seemingly plausible but wrong. What is more, misuses such as disregard of evident multiple testing and $p$-hacking as well as selective publishing may inflate statistical significance claims and distort the published body of evidence. Many reformatory measures have been suggested in the literature to mitigate these problems. The most obvious one is the call for an increase of statistical literacy through better teaching. Others include changes of research standards and incentives that would reduce the publish-statistically-significant-results-or-perish pressure that still dominates many disciplinary cultures. The most notable examples of reform are increasingly adopted policies of preregistration, replication, sharing (of data and analytical protocols), and unbiased publishing of both positive and negative results.

While some of these policies are being slowly introduced into economics, there are a number of specific concerns that need to be addressed before approaches from other fields can be successfully transplanted to (non-experimental) economic research. These concerns stem from the prevalence of observational studies and multiple regressions, and the data-driven specification search that characterizes most econometric analyses. A fundamental question in this context is whether we can draw a dividing line between the search for an adequate model specification (among many) and $p$-hacking. We would also need to ask where on the continuum between strictly theory-based modeling and purely data-driven modeling econometric study changes its nature from being confirmatory to being exploratory to being merely a description of the idiosyncratic structure (including noise) of the data set at hand (overfitting). And by which criteria should the latter be avoided and a study be declared either confirmatory or exploratory? Furthermore, if a study can be declared confirmatory, how should we adjust for multiple testing given the fact that any data-based model contains more variables than just the initial focal predictors and that is has been retained as the final model after an often large number of different model specifications have been tested.

Even widely applauded measures aimed at improving research practices need to be scrutinized regarding their viability and effectiveness in econometric research. It is not clear, for example, how measures such as the preregistration of analytical designs, the replication of studies, and the correction for multiple testing would have to look like within a scientific culture where it is common practice to specify statistical models that fit the data. Related to that, the question arises of how to meet the requirement of considering the body of evidence instead of focusing on single studies. How can we assess the body of evidence regarding a specific scientific issue when comparability across studies is impeded because there are often (nearly) as many data-dependent model specifications as there are studies? Another question is how, in view of the $p$-value's inferential limitations, we can provide an interpretative orientation vis-à-vis the often large numbers of regression coefficient estimates. This is an especially urging issue since data-based economic models are often heavily populated, besides the original variables of interest, by interaction terms, (log)transformed variables, lagged variables, instrumental variables, higher-order polynomials, and control variables. Can Bayesian approaches provide a practically feasibly solution in such a context? And if so, how can we specify Bayesian priors given the regularly lacking comparability of studies and the often large variable sets?

**References**

Altman, N., Krzywinski, M. (2017): P values and the search for significance. Nature Methods 14(1): 3-4.

Armstrong, J.S. (2007): Significance tests harm progress in forecasting. International Journal of Forecasting 23(2): 321-327.

Auspurg, K., Hinz, T. (2011): What Fuels Publication Bias? Theoretical and Empirical Analyses of Risk Factors Using the Caliper Test. Journal of Economics and Statistics 231(5-6): 636-660.

Baker, M. (2016): Statisticians issue warning on *P* values. Nature 531(7593): 151.

Becker, B.J., Wu, M-J. (2007): The Synthesis of Regression Slopes in Meta-Analysis. Statistical Science 22(3): 414-429.

Benjamini, Y., Hochberg, Y. (1995): Controlling the false discovery rate: A practical and powerful approach to multiple testing. Journal of the Royal Statistical Society B 57(1): 289-300.

Bennett, D.A., Latham, N.K., Stretton, C., Anderson, C.S. (2004): Capture-recapture is a potentially useful method for assessing publication bias. Journal of Clinical Epidemiology 57(4): 349-357.

Berning, C., Weiß, B. (2016): Publication Bias in the German Social Sciences: An Application of the Caliper Test to Three Top-Tier German Social Science Journals. Quality & Quantity 50(2): 901-917.

Berry, D.A. (2016): P-Values Are Not What They're Cracked Up to Be. Online Discussion: ASA Statement on Statistical Significance and P-values. The American Statistician 70(2): 1-2.

Boos, D.D., Stefanski, L.A. (2011): P-Value Precision and Reproducibility. The American Statistician, 65(4): 213-221.

Borenstein, M., Hedges, L.V., Higgins, J.P.T., Rothstein, H.R. (2009): Introduction to Meta-Analysis. Chichester: John Wiley & Sons.

Bretz, F., Hothorn, T., Westfall, P. (2010): Multiple comparisons using R. Boca Raton: CRC Press.

Brodeur, A., Lé, M., Sangnier, M., Zylberberg, Y. (2016): Star Wars: The Empirics Strike Back. American Economic Journal: Applied Economics 8(1): 1-32.

Bruns, S.B. (2017): Meta-Regression Models and Observational Research. Oxford Bulletin of Economics and Statistics 0305–9049, doi: 10.1111/obes.12172.

Card, D., Krueger, A. B. (1995): Time-series minimum-wage studies: A meta-analysis. American Economic Review (AEA Papers and Proceedings) 85: 238-243.

Card, N. A. (2012): Applied meta-analysis for social science research. New York: Guilford Press.

Cohen, J. (1994): The earth is round (p < 0.05). American Psychologist 49(12): 997-1003.

Colquhoun, D. (2014): An investigation of the false discovery rate and the misinterpretation of *p*-values. Royal Society Open Science 1:140216; http://dx.doi.org/10.1098/rsos.140216: 1-16.

Cooper, D.J., Dutcher, E.G. (2011): The dynamics of responder behavior in ultimatum games: a meta-study Experimental Economics 14(4): 519-546.

Cooper, H., Hedges, L., Valentine. J. (eds.) (2009): The handbook or research synthesis and meta-analysis. 2nd ed., Russell Sage Foundation, New York.

Crouch, G.I. (1995): A meta-analysis of tourism demand. Annals of tourism research 22(1): 103-118.

Cumming, G. (2008): Replication and p intervals: p values predict the future only vaguely, but confidence intervals do much better. Perspectives on Psychological Science 3(4): 286-300.

Denton, F.T. (1985): Data Mining as an Industry. Review of Economics and Statistics 67(1): 124-27.

Denton, F.T. (1988): The significance of significance: Rhetorical aspects of statistical hypothesis testing in economics. In: Klamer, A., McCloskey, D.N., Solow, R.M. (eds.): The consequences of economic rhetoric. Cambridge: Cambridge University Press: 163-193.

Didelez, V., Pigeot, I., Walter, P. (2006): Modifications of the Bonferroni-Holm procedure for a multi-way ANOVA. Statistical Papers 47: 181-209.

Duvendack, M., Palmer-Jones, R., Reed, W.R. (2015): Replications in Economics: A Progress Report. Econ Journal Watch 12(2): 164-191.

Duvendack, M., Palmer-Jones, R., Reed, W.R. (2017): What Is Meant by "Replication" and Why Does It Encounter Resistance in Economics? American Economic Review: Papers & Proceedings 2017: 107(5): 46-51.

Egger, M., Smith, G.D., Schneider, M., Minder, C. (1997): Bias in meta-analysis detected by a simple, graphical test. British Medical Journal 315 (7109): 629-634.

Engel, C. (2011): Dictator games: a meta study. Experimental Economics 14(4): 583-610.

Evanschitzky, H., Armstrong, J.S. (2010): Replications of forecasting research. International Journal of Forecasting 26: 4-8.

Fanelli, D. (2010): Positive'' results increase down the hierarchy of the sciences. PLoS One 5 (4): e10068.

Fanelli, D. (2011): Negative results are disappearing from most disciplines and countries. Scientometrics 90(3): 891-904.

Fisher, R.A. (1935): The design of experiments. Edinburgh: Oliver & Boyd.

Fitzpatrick, L., Parmeter, C.F., Agar, J. (2017): Threshold Effects in Meta-Analyses With Application to Benefit Transfer for Coral Reef Valuation. Ecological Economics 133: 74-85.

Gerber, A. S., N. Malhotra (2008): Publication Bias in Empirical Sociological Research. Do Arbitrary Significance Levels Distort Published Results? Sociological Methods & Research 37(1): 3-30.

Gerber, A.S., Malhotra, N., Dowling, C.M., Doherty, D. (2010): Publication Bias in Two Political Behavior Literatures. American Politics Research 38(4): 591-613.

Gigerenzer, G., Krauss, S., Vitouch, O. (2004): The null ritual: what you always wanted to know about significance testing but were afraid to ask. In: Kaplan, D. (ed.): The SAGE handbook of quantitative methodology for the social sciences (Chapter 21). Thousand Oaks: Sage.

Goodman, S. (2008): A dirty dozen: Twelve *p*-value Misconceptions. Seminars in Hematology 45: 135-140.

Goodman, S.N. (1992): A Comment of Replication, P-Values and Evidence. Statistics in Medicine 11: 875-879.

Greenland, S., Senn, S.J., Rothman, K.J., Carlin, J.B., Poole, C., Goodman, S.N., Altman, D.G. (2016): Statistical tests, P values, confidence intervals, and power: a guide to misinterpretations. European Journal of Epidemiology 31(4): 337-350.

Haller, H., Krauss, S. (2002): Misinterpretations of Significance: A Problem Students Share with Their Teachers? Methods of Psychological Research Online 7(1): 1-20.

Halsey, L.G., Curran-Everett, D., Vowler, S.L., Drummond, B. (2015): The fickle P value generates irreproducible results. Nature Methods 12(3): 179-185.

Hartung, J., Knapp, G., Sinha, B.K. (2008): Statistical Meta-Analysis with Applications. Hoboken: John Wiley & Sons.

Head, M.L, Holman, L., Lanfear, R., Kahn, A.T., Jennions, M.D. (2015): The Extent and Consequences of P-Hacking in Science. PLoS Biology 13(3): e1002106. doi:10.1371/journal.pbio.1002106.

Hirschauer, N., Mußhoff, O., Grüner, S., Frey, U., Theesfeld, I., Wagner, P. (2016): Inferential misconceptions and replication crisis. Journal of Epidemiology, Biostatistics, and Public Health 13(4): e12066-1-e12066-16.

Hochberg, Y., Tamhane, A.C. (1987). Multiple comparison procedures. New York: Wiley.

Holm, S. (1979). A simple sequentially rejective multiple test procedure. Scandinavian Journal of Statistics 6(2): 65-70.

Howard, G.S., Maxwell, S.E., Fleming, K.J. (2000): The proof of the pudding: An illustration of the relative strengths of null hypothesis, meta-analysis, and Bayesian analysis. Psychological Methods 5: 315-332.

Ioannidis, J., Doucouliagos, C. (2013): What's to know about the credibility of empirical economics? Journal of Economic Surveys 27(5): 997-1004.

Ioannidis, J.P.A. (2005): Why Most Published Research Findings are False. PLoS Medicine 2(8): e124: 0696-0701.

Joober, R., Schmitz, N., Dipstat, L.A., Boksa, P. (2012): Publication bias: What are the challenges and can they be overcome? Journal of Psychiatry & Neuroscience: 37(3): 149-152.

Kerr, N.L. (1998): HARKing: Hypothesizing after the results are known. Personality and Social Psychology Review 2(3): 196-217.

Kicinski, M, Springate, D.A., Kontopantelis, E. (2015): Publication bias in meta-analyses from the Cochrane Database of Systematic Reviews. Statistics in Medicine 34: 2781-2793.

Kline, R.B. (2013): Beyond Significance Testing: Statistics Reform in the Behavioral Sciences. Washington: American Psychological Association.

Krämer, W. (2011): The Cult of Statistical Significance – What Economists Should and Should Not Do to Make their Data Talk. Schmollers Jahrbuch 131(3): 455-468.

Lange, T. (2016): Discrimination in the laboratory: A meta-analysis of economics experiments. European Economic Review 90: 375-402.

Leamer, E.E. (1978): Specification Searches: Ad Hoc Inference with Nonexperimental Data. New York, Wiley.

Lecoutre, B., Poitevineau, J. (2014): The Significance Test Controversy Revisited. The Fiducial Bayesian Alternative. Heidelberg: Springer.

Light, R.J., Pillemer, D.B. (1984): Summing Up: The Science of Reviewing Research. Cambridge: Harvard University Press.

List, J.A., Shaikh, A.M., Xu, Y. (2016): Multiple Hypothesis Testing in Experimental Economics. No. w21875. National Bureau of Economic Research, Working Paper No. 21875.

Loomis, J.B., White, D.S. (1996): Economic benefits of rare and endangered species: summary and meta-analysis. Ecological Economics 18(3): 197-206.

Lovell, M.C. (1983). Data Mining. Review of Economics and Statistics 65(1): 1-12.

McCloskey, D.N., Ziliak, S.T. (1996): The Standard Error of Regressions. Journal of Economic Literature 34(1): 97-114.

Motulsky, J.J. (2014): Common Misconceptions about Data Analysis and Statistics. The Journal of Pharmacology and Experimental Theurapeutics 351(8): 200-205.

Munafò, M.R., Nosek, B.A., Bishop, D.V.M., Button, K.S., Chambers, C.D., du Sert, N.P., Simonsohn, U., Wagenmakers, E-J., Ware, J.J., Ioannidis, J.P.A. (2017): A manifesto for reproducible science. Nature Human Behaviour 1(0021): 1-8.

Nickerson, R.S. (2000): Null hypothesis significance testing: A review of an old and continuing controversy. Psychological Methods 5(2): 241-301.

Nuzzo, R. (2014): Statistical Errors. *P*-values, the 'gold standard' of statistical validity, are not as reliable as many scientists assume. Nature 506(7487): 150-152.

Oakes, M. (1986): Statistical inference: A commentary for the social and behavioural sciences. New York: Wiley.

Pigeot, I. (2000): Basic concepts of multiple tests – A survey. Invited paper. Statistical Papers 41: 3-36.

Pitchforth, J.O., Mengersen, K.L. (2013): Bayesian Meta-Analysis. In: Alston, C.L., Mengersen, K.L., Pettitt, A.N. (eds.): Case Studies in Bayesian Statistical Modelling and Analysis. Chichester: John Wiley & Sons, Ltd., 118-140.

Poorolajal, J., Haghdoost, A.A., Mahmoodi, M., Majdzadeh, R., Nasseri-Moghaddam, S., Fotouhi, A. (2010): Capture-recapture method for assessing publication bias. Journal of Research in Medical Sciences : The Official Journal of Isfahan University of Medical Sciences 15(2): 107-115.

Roberts, C.J. (2005): Issues in meta-regression analysis: An overview. Journal of Economic Surveys 19(3): 295-298.

Rosenberg, M.S. (2005): The File-drawer Problem Revisited: A General Weighted Method for Calculating Fail-Safe Numbers in Meta-Analysis. Evolution 59(2): 464-468.

Rosenthal, R. (1979): The file drawer problem and tolerance for null results. Psychological Bulletin 86(3): 638-641.

Rothstein, H., Sutton, A.J., Borenstein, M. (2005): Publication Bias in Meta-Analysis. Prevention, Assessment and Adjustments. Sussex: Wiley.

Schmidt, F.L., Hunter, J.E. (2014): Methods of meta-analysis: Correcting error and bias in research findings. Los Angeles: Sage publications.

Silliman, N. (1997): Hierarchical selection models with applications in meta-analysis. Journal of American Statistical Association 92(439): 926-936.

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.

Simonsohn, U., Nelson, L.D., Simmons, J.P. (2014): *P*-Curve: A Key to the File-Drawer. Journal of Experimental Psychology 143(2): 534-547.

Smith, M.L. (1980): Publication bias and meta-analysis. Evaluation in Education 4: 22-24.

Song, F., Eastwood, A.J., Gilbody, S., Duley, L., Sutton, A.J. (2000): Publication and related biases. Southampton: The National Coordinating Centre for Health Technology Assessment.

Song, F., Hooper, L., Loke, Y.K. (2013): Publication bias: what is it? How do we measure it? How do we avoid it? Open Access Journal of Clinical Trials 5: 71-81.

Stanley, T.D., Jarrell, S. B. (1989): Meta-regression analysis: A quantitative method of literature surveys. Journal of Economic Surveys 3(2): 161-170.

Stanley, T.D., Doucouliagos, H. (2012): Meta-Regression Analysis in Economics and Business. London: Routledge.

Sterling, T.D. (1959): Publication Decisions and their Possible Effects on Inferences Drawn from Tests of Significance–Or Vice Versa. Journal of the American Statistical Association 54(285): 30-34.

Sterne, J.A.C., Egger, M. (2005): Regression Methods to Detect Publication and Other Bias in Meta-Analysis. In: Rothstein, H.R., Sutton, A.J., Borenstein, M. (eds.): Publication Bias in Meta-Analysis. Prevention, Assessment and Adjustments. Chichester: Wiley: 99-110.

Sterne, J.A.C., Egger, M., Moher, D. (2008): Addressing reporting biases. In: Higgins, J.P.T., Green, S. (Eds.): Cochrane handbook for systematic reviews of interventions: 297-333. Chichester: Wiley.

Van Houtven, G.L., Pattanayak, S.K., Usmani, F., Yang, J.C. (2017): What are Households Willing to Pay for Improved Water Access? Results from a Meta-Analysis. Ecological Economics 136: 126-135.

Wasserstein, R.L., Lazar N.A. (2016): The ASA's statement on p-values: context, process, and purpose, The American Statistician 70(2): 129-133.

Weiß, B., Wagner, M. (2011): The identification and prevention of publication bias in the social sciences and economics. Jahrbücher für Nationalökonomie und Statistik 231(5-6): 661-684.

Westfall, P., Tobias, R., Wolfinger, R. (2011): Multiple comparisons and multiple testing using SAS. Cary: SAS Institute.

Zelmer, J. (2003): Linear public goods experiments: A meta-analysis. Experimental Economics 6(3): 299-310.

Ziliak, S.T., McCloskey, D.N. (2008): The Cult of Statistical Significance. How the Standard Error Costs Us Jobs, Justice, and Lives. Ann Arbor: The University of Michigan Press.

Zyphur, M.J., Oswald, F.L. (2015): Bayesian Estimation and Inference: A User's Guide. Bayesian Probability and Statistics in Management Research, (Special Issue of the) Journal of Management 41(2): 390-420.