



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

## Null hypothesis significance testing: a short tutorial

**Citation for published version:**

Pernet, C 2017, 'Null hypothesis significance testing: a short tutorial', *F1000Research*, vol. 4, pp. 621.  
<https://doi.org/10.12688/f1000research.6963.3>

**Digital Object Identifier (DOI):**

[10.12688/f1000research.6963.3](https://doi.org/10.12688/f1000research.6963.3)

**Link:**

[Link to publication record in Edinburgh Research Explorer](#)

**Document Version:**

Other version

**Published In:**

F1000Research

**Publisher Rights Statement:**

Copyright: © 2016 Pernet C. This is an open access article distributed under the terms of the Creative Commons Attribution Licence, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

**General rights**

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

**Take down policy**

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact [openaccess@ed.ac.uk](mailto:openaccess@ed.ac.uk) providing details, and we will remove access to the work immediately and investigate your claim.





OPINION ARTICLE

**REVISED** Null hypothesis significance testing: a short tutorial  
 [version 3; referees: 1 approved with reservations, 2 not approved]

Cyril Pernet

Centre for Clinical Brain Sciences (CCBS), Neuroimaging Sciences, The University of Edinburgh, Edinburgh, UK

**v3** **First published:** 25 Aug 2015, 4:621 (doi: [10.12688/f1000research.6963.1](https://doi.org/10.12688/f1000research.6963.1))  
**Second version:** 13 Jul 2016, 4:621 (doi: [10.12688/f1000research.6963.2](https://doi.org/10.12688/f1000research.6963.2))  
**Latest published:** 10 Oct 2016, 4:621 (doi: [10.12688/f1000research.6963.3](https://doi.org/10.12688/f1000research.6963.3))

**Abstract**

Although thoroughly criticized, null hypothesis significance testing (NHST) remains the statistical method of choice used to provide evidence for an effect, in biological, biomedical and social sciences. In this short tutorial, I first summarize the concepts behind the method, distinguishing test of significance (Fisher) and test of acceptance (Newman-Pearson) and point to common interpretation errors regarding the p-value. I then present the related concepts of confidence intervals and again point to common interpretation errors. Finally, I discuss what should be reported in which context. The goal is to clarify concepts to avoid interpretation errors and propose reporting practices.

**Open Peer Review**

**Referee Status:**

	Invited Referees		
	1	2	3
<b>REVISED</b> version 3 published 10 Oct 2016			
<b>REVISED</b> version 2 published 13 Jul 2016			report
version 1 published 25 Aug 2015	report	report	

- 1 **Daniel Lakens**, Eindhoven University of Technology Netherlands
- 2 **Marcel ALM van Assen**, Tilburgh University Netherlands
- 3 **Stephen J. Senn**, Luxembourg Institute of Health Luxembourg

**Discuss this article**

Comments (0)

**Corresponding author:** Cyril Pernet ([cyril.pernet@ed.ac.uk](mailto:cyril.pernet@ed.ac.uk))

**How to cite this article:** Pernet C. **Null hypothesis significance testing: a short tutorial [version 3; referees: 1 approved with reservations, 2 not approved]** *F1000Research* 2016, 4:621 (doi: [10.12688/f1000research.6963.3](https://doi.org/10.12688/f1000research.6963.3))

**Copyright:** © 2016 Pernet C. This is an open access article distributed under the terms of the [Creative Commons Attribution Licence](#), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

**Grant information:** The author(s) declared that no grants were involved in supporting this work.

**Competing interests:** No competing interests were disclosed.

**First published:** 25 Aug 2015, 4:621 (doi: [10.12688/f1000research.6963.1](https://doi.org/10.12688/f1000research.6963.1))

**REVISED Amendments from Version 2**

This version 3 includes minor changes that reflect the 3rd reviewers' comments - in particular the theoretical vs. practical difference between Fisher and Neyman-Pearson. Additional information and reference is also included regarding the interpretation of p-value for low powered studies.

See referee reports

**The Null Hypothesis Significance Testing framework**

NHST is a method of statistical inference by which an experimental factor is tested against a hypothesis of no effect or no relationship based on a given observation. The method is a combination of the concepts of significance testing developed by Fisher in 1925 and of acceptance based on critical rejection regions developed by Neyman & Pearson in 1928. In the following I am first presenting each approach, highlighting the key differences and common misconceptions that result from their combination into the NHST framework (for a more mathematical comparison, along with the Bayesian method, see Christensen, 2005). I next present the related concept of confidence intervals. I finish by discussing practical aspects in using NHST and reporting practice.

**Fisher, significance testing, and the p-value**

The method developed by (Fisher, 1934; Fisher, 1955; Fisher, 1959) allows to compute the probability of observing a result at least as extreme as a test statistic (e.g. t value), assuming the null hypothesis of no effect is true. This probability or p-value reflects (1) the conditional probability of achieving the observed outcome or larger:  $p(\text{Obs} \geq t | H_0)$ , and (2) is therefore a cumulative probability rather than a point estimate. It is equal to the area under the null probability distribution curve from the observed test statistic to the tail of the null distribution (Turkheimer *et al.*, 2004). The approach proposed is of 'proof by contradiction' (Christensen, 2005), we pose the null model and test if data conform to it.

In practice, it is recommended to set a *level of significance* (a theoretical p-value) that acts as a reference point to identify significant results, that is to identify results that differ from the null-hypothesis of no effect. Fisher recommended using  $p=0.05$  to judge whether an effect is significant or not as it is roughly two standard deviations away from the mean for the normal distribution (Fisher, 1934 page 45: 'The value for which  $p=.05$ , or 1 in 20, is 1.96 or nearly 2; it is convenient to take this point as a limit in judging whether a deviation is to be considered significant or not'). A key aspect of Fisher's theory is that only the null-hypothesis is tested, and therefore p-values are meant to be used in a graded manner to decide whether the evidence is worth additional investigation and/or replication (Fisher, 1971 page 13: 'it is open to the experimenter to be more or less exacting in respect of the smallness of the probability he would require [...] and 'no isolated experiment, however significant in itself, can suffice for the experimental demonstration of any natural phenomenon'). How small the level of significance is, is thus left to researchers.

**What is not a p-value? Common mistakes**

The p-value is not an indication of the strength or magnitude of an effect. Any interpretation of the p-value in relation to the effect

under study (strength, reliability, probability) is wrong, since p-values are conditioned on  $H_0$ . In addition, while p-values are randomly distributed (if all the assumptions of the test are met) when there is no effect, their distribution depends of both the population effect size and the number of participants, making impossible to infer strength of effect from them.

Similarly,  $1-p$  is not the probability to replicate an effect. Often, a small value of p is considered to mean a strong likelihood of getting the same results on another try, but again this cannot be obtained because the p-value is not informative on the effect itself (Miller, 2009). Because the p-value depends on the number of subjects, it can only be used in high powered studies to interpret results. In low powered studies (typically small number of subjects), the p-value has a large variance across repeated samples, making it unreliable to estimate replication (Halsey *et al.*, 2015).

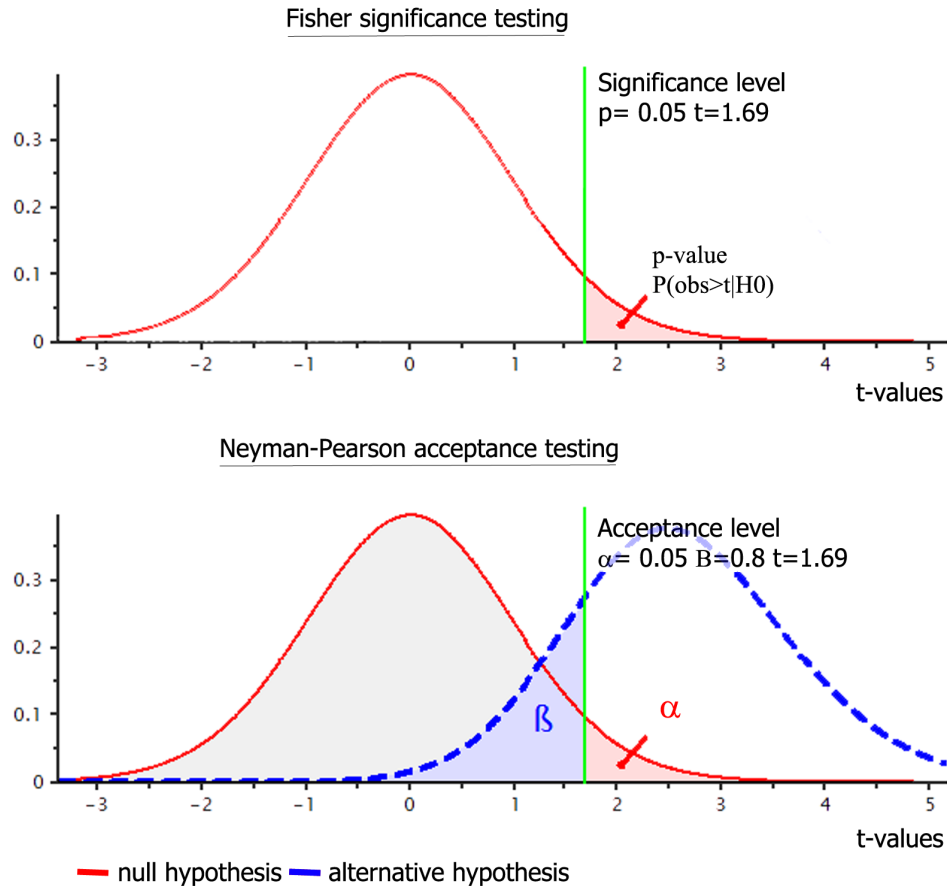
A (small) p-value is not an indication favouring a given hypothesis. Because a low p-value only indicates a misfit of the null hypothesis to the data, it cannot be taken as evidence in favour of a specific alternative hypothesis more than any other possible alternatives such as measurement error and selection bias (Gelman, 2013). Some authors have even argued that the more (a priori) implausible the alternative hypothesis, the greater the chance that a finding is a false alarm (Krzywinski & Altman, 2013; Nuzzo, 2014).

The p-value is not the probability of the null hypothesis  $p(H_0)$ , of being true, (Krzywinski & Altman, 2013). This common misconception arises from a confusion between the probability of an observation given the null  $p(\text{Obs} \geq t | H_0)$  and the probability of the null given an observation  $p(H_0 | \text{Obs} \geq t)$  that is then taken as an indication for  $p(H_0)$  (see Nickerson, 2000).

**Neyman-Pearson, hypothesis testing, and the  $\alpha$ -value**

Neyman & Pearson (1933) proposed a framework of statistical inference for applied decision making and quality control. In such framework, two hypotheses are proposed: the null hypothesis of no effect and the alternative hypothesis of an effect, along with a control of the long run probabilities of making errors. The first key concept in this approach, is the establishment of an *alternative hypothesis* along with an a priori effect size. This differs markedly from Fisher who proposed a general approach for scientific inference conditioned on the null hypothesis only. The second key concept is the *control of error rates*. Neyman & Pearson (1928) introduced the notion of critical intervals, therefore dichotomizing the space of possible observations into correct vs. incorrect zones. This dichotomization allows distinguishing correct results (rejecting  $H_0$  when there is an effect and not rejecting  $H_0$  when there is no effect) from errors (rejecting  $H_0$  when there is no effect, the type I error, and not rejecting  $H_0$  when there is an effect, the type II error). In this context, alpha is the probability of committing a Type I error in the long run. Alternatively, Beta is the probability of committing a Type II error in the long run.

The (theoretical) difference in terms of hypothesis testing between Fisher and Neyman-Pearson is illustrated on Figure 1. In the 1<sup>st</sup> case, we choose a level of significance for observed data of 5%, and compute the p-value. If the p-value is below the level of significance, it is used to reject  $H_0$ . In the 2<sup>nd</sup> case, we set a critical interval



**Figure 1. Illustration of the difference between the Fisher and Neyman-Pearson procedures.** The figure was prepared with G-power for a one-sided one-sample t-test, with a sample size of 32 subjects, an effect size of 0.45, and error rates  $\alpha=0.049$  and  $\beta=0.80$ . In Fisher's procedure, only the null hypothesis is posed, and the observed p-value is compared to an a priori level of significance. If the observed p-value is below this level (here  $p=0.05$ ), one rejects  $H_0$ . In Neyman-Pearson's procedure, the null and alternative hypotheses are specified along with an a priori level of acceptance. If the observed statistical value is outside the critical region (here  $[-\infty + 1.69]$ ), one rejects  $H_0$ .

based on the a priori effect size and error rates. If an observed statistic value is below and above the critical values (the bounds of the confidence region), it is deemed significantly different from  $H_0$ . In the NHST framework, the level of significance is (in practice) assimilated to the alpha level, which appears as a simple decision rule: if the p-value is less or equal to alpha, the null is rejected. It is however a common mistake to assimilate these two concepts. The level of significance set for a given sample is not the same as the frequency of acceptance alpha found on repeated sampling because alpha (a point estimate) is meant to reflect the long run probability whilst the p-value (a cumulative estimate) reflects the current probability (Fisher, 1955; Hubbard & Bayarri, 2003).

#### Acceptance or rejection of $H_0$ ?

The acceptance level  $\alpha$  can also be viewed as the maximum probability that a test statistic falls into the rejection region when the null hypothesis is true (Johnson, 2013). Therefore, one can only reject the null hypothesis if the test statistics falls into the critical region(s), or fail to reject this hypothesis. In the latter case, all we can say is that no significant effect was observed, but one cannot conclude that the null hypothesis is true. This is another

common mistake in using NHST: there is a profound difference between accepting the null hypothesis and simply failing to reject it (Killeen, 2005). By failing to reject, we simply continue to assume that  $H_0$  is true, which implies that one cannot argue against a theory from a non-significant result (absence of evidence is not evidence of absence). To accept the null hypothesis, tests of equivalence (Walker & Nowacki, 2011) or Bayesian approaches (Dienes, 2014; Kruschke, 2011) must be used.

#### Confidence intervals

Confidence intervals (CI) are builds that fail to cover the true value at a rate of alpha, the Type I error rate (Morey & Rouder, 2011) and therefore indicate if observed values can be rejected by a (two tailed) test with a given alpha. CI have been advocated as alternatives to p-values because (i) they allow judging the statistical significance and (ii) provide estimates of effect size. Assuming the CI (a)symmetry and width are correct (but see Wilcox, 2012), they also give some indication about the likelihood that a similar value can be observed in future studies. For future studies of the same sample size, 95% CI give about 83% chance of replication success (Cumming & Maillardet, 2006). If sample sizes however differ between studies, CI do not however warranty any a priori coverage.

Although CI provide more information, they are not less subject to interpretation errors (see [Savalei & Dunn, 2015](#) for a review). The most common mistake is to interpret CI as the probability that a parameter (e.g. the population mean) will fall in that interval X% of the time. The correct interpretation is that, for repeated measurements with the same sample sizes, taken from the same population, X% of times the CI obtained will contain the true parameter value ([Tan & Tan, 2010](#)). The alpha value has the same interpretation as testing against H<sub>0</sub>, i.e. we accept that 1-alpha CI are wrong in alpha percent of the times in the long run. This implies that CI do not allow to make strong statements about the parameter of interest (e.g. the mean difference) or about H<sub>1</sub> ([Hoekstra et al., 2014](#)). To make a statement about the probability of a parameter of interest (e.g. the probability of the mean), Bayesian intervals must be used.

### The (correct) use of NHST

NHST has always been criticized, and yet is still used every day in scientific reports ([Nickerson, 2000](#)). One question to ask oneself is what is the goal of a scientific experiment at hand? If the goal is to establish a discrepancy with the null hypothesis and/or establish a pattern of order, because both requires ruling out equivalence, then NHST is a good tool ([Frick, 1996](#); [Walker & Nowacki, 2011](#)). If the goal is to test the presence of an effect and/or establish some quantitative values related to an effect, then NHST is not the method of choice since testing is conditioned on H<sub>0</sub>.

While a Bayesian analysis is suited to estimate that the probability that a hypothesis is correct, like NHST, it does not prove a theory on itself, but adds its plausibility ([Lindley, 2000](#)). No matter what testing procedure is used and how strong results are, ([Fisher, 1959](#) p13) reminds us that ‘[...] no isolated experiment, however significant in itself, can suffice for the experimental demonstration of any natural phenomenon’. Similarly, the recent statement of the American Statistical Association ([Wasserstein & Lazar, 2016](#)) makes it clear that conclusions should be based on the researchers understanding of the problem in context, along with all summary data and tests, and that no single value (being p-values, Bayesian factor or else) can be used support or invalidate a theory.

### What to report and how?

Considering that quantitative reports will always have more information content than binary (significant or not) reports, we can always argue that raw and/or normalized effect size, confidence

intervals, or Bayes factor must be reported. Reporting everything can however hinder the communication of the main result(s), and we should aim at giving only the information needed, at least in the core of a manuscript. Here I propose to adopt optimal reporting in the result section to keep the message clear, but have detailed supplementary material. When the hypothesis is about the presence/absence or order of an effect, and providing that a study has sufficient power, NHST is appropriate and it is sufficient to report in the text the actual p-value since it conveys the information needed to rule out equivalence. When the hypothesis and/or the discussion involve some quantitative value, and because p-values do not inform on the effect, it is essential to report on effect sizes ([Lakens, 2013](#)), preferably accompanied with confidence or credible intervals. The reasoning is simply that one cannot predict and/or discuss quantities without accounting for variability. For the reader to understand and fully appreciate the results, nothing else is needed.

Because science progress is obtained by cumulating evidence ([Rosenthal, 1991](#)), scientists should also consider the secondary use of the data. With today’s electronic articles, there are no reasons for not including all of derived data: mean, standard deviations, effect size, CI, Bayes factor should always be included as supplementary tables (or even better also share raw data). It is also essential to report the context in which tests were performed – that is to report all of the tests performed (all t, F, p values) because of the increase type one error rate due to selective reporting (multiple comparisons and p-hacking problems - [Ioannidis, 2005](#)). Providing all of this information allows (i) other researchers to directly and effectively compare their results in quantitative terms (replication of effects beyond significance, [Open Science Collaboration, 2015](#)), (ii) to compute power to future studies ([Lakens & Evers, 2014](#)), and (iii) to aggregate results for meta-analyses whilst minimizing publication bias ([van Assen et al., 2014](#)).

### Competing interests

No competing interests were disclosed.

### Grant information

The author(s) declared that no grants were involved in supporting this work.

## References

- Christensen R: **Testing Fisher, Neyman, Pearson, and Bayes.** *The American Statistician*. 2005; **59**(2): 121–126.  
[Publisher Full Text](#)
- Cumming G, Maillardet R: **Confidence intervals and replication: Where will the next mean fall?** *Psychological Methods*. 2006; **11**(3): 217–227.  
[PubMed Abstract](#) | [Publisher Full Text](#)
- Dienes Z: **Using Bayes to get the most out of non-significant results.** *Front Psychol*. 2014; **5**: 781.  
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)

- Fisher RA: **Statistical Methods for Research Workers.** (Vol. 5th Edition). Edinburgh, UK: Oliver and Boyd. 1934.  
[Reference Source](#)
- Fisher RA: **Statistical Methods and Scientific Induction.** *Journal of the Royal Statistical Society, Series B*. 1955; **17**(1): 69–78.  
[Reference Source](#)
- Fisher RA: **Statistical methods and scientific inference.** (2nd ed.). New York: Hafner Publishing, 1959.  
[Reference Source](#)

Fisher RA: **The Design of Experiments**. Hafner Publishing Company, New-York. 1971.

[Reference Source](#)

Frick RW: **The appropriate use of null hypothesis testing**. *Psychol Methods*. 1996; **1**(4): 379–390.

[Publisher Full Text](#)

Gelman A: **P values and statistical practice**. *Epidemiology*. 2013; **24**(1): 69–72.

[PubMed Abstract](#) | [Publisher Full Text](#)

Halsey LG, Curran-Everett D, Vowler SL, *et al.*: **The fickle P value generates irreproducible results**. *Nat Methods*. 2015; **12**(3): 179–85.

[PubMed Abstract](#) | [Publisher Full Text](#)

Hoekstra R, Morey RD, Rouder JN, *et al.*: **Robust misinterpretation of confidence intervals**. *Psychon Bull Rev*. 2014; **21**(5): 1157–1164.

[PubMed Abstract](#) | [Publisher Full Text](#)

Hubbard R, Bayarri MJ: **Confusion over measures of evidence (p's) versus errors ( $\alpha$ 's) in classical statistical testing**. *The American Statistician*. 2003; **57**(3): 171–182.

[Publisher Full Text](#)

Ioannidis JP: **Why most published research findings are false**. *PLoS Med*. 2005; **2**(8): e124.

[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)

Johnson VE: **Revised standards for statistical evidence**. *Proc Natl Acad Sci U S A*. 2013; **110**(48): 19313–19317.

[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)

Killeen PR: **An alternative to null-hypothesis significance tests**. *Psychol Sci*. 2005; **16**(5): 345–353.

[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)

Kruschke JK: **Bayesian Assessment of Null Values Via Parameter Estimation and Model Comparison**. *Perspect Psychol Sci*. 2011; **6**(3): 299–312.

[PubMed Abstract](#) | [Publisher Full Text](#)

Krzywinski M, Altman N: **Points of significance: Significance, P values and t-tests**. *Nat Methods*. 2013; **10**(11): 1041–1042.

[PubMed Abstract](#) | [Publisher Full Text](#)

Lakens D: **Calculating and reporting effect sizes to facilitate cumulative science: a practical primer for t-tests and ANOVAs**. *Front Psychol*. 2013; **4**: 863.

[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)

Lakens D, Evers ER: **Sailing From the Seas of Chaos Into the Corridor of Stability: Practical Recommendations to Increase the Informational Value of Studies**. *Perspect Psychol Sci*. 2014; **9**(3): 278–292.

[PubMed Abstract](#) | [Publisher Full Text](#)

Lindley D: **The philosophy of statistics**. *Journal of the Royal Statistical Society*. 2000; **49**(3): 293–337.

[Publisher Full Text](#)

Miller J: **What is the probability of replicating a statistically significant effect?**

*Psychon Bull Rev*. 2009; **16**(4): 617–640.

[PubMed Abstract](#) | [Publisher Full Text](#)

Morey RD, Rouder JN: **Bayes factor approaches for testing interval null hypotheses**. *Psychol Methods*. 2011; **16**(4): 406–419.

[PubMed Abstract](#) | [Publisher Full Text](#)

Neyman J, Pearson ES: **On the Use and Interpretation of Certain Test Criteria for Purposes of Statistical Inference: Part I**. *Biometrika*. 1928; **20A**(1/2): 175–240.

[Publisher Full Text](#)

Neyman J, Pearson ES: **On the problem of the most efficient tests of statistical hypotheses**. *Philos Trans R Soc Lond Ser A*. 1933; **231**(694–706): 289–337.

[Publisher Full Text](#)

Nickerson RS: **Null hypothesis significance testing: a review of an old and continuing controversy**. *Psychol Methods*. 2000; **5**(2): 241–301.

[PubMed Abstract](#) | [Publisher Full Text](#)

Nuzzo R: **Scientific method: statistical errors**. *Nature*. 2014; **506**(7487): 150–152.

[PubMed Abstract](#) | [Publisher Full Text](#)

Open Science Collaboration. **PSYCHOLOGY. Estimating the reproducibility of psychological science**. *Science*. 2015; **349**(6251): aac4716.

[PubMed Abstract](#) | [Publisher Full Text](#)

Rosenthal R: **Cumulating psychology: an appreciation of Donald T. Campbell**. *Psychol Sci*. 1991; **2**(4): 213–221.

[Publisher Full Text](#)

Savalei V, Dunn E: **Is the call to abandon p-values the red herring of the replicability crisis?** *Front Psychol*. 2015; **6**: 245.

[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)

Tan SH, Tan SB: **The Correct Interpretation of Confidence Intervals**. *Proceedings of Singapore Healthcare*. 2010; **19**(3): 276–278.

[Publisher Full Text](#)

Turkheimer FE, Aston JA, Cunningham VJ: **On the logic of hypothesis testing in functional imaging**. *Eur J Nucl Med Mol Imaging*. 2004; **31**(5): 725–732.

[PubMed Abstract](#) | [Publisher Full Text](#)

van Assen MA, van Aert RC, Nuijten MB, *et al.*: **Why Publishing Everything Is More Effective than Selective Publishing of Statistically Significant Results**. *PLoS One*. 2014; **9**(1): e84896.

[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)

Walker E, Nowacki AS: **Understanding equivalence and noninferiority testing**. *J Gen Intern Med*. 2011; **26**(2): 192–196.

[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)

Wasserstein RL, Lazar NA: **The ASA's Statement on p-Values: Context, Process, and Purpose**. *The American Statistician*. 2016; **70**(2): 129–133.

[Publisher Full Text](#)

Wilcox R: **Introduction to Robust Estimation and Hypothesis Testing**. Edition 3, Academic Press, Elsevier: Oxford, UK, ISBN: 978-0-12-386983-8. 2012.

[Reference Source](#)

# Open Peer Review

Current Referee Status:



Version 2

Referee Report 28 September 2016

doi:10.5256/f1000research.9903.r16257



**Stephen J. Senn**

Luxembourg Institute of Health, Strassen, L-1445, Luxembourg

On the whole I think that this article is reasonable, my main reservation being that I have my doubts on whether the literature needs yet another tutorial on this subject.

A further reservation I have is that the author, following others, stresses what in my mind is a relatively unimportant distinction between the Fisherian and Neyman-Pearson (NP) approaches. The distinction stressed by many is that the NP approach leads to a dichotomy accept/reject based on probabilities established in advance, whereas the Fisherian approach uses tail area probabilities calculated from the observed statistic. I see this as being unimportant and not even true. Unless one considers that the person carrying out a hypothesis test (original tester) is mandated to come to a conclusion on behalf of all scientific posterity, then one must accept that any remote scientist can come to his or her conclusion depending on the personal type I error favoured. To operate the results of an NP test carried out by the original tester, the remote scientist then needs to know the p-value. The type I error rate is then compared to this to come to a personal accept or reject decision (1). In fact Lehmann (2), who was an important developer of and proponent of the NP system, describes exactly this approach as being good practice. (See *Testing Statistical Hypotheses*, 2nd edition P70). Thus using tail-area probabilities calculated from the observed statistics does not constitute an operational difference between the two systems.

A more important distinction between the Fisherian and NP systems is that the former does not use alternative hypotheses(3). Fisher's opinion was that the null hypothesis was more primitive than the test statistic but that the test statistic was more primitive than the alternative hypothesis. Thus, alternative hypotheses could not be used to justify choice of test statistic. Only experience could do that.

Further distinctions between the NP and Fisherian approach are to do with conditioning and whether a null hypothesis can ever be accepted.

I have one minor quibble about terminology. As far as I can see, the author uses the usual term 'null hypothesis' and the eccentric term 'nil hypothesis' interchangeably. It would be simpler if the latter were abandoned.

## References

1. Senn S: A comment on replication, p-values and evidence S.N. Goodman, *Statistics in Medicine* 1992;11:875-879. *Statistics in Medicine*. 2002; **21** (16): 2437-2444 [Publisher Full Text](#)
2. Lehmann E L: *Testing Statistical Hypotheses*, 2nd edition. *Chapman and Hall*. 1993.



3. Senn S: You may believe you are a Bayesian but you are probably wrong. *RMM*. 2011; **2**: 41-66

[Reference Source](#)

**I have read this submission. I believe that I have an appropriate level of expertise to confirm that it is of an acceptable scientific standard, however I have significant reservations, as outlined above.**

**Competing Interests:** No competing interests were disclosed.

---

### Version 1

Referee Report 10 November 2015

doi:[10.5256/f1000research.7499.r11036](https://doi.org/10.5256/f1000research.7499.r11036)



**Marcel ALM van Assen**

Department of Methodology and Statistics, Tilburgh University, Tilburg, Netherlands

Null hypothesis significance testing (NHST) is a difficult topic, with misunderstandings arising easily. Many texts, including basic statistics books, deal with the topic, and attempt to explain it to students and anyone else interested. I would refer to a good basic text book, for a detailed explanation of NHST, or to a specialized article when wishing an explaining the background of NHST. So, what is the added value of a new text on NHST? In any case, the added value should be described at the start of this text. Moreover, the topic is so delicate and difficult that errors, misinterpretations, and disagreements are easy. I attempted to show this by giving comments to many sentences in the text.

Abstract: “null hypothesis significance testing is the statistical method of choice in biological, biomedical and social sciences to investigate if an effect is likely”. No, NHST is the method to test the hypothesis of no effect.

Intro: “Null hypothesis significance testing (NHST) is a method of statistical inference by which an observation is tested against a hypothesis of no effect or no relationship.” What is an ‘observation’? NHST is difficult to describe in one sentence, particularly here. I would skip this sentence entirely, here.

Section on Fisher; also explain the one-tailed test.

Section on Fisher;  $p(\text{Obs}|H_0)$  does not reflect the verbal definition (the ‘or more extreme’ part).

Section on Fisher; use a reference and citation to Fisher’s interpretation of the p-value

Section on Fisher; “This was however only intended to be used as an indication that there is something in the data that deserves further investigation. The reason for this is that only  $H_0$  is tested whilst the effect under study is not itself being investigated.” First sentence, can you give a reference? Many people say a lot about Fisher’s intentions, but the good man is dead and cannot reply... Second sentence is a bit awkward, because the effect is investigated in a way, by testing the  $H_0$ .

Section on p-value; Layout and structure can be improved greatly, by first again stating what the p-value is, and then statement by statement, what it is not, using separate lines for each statement. Consider

adding that the p-value is randomly distributed under H0 (if all the assumptions of the test are met), and that under H1 the p-value is a function of population effect size and N; the larger each is, the smaller the p-value generally is.

Skip the sentence "If there is no effect, we should replicate the absence of effect with a probability equal to 1-p". Not insightful, and you did not discuss the concept 'replicate' (and do not need to).

Skip the sentence "The total probability of false positives can also be obtained by aggregating results (Ioannidis, 2005)." Not strongly related to p-values, and introduces unnecessary concepts 'false positives' (perhaps later useful) and 'aggregation'.

Consider deleting; "If there is an effect however, the probability to replicate is a function of the (unknown) population effect size with no good way to know this from a single experiment (Killeen, 2005)."

The following sentence; "Finally, a (small) p-value *is not an indication favouring a hypothesis*. A low p-value indicates a misfit of the null hypothesis to the data and cannot be taken as evidence in favour of a specific alternative hypothesis more than any other possible alternatives such as measurement error and selection bias (Gelman, 2013)." is surely not mainstream thinking about NHST; I would surely delete that sentence. In NHST, a p-value is used for testing the H0. Why did you not yet discuss significance level? Yes, before discussing what is not a p-value, I would explain NHST (i.e., what it is and how it is used).

Also the next sentence "The more (a priori) implausible the alternative hypothesis, the greater the chance that a finding is a false alarm (Krzywinski & Altman, 2013; Nuzzo, 2014)." is not fully clear to me. This is a Bayesian statement. In NHST, no likelihoods are attributed to hypotheses; the reasoning is "IF H0 is true, then...".

Last sentence: "As Nickerson (2000) puts it 'theory corroboration requires the testing of multiple predictions because the chance of getting statistically significant results for the wrong reasons in any given case is high'." What is relation of this sentence to the contents of this section, precisely?

Next section: "For instance, we can estimate that the probability of a given F value to be in the critical interval [+2 +∞] is less than 5%" This depends on the degrees of freedom.

"When there is no effect (H0 is true), the erroneous rejection of H0 is known as type I error and is equal to the p-value." Strange sentence. The Type I error is the probability of erroneously rejecting the H0 (so, when it is true). The p-value is ... well, you explained it before; it surely does not equal the Type I error.

Consider adding a figure explaining the distinction between Fisher's logic and that of Neyman and Pearson.

"When the test statistics falls outside the critical region(s)" What is outside?

"There is a profound difference between accepting the null hypothesis and simply failing to reject it (Killeen, 2005)" I agree with you, but perhaps you may add that some statisticians simply define "accept H0" as obtaining a p-value larger than the significance level. Did you already discuss the significance level, and it's mostly used values?

"To accept or reject equally the null hypothesis, Bayesian approaches (Dienes, 2014; Kruschke, 2011) or confidence intervals must be used." Is 'reject equally' appropriate English? Also using Cis, one cannot

accept the H0.

Do you start discussing alpha only in the context of CIs?

“CI also indicates the precision of the estimate of effect size, but unless using a percentile bootstrap approach, they require assumptions about distributions which can lead to serious biases in particular regarding the symmetry and width of the intervals (Wilcox, 2012).” Too difficult, using new concepts. Consider deleting.

“Assuming the CI (a)symmetry and width are correct, this gives some indication about the likelihood that a similar value can be observed in future studies, with 95% CI giving about 83% chance of replication success (Lakens & Evers, 2014).” This statement is, in general, completely false. It very much depends on the sample sizes of both studies. If the replication study has a much, much, much larger N, then the probability that the original CI will contain the effect size of the replication approaches  $(1-\alpha)*100\%$ . If the original study has a much, much, much larger N, then the probability that the original CI will contain the effect size of the replication study approaches 0%.

“Finally, contrary to p-values, CI can be used to accept H0. Typically, if a CI includes 0, we cannot reject H0. If a critical null region is specified rather than a single point estimate, for instance  $[-2, +2]$  and the CI is included within the critical null region, then H0 can be accepted. Importantly, the critical region must be specified a priori and cannot be determined from the data themselves.” No. H0 cannot be accepted with CIs.

“The (posterior) probability of an effect can however not be obtained using a frequentist framework.” Frequentist framework? You did not discuss that, yet.

“X% of times the CI obtained will contain the same parameter value”. The same? True, you mean?

“e.g. X% of the times the CI contains the same mean” I do not understand; which mean?

“The alpha value has the same interpretation as when using H0, i.e. we accept that 1-alpha CI are wrong in alpha percent of the times. “ What do you mean, CI are wrong? Consider rephrasing.

“To make a statement about the probability of a parameter of interest, likelihood intervals (maximum likelihood) and credibility intervals (Bayes) are better suited.” ML gives the likelihood of the data given the parameter, not the other way around.

“Many of the disagreements are not on the method itself but on its use.” Bayesians may disagree.

“If the goal is to establish the likelihood of an effect and/or establish a pattern of order, because both requires ruling out equivalence, then NHST is a good tool (Frick, 1996)” NHST does not provide evidence on the likelihood of an effect.

“If the goal is to establish some quantitative values, then NHST is not the method of choice.” P-values are also quantitative... this is not a precise sentence. And NHST may be used in combination with effect size estimation (this is even recommended by, e.g., the American Psychological Association (APA)).

“Because results are conditioned on H0, NHST cannot be used to establish beliefs.” It can reinforce some beliefs, e.g., if H0 or any other hypothesis, is true.

“To estimate the probability of a hypothesis, a Bayesian analysis is a better alternative.” It is the only alternative?

“Note however that even when a specific quantitative prediction from a hypothesis is shown to be true (typically testing H1 using Bayes), it does not prove the hypothesis itself, it only adds to its plausibility.” How can we *show* something is true?

I do not agree on the contents of the last section on ‘minimal reporting’. I prefer ‘optimal reporting’ instead, i.e., the reporting the information that is essential to the interpretation of the result, to any ready, which may have other goals than the writer of the article. This reporting includes, for sure, an estimate of effect size, and preferably a confidence interval, which is in line with recommendations of the APA.

**I have read this submission. I believe that I have an appropriate level of expertise to state that I do not consider it to be of an acceptable scientific standard, for reasons outlined above.**

**Competing Interests:** No competing interests were disclosed.

Author Response 06 Jul 2016

**Cyril Pernet**, The University of Edinburgh, UK

- *Null hypothesis significance testing (NHST) is a difficult topic, with misunderstandings arising easily. Many texts, including basic statistics books, deal with the topic, and attempt to explain it to students and anyone else interested. I would refer to a good basic text book, for a detailed explanation of NHST, or to a specialized article when wishing an explaining the background of NHST. So, what is the added value of a new text on NHST? In any case, the added value should be described at the start of this text. Moreover, the topic is so delicate and difficult that errors, misinterpretations, and disagreements are easy. I attempted to show this by giving comments to many sentences in the text.*

The idea of this short review was to point to common interpretation errors (stressing again and again that we are under H0) being in using p-values or CI, and also proposing reporting practices to avoid bias. This is now stated at the end of abstract.

Regarding text books, it is clear that many fail to clearly distinguish Fisher/Pearson/NHST, see Glinet et al (2012) J. Exp Education 71, 83-92. If you have 1 or 2 in mind that you know to be good, I'm happy to include them.

- *Abstract: “null hypothesis significance testing is the statistical method of choice in biological, biomedical and social sciences to investigate if an effect is likely”. No, NHST is the method to test the hypothesis of no effect.*

I agree – yet people use it to investigate (not test) if an effect is likely. The issue here is wording. What about adding this distinction at the end of the sentence?: ‘null hypothesis significance testing is the statistical method of choice in biological, biomedical and social sciences used to investigate if an effect is likely, even though it actually tests for the hypothesis of no effect’.

- *Intro: “Null hypothesis significance testing (NHST) is a method of statistical inference by which an observation is tested against a hypothesis of no effect or no relationship.” What is an ‘observation’? NHST is difficult to describe in one sentence, particularly here. I would skip this sentence entirely, here.*

I think a definition is needed, as it offers a starting point. What about the following: ‘NHST is a method of statistical inference by which an experimental factor is tested against a hypothesis of no effect or no relationship based on a given observation’

- *Section on Fisher; also explain the one-tailed test.  
Section on Fisher;  $p(\text{Obs} \geq t|H_0)$  does not reflect the verbal definition (the ‘or more extreme’ part).  
Section on Fisher; use a reference and citation to Fisher’s interpretation of the p-value  
Section on Fisher; “This was however only intended to be used as an indication that there is something in the data that deserves further investigation. The reason for this is that only  $H_0$  is tested whilst the effect under study is not itself being investigated.” First sentence, can you give a reference? Many people say a lot about Fisher’s intentions, but the good man is dead and cannot reply... Second sentence is a bit awkward, because the effect is investigated in a way, by testing the  $H_0$ .*

The section on Fisher has been modified (more or less) as suggested: (1) avoiding talking about one or two tailed tests (2) updating for  $p(\text{Obs} \geq t|H_0)$  and (3) referring to Fisher more explicitly (ie pages from articles and book) ; I cannot tell his intentions but these quotes leave little space to alternative interpretations.

- *Section on p-value; Layout and structure can be improved greatly, by first again stating what the p-value is, and then statement by statement, what it is not, using separate lines for each statement. Consider adding that the p-value is randomly distributed under  $H_0$  (if all the assumptions of the test are met), and that under  $H_1$  the p-value is a function of population effect size and  $N$ ; the larger each is, the smaller the p-value generally is.*

Done

- *Skip the sentence “If there is no effect, we should replicate the absence of effect with a probability equal to  $1-p$ ”. Not insightful, and you did not discuss the concept ‘replicate’ (and do not need to).*

Done

- *Skip the sentence “The total probability of false positives can also be obtained by aggregating results (Ioannidis, 2005).” Not strongly related to p-values, and introduces unnecessary concepts ‘false positives’ (perhaps later useful) and ‘aggregation’.*

Done

- *Consider deleting; “If there is an effect however, the probability to replicate is a function of the (unknown) population effect size with no good way to know this from a single experiment (Killeen, 2005).”*

Done

- *The following sentence; “Finally, a (small) p-value is not an indication favouring a hypothesis. A low p-value indicates a misfit of the null hypothesis to the data and cannot be taken as evidence in favour of a specific alternative hypothesis more than any other possible alternatives such as measurement error and selection bias (Gelman, 2013).” is surely not mainstream thinking about NHST; I would surely delete that sentence. In NHST, a p-value is used for testing the H0. Why did you not yet discuss significance level? Yes, before discussing what is not a p-value, I would explain NHST (i.e., what it is and how it is used).*

*Also the next sentence “The more (a priori) implausible the alternative hypothesis, the greater the chance that a finding is a false alarm (Krzywinski & Altman, 2013; Nuzzo, 2014).” is not fully clear to me. This is a Bayesian statement. In NHST, no likelihoods are attributed to hypotheses; the reasoning is “IF H0 is true, then...”.*

The reasoning here is as you state yourself, part 1: ‘a p-value is used for testing the H0; and part 2: ‘no likelihoods are attributed to hypotheses’ it follows we cannot favour a hypothesis. It might seem contentious but this is the case that all we can do is to reject the null – how could we favour a specific alternative hypothesis from there? This is explored further down the manuscript (and I now point to that) – note that we do not need to be Bayesian to favour a specific H1, all I’m saying is this cannot be attained with a p-value.

- *Last sentence: “As Nickerson (2000) puts it ‘theory corroboration requires the testing of multiple predictions because the chance of getting statistically significant results for the wrong reasons in any given case is high’.” What is relation of this sentence to the contents of this section, precisely?*

The point was to emphasise that a p value is not there to tell us a given H1 is true and can only be achieved through multiple predictions and experiments. I deleted it for clarity.

- *Next section: “For instance, we can estimate that the probability of a given F value to be in the critical interval  $[+2, +\infty]$  is less than 5%” This depends on the degrees of freedom.*

This sentence has been removed

- *“When there is no effect (H0 is true), the erroneous rejection of H0 is known as type I error and is equal to the p-value.” Strange sentence. The Type I error is the probability of erroneously rejecting the H0 (so, when it is true). The p-value is ... well, you explained it before; it surely does not equal the Type I error.*

Indeed, you are right and I have modified the text accordingly. When there is no effect (H0 is true), the erroneous rejection of H0 is known as type 1 error. Importantly, the type 1 error rate, or alpha value is determined a priori. It is a common mistake but the level of significance (for a given sample) is not the same as the frequency of acceptance alpha found on repeated sampling (Fisher, 1955).

- *Consider adding a figure explaining the distinction between Fisher’s logic and that of Neyman and Pearson.*

A figure is now presented – with levels of acceptance, critical region, level of significance and p-value.

- “When the test statistics falls outside the critical region(s)” What is outside?

*“There is a profound difference between accepting the null hypothesis and simply failing to reject it (Killeen, 2005)” I agree with you, but perhaps you may add that some statisticians simply define “accept  $H_0$ ” as obtaining a p-value larger than the significance level. Did you already discuss the significance level, and it’s mostly used values?*

*“To accept or reject equally the null hypothesis, Bayesian approaches (Dienes, 2014; Kruschke, 2011) or confidence intervals must be used.” Is ‘reject equally’ appropriate English? Also using Cis, one cannot accept the  $H_0$ .*

I should have clarified further here – as I was having in mind tests of equivalence. To clarify, I simply states now: ‘To accept the null hypothesis, tests of equivalence or Bayesian approaches must be used.’

- Do you start discussing alpha only in the context of Cis?

It is now presented in the paragraph before.

- “CI also indicates the precision of the estimate of effect size, but unless using a percentile bootstrap approach, they require assumptions about distributions which can lead to serious biases in particular regarding the symmetry and width of the intervals (Wilcox, 2012).” Too difficult, using new concepts. Consider deleting.

Done

- “Assuming the CI (a)symmetry and width are correct, this gives some indication about the likelihood that a similar value can be observed in future studies, with 95% CI giving about 83% chance of replication success (Lakens & Evers, 2014).” This statement is, in general, completely false. It very much depends on the sample sizes of both studies. If the replication study has a much, much, much larger N, then the probability that the original CI will contain the effect size of the replication approaches  $(1-\alpha)*100\%$ . If the original study has a much, much, much larger N, then the probability that the original Ci will contain the effect size of the replication study approaches 0%.

Yes, you are right, I completely overlooked this problem. The corrected sentence (with more accurate ref) is now “Assuming the CI (a)symmetry and width are correct, this gives some indication about the likelihood that a similar value can be observed in future studies. For future studies of the same sample size, 95% CI giving about 83% chance of replication success (Cumming and Mallardet, 2006). If sample sizes differ between studies, CI do not however warranty any a priori coverage”.

- “Finally, contrary to p-values, CI can be used to accept  $H_0$ . Typically, if a CI includes 0, we cannot reject  $H_0$ . If a critical null region is specified rather than a single point estimate, for instance  $[-2 +2]$  and the CI is included within the critical null region, then  $H_0$  can be

*accepted. Importantly, the critical region must be specified a priori and cannot be determined from the data themselves.” No.  $H_0$  cannot be accepted with Cis.*

Again, I had in mind equivalence testing, but in both cases you are right we can only reject and I therefore removed that sentence.

- *“The (posterior) probability of an effect can however not be obtained using a frequentist framework.” Frequentist framework? You did not discuss that, yet.*

Removed

- *“X% of times the CI obtained will contain the same parameter value”. The same? True, you mean?*

*“e.g. X% of the times the CI contains the same mean” I do not understand; which mean?*

*“The alpha value has the same interpretation as when using  $H_0$ , i.e. we accept that 1-alpha CI are wrong in alpha percent of the times. “ What do you mean, CI are wrong? Consider rephrasing.*

*“To make a statement about the probability of a parameter of interest, likelihood intervals (maximum likelihood) and credibility intervals (Bayes) are better suited.” ML gives the likelihood of the data given the parameter, not the other way around.*

corrected

- *“Many of the disagreements are not on the method itself but on its use.” Bayesians may disagree.*

removed

- *“If the goal is to establish the likelihood of an effect and/or establish a pattern of order, because both requires ruling out equivalence, then NHST is a good tool (Frick, 1996)” NHST does not provide evidence on the likelihood of an effect.*

*“If the goal is to establish some quantitative values, then NHST is not the method of choice.” P-values are also quantitative... this is not a precise sentence. And NHST may be used in combination with effect size estimation (this is even recommended by, e.g., the American Psychological Association (APA)).*

Yes, p-values must be interpreted in context with effect size, but this is not what people do. The point here is to be pragmatic, does and don't. The sentence was changed.

- *“Because results are conditioned on  $H_0$ , NHST cannot be used to establish beliefs.” It can reinforce some beliefs, e.g., if  $H_0$  or any other hypothesis, is true.*

*“To estimate the probability of a hypothesis, a Bayesian analysis is a better alternative.” It is the only alternative?*



Not for testing, but for probability, I am not aware of anything else.

- *“Note however that even when a specific quantitative prediction from a hypothesis is shown to be true (typically testing  $H_1$  using Bayes), it does not prove the hypothesis itself, it only adds to its plausibility.” How can we show something is true?*

Cumulative evidence is, in my opinion, the only way to show it. Even in hard science like physics multiple experiments. In the recent CERN study on finding Higgs bosons, 2 different and complementary experiments ran in parallel – and the cumulative evidence was taken as a proof of the true existence of Higgs bosons.

**Competing Interests:** No competing interests were disclosed.

Referee Report 30 October 2015

doi:10.5256/f1000research.7499.r10159



### Daniel Lakens

School of Innovation Sciences, Eindhoven University of Technology, Eindhoven, Netherlands

I appreciate the author's attempt to write a short tutorial on NHST. Many people don't know how to use it, so attempts to educate people are always worthwhile. However, I don't think the current article reaches its aim. For one, I think it might be practically impossible to explain a lot in such an ultra short paper - every section would require more than 2 pages to explain, and there are many sections. Furthermore, there are some excellent overviews, which, although more extensive, are also much clearer (e.g., [Nickerson, 2000](#)). Finally, I found many statements to be unclear, and perhaps even incorrect (noted below). Because there is nothing worse than creating more confusion on such a topic, I have extremely high standards before I think such a short primer should be indexed. I note some examples of unclear or incorrect statements below. I'm sorry I can't make a more positive recommendation.

“investigate if an effect is likely” – ambiguous statement. I think you mean, whether the observed DATA is probable, assuming there is no effect?

The Fisher (1959) reference is not correct – Fischer developed his method much earlier.

“This p-value thus reflects the conditional probability of achieving the observed outcome or larger,  $p(\text{Obs}|H_0)$ ” – please add 'assuming the null-hypothesis is true'.

“ $p(\text{Obs}|H_0)$ ” – explain this notation for novices.

“Following Fisher, the smaller the p-value, the greater the likelihood that the null hypothesis is false.” This is wrong, and any statement about this needs to be much more precise. I would suggest direct quotes.

“there is something in the data that deserves further investigation” –unclear sentence.

“The reason for this” – unclear what ‘this’ refers to.

“not the probability of the null hypothesis of being true,  $p(H_0)$ ” – second of can be removed?

“Any interpretation of the p-value in relation to the effect under study (strength, reliability, probability) is indeed

wrong, since the p-value is conditioned on  $H_0$ ” - incorrect. A big problem is that it depends on the sample size, and that the probability of a theory depends on the prior.

“If there is no effect, we should replicate the absence of effect with a probability equal to  $1-p$ .” I don’t understand this, but I think it is incorrect.

“The total probability of false positives can also be obtained by aggregating results (Ioannidis, 2005).” Unclear, and probably incorrect.

“By failing to reject, we simply continue to assume that  $H_0$  is true, which implies that one cannot, from a nonsignificant result, argue against a theory” – according to which theory? From a NP perspective, you can ACT as if the theory is false.

“(Lakens & Evers, 2014)” – we are not the original source, which should be cited instead.

“Typically, if a CI includes 0, we cannot reject  $H_0$ .” - when would this not be the case? This assumes a CI of  $1-\alpha$ .

“If a critical null region is specified rather than a single point estimate, for instance  $[-2 +2]$  and the CI is included within the critical null region, then  $H_0$  can be accepted.” – you mean practically, or formally? I’m pretty sure only the former.

The section on ‘The (correct) use of NHST’ seems to conclude only Bayesian statistics should be used. I don’t really agree.

“we can always argue that effect size, power, etc. must be reported.” – which power? Post-hoc power? Surely not? Other types are unknown. So what do you mean?

The recommendation on what to report remains vague, and it is unclear why what should be reported.

**I have read this submission. I believe that I have an appropriate level of expertise to state that I do not consider it to be of an acceptable scientific standard, for reasons outlined above.**

**Competing Interests:** No competing interests were disclosed.

Author Response 06 Jul 2016

**Cyril Pernet**, The University of Edinburgh, UK

- “investigate if an effect is likely” – ambiguous statement. I think you mean, whether the observed DATA is probable, assuming there is no effect?

This sentence was changed, following as well the other reviewer, to ‘null hypothesis significance testing is the statistical method of choice in biological, biomedical and social

sciences to investigate if an effect is likely, even though it actually tests whether the observed data are probable, assuming there is no effect'

- *The Fisher (1959) reference is not correct – Fischer developed his method much earlier.*

Changed, refers to Fisher 1925

- *“This p-value thus reflects the conditional probability of achieving the observed outcome or larger,  $p(\text{Obs}|\text{H}_0)$ ” – please add 'assuming the null-hypothesis is true'. “ $p(\text{Obs}|\text{H}_0)$ ” – explain this notation for novices.*

I changed a little the sentence structure, which should make explicit that this is the condition probability.

- *“Following Fisher, the smaller the p-value, the greater the likelihood that the null hypothesis is false.” This is wrong, and any statement about this needs to be much more precise. I would suggest direct quotes.*

This sentence has been removed

- *“there is something in the data that deserves further investigation” –unclear sentence. “The reason for this” – unclear what ‘this’ refers to.*

This has been changed to ‘[...] to decide whether the evidence is worth additional investigation and/or replication (Fisher, 1971 p13)’

- *“not the probability of the null hypothesis of being true,  $p(\text{H}_0)$ ” – second of can be removed?*

my mistake – the sentence structure is now ‘not the probability of the null hypothesis  $p(\text{H}_0)$ , of being true,’ ; hope this makes more sense (and this way refers back to  $p(\text{Obs}>t|\text{H}_0)$ )

- *“Any interpretation of the p-value in relation to the effect under study (strength, reliability, probability) is indeed wrong, since the p-value is conditioned on  $\text{H}_0$ ” - incorrect. A big problem is that it depends on the sample size, and that the probability of a theory depends on the prior.*

Fair enough – my point was to stress the fact that p value and effect size or  $\text{H}_1$  have very little in common, but yes that the part in common has to do with sample size. I left the conditioning on  $\text{H}_0$  but also point out the dependency on sample size.

- *“If there is no effect, we should replicate the absence of effect with a probability equal to  $1-p$ .” I don’t understand this, but I think it is incorrect.*

Removed

- *“The total probability of false positives can also be obtained by aggregating results (Ioannidis, 2005).” Unclear, and probably incorrect.*

Removed

- *“By failing to reject, we simply continue to assume that  $H_0$  is true, which implies that one cannot, from a nonsignificant result, argue against a theory” – according to which theory? From a NP perspective, you can ACT as if the theory is false.*

The whole paragraph was changed to reflect a more philosophical take on scientific induction/reasoning. I hope this is clearer.

- *“(Lakens & Evers, 2014)” – we are not the original source, which should be cited instead.*

done

- *“Typically, if a CI includes 0, we cannot reject  $H_0$ .” - when would this not be the case? This assumes a CI of 1-alpha. “If a critical null region is specified rather than a single point estimate, for instance [-2 +2] and the CI is included within the critical null region, then  $H_0$  can be accepted.” – you mean practically, or formally? I’m pretty sure only the former.*

Changed to refer to equivalence testing

- *The section on ‘The (correct) use of NHST’ seems to conclude only Bayesian statistics should be used. I don’t really agree.*

I rewrote this, as to show frequentist analysis can be used - I’m trying to sell Bayes more than any other approach.

- *“we can always argue that effect size, power, etc. must be reported.” – which power? Post-hoc power? Surely not? Other types are unknown. So what do you mean? The recommendation on what to report remains vague, and it is unclear why what should be reported.*

I’m arguing we should report it all, that’s why there is no exhausting list – I can if needed.

**Competing Interests:** No competing interests were disclosed.