

(Excerpts from) 'P-Value Thresholds: Forfeit at Your Peril' (free access)

by [Mayo](#)



A key recognition among those who write on the statistical crisis in science is that the pressure to publish attention-getting articles can incentivize researchers to produce eye-catching but inadequately scrutinized claims. We may see much the same sensationalism in broadcasting metastatistical research, especially if it takes the form of scapegoating or banning statistical significance. A lot of excitement was generated recently when Ron Wasserstein, Executive Director of the American Statistical Association (ASA), and co-editors A. Schirm and N. Lazar, updated the [2016 ASA Statement on P-Values and Statistical Significance \(ASA I\)](#). In their 2019 interpretation, ASA I “stopped just short of recommending that declarations of ‘statistical significance’ be abandoned,” and in [their new statement \(ASA II\)](#) announced: “We take that step here...‘statistically significant’ –don’t say it and don’t use it”. To herald the ASA II, and the special issue “Moving to a world beyond ‘ $p < 0.05$ ’”, the journal *Nature* requisitioned a commentary from Amrhein, Greenland and McShane “[Retire Statistical Significance](#)” ([AGM](#)). With over 800 signatories, the commentary received the imposing title “Scientists rise up against significance tests”!

[Tom Hardwicke and John Ioannidis](#) surveyed those signatories and give a report on the respondents ([Hardwicke and Ioannidis 2019](#)). I was invited to write an editorial on any aspect of the episode (“[P-value thresholds: Forfeit at your peril](#)”)–the opening of which is above. [Hardwicke and Ioannidis](#), a preprint of my editorial, and an editorial by Andrew [Gelman](#) are currently “free access” in the European Journal of

Clinical Investigation. I guess that means these version are currently freely accessible.

My article continues:

Note: By "[ASA II](#)" I allude only to the authors' general recommendations, not their summaries of the 43 papers in the issue.)

[Hardwicke and Ioannidis](#) (2019) worry that recruiting signatories on such a paper politicizes the process of evaluating a stance on scientific method, and fallaciously appeals to popularity (*argumentum ad populum*) "because it conflates justification of a belief with the acceptance of a belief by a given group of people". Opposing viewpoints are not given a similar forum. Fortunately, John Ioannidis (2019) can come out with [a note in JAMA](#) challenging ASA II and AGM, but the vast majority of stakeholders in the debate go unheard. Appealing to popularity gives a *prudential* reason to go along, it is risky to stand in opposition to the hundreds who signed, not to mention, the thought leaders at the ASA. There is also an appeal to fear, with the result that many will fear using statistical significance tests altogether. Why risk using a method that is persecuted with such zeal and fanfare?

Ioannidis (2019) points out what may not be obvious at first: it is not just a word ban but a gatekeeper ban:

Many fields of investigation ... have major gaps in the ways they conduct, analyze, and report studies and lack protection from bias. Instead of trying to fix what is lacking and set better and clearer rules, one reaction is to overturn the tables and abolish any gatekeeping rules (such as removing the term statistical significance). However, potential for falsification is a prerequisite for science. Fields that obstinately resist refutation can hide behind the abolition of statistical significance but risk becoming self-ostracized from the remit of science.

Among the top-cited signatories who respond to their questionnaire, Hardwicke and Ioannidis find a heavy representation of fields with

prevalent concerns about low reproducibility. Yet “abandoning the concept of statistical significance would make claims of ‘irreproducibility’ difficult if not impossible to make. In our opinion this approach may give bias a free pass”.

I agree, and will show why.

I continue with (excerpts of a preprint of) my article; references are formatted in the usual way. You can read the “free access” version [here](#).

....

It might be assumed I would agree to “retire significance” since I often claim “the crude dichotomy of ‘pass/fail’ or ‘significant or not’ will scarcely do” and because I reformulate tests so as to “determine the magnitudes (and directions) of any statistical discrepancies warranted, and the limits to any substantive claims you may be entitled to infer from the statistical ones.” (Mayo 2018) [Genuine effects, as Fisher insisted, require not isolated small P-values, but a reliable method to successfully generate them.] We should not confuse prespecifying minimal thresholds in each test, which I would uphold, with fixing a value to habitually use (which I would not). N-P tests called for the practitioner to balance error probabilities according to context, not rigidly fix a value like .05. Nor does having a minimal P-value threshold mean we do not report the attained P-value: we should, and N-P agreed!

The “no threshold” view is not merely to never use the S word and report continuous P-values

These two rules alone would not lead Hardwicke and Ioannidis to charge, correctly, in my judgment, that “this approach may give bias a free pass”. ASA II and AGM decry using any prespecified P-value threshold as the basis for categorizing data in some way, such as inferring that results are, or are not, evidence of a genuine effect.

- “Decisions to interpret or to publish results will not be based on statistical thresholds” (AGM).

- “Whether a p-value passes any arbitrary threshold should not be considered at all” in interpreting data ([ASA II](#)).

Consider how far reaching the “no threshold” view is for interpreting data. For example, according to ASA II, in order for the U.S. Food and Drug Administration (FDA) to comply with its “no threshold” position, it does not suffice that they report continuous P-values and confidence intervals. The FDA would have to end its “long established drug review procedures that involve comparing p -values to significance thresholds for Phase III drug trials”.

The [New England Journal of Medicine \(NEJM\)](#) responds (2019) to the ASA call to revise their guidelines, but insists that a central premise on which their revisions are based is “the use of statistical thresholds for claiming an effect or association should be limited to analyses for which the analysis plan outlined a method for controlling type I error”. In the [article accompanying](#) the revised guidelines:

“A well-designed randomized or observational study will have a primary hypothesis and a prespecified method of analysis, and the significance level from that analysis is a reliable indicator of the extent to which the observed data contradict a null hypothesis of no association between an intervention or an exposure and a response. Clinicians and regulatory agencies must make decisions about which treatment to use or to allow to be marketed, and P values interpreted by reliably calculated thresholds subjected to appropriate adjustments [for multiple trials] have a role in those decisions”.

Specifying “thresholds that have a strong theoretical and empirical justification” escapes the ASA II ruling: “Don’t conclude anything about scientific ...importance based on statistical significance”.

Although less well advertised, the “no thresholds” view also torpedoes common uses of confidence intervals and Bayes Factor standards.

[T]he problem is not that of having only two labels. Results should not

be trichotomized, or indeed categorized into any number of groups. Similarly, we need to stop using confidence intervals [CIs] as another means of dichotomizing. ([ASA II](#))

AGM's "compatibility intervals" are redolent of the consonance intervals of Kempthorne and Folks(1971) , except that the latter use many thresholds, one for each of several consonance levels. Even these would seem to violate the rule that results should not be "categorized into any number of groups".

...Nor could Bayes factor thresholds be used, as they often are, to test a null against an alternative. It is not clear how any statistical tests survive. A claim has not passed a genuine test, if none of the results are allowed to count against it. We are not told what happens to the use of significance tests to check if statistical model assumptions hold approximately, or not—essential across methodologies. As George Box, a Bayesian, remarks, "diagnostic checks and tests of fit ... require frequentist theory significance tests for their formal justification" (1983, p. 57).

What arguments are given to accept the no threshold view?

Getting past the appeals to popularity and fear, the reasons ASA II and AGM give are that thresholds can lead to well-known fallacies, and even to some howlers more extreme than those long lampooned. Of course it's true:

a statistically non-significant result does not 'prove' the null hypothesis (the hypothesis that there is no difference between groups or no effect of a treatment ...). Nor do statistically significant results 'prove' some other hypothesis. (AGM)

It is easy to be swept up in their outrage, but the argument: "significance thresholds can be used very badly, therefore remove significance thresholds" is a very bad argument. Moreover, it would remove the very standards we need to call out the fallacies. A rule that went from any

non-significant result to inferring no effect was proved, or to take something less extreme, to inferring it is well warranted or the like, would have extremely high Type II error probabilities. They deal with a point null hypothesis, which makes it even worse.

...The "free access" version is [here](#).

Giving Data Dredgers a Free Pass

The danger of removing thresholds on grounds they could be badly used is that they are not there when you need them. Ioannidis zeroes in on the problem:

The proposal to entirely remove the barrier does not mean that scientists will not often still wish to interpret their results as showing important signals and fit preconceived notions and biases. With the gatekeeper of statistical significance, eager investigators whose analyses yield, for example, $P = .09$ have to either manipulate their statistics to get to $P < .05$ or add spin to their interpretation to suggest that results point to an important signal through an observed "trend." When that gate keeper is removed, any result may be directly claimed to reflect an important signal or fit to a preexisting narrative.

As against Ioannidis' anything goes charge, it might be said that even in a world without thresholds a largish P-value could not be taken as evidence of a genuine effect. For to do so would be to say something nonsensical. It would be to say: Even though larger differences would frequently be expected by chance variability alone (i.e., even though the P-value is largish), I maintain the data provide evidence they are not due to chance variability.

But such a response turns on appealing to a threshold to block it, minimally requiring the P-value be rather small e.g., $< .1$? (It also shows why P-values are apt measures for the job of distinguishing random error.) Thus, our eager investigators, facing a non-small P-value, are still incentivized to manipulate their statistics. Say they ransack the data until

finding a non-prespecified subgroup that provides a nominally small enough P-value. In a world without thresholds, we would be hamstrung from highlighting, critically, P-values that breach (as opposed to uphold) preset thresholds.

“Whether a p-value passes any arbitrary threshold should not be considered *at all* when deciding which results to present or highlight” (my emphasis, [ASA II](#)).

More important than keeping a specific word is keeping a filter for error control. The 2016 ASA I warned in Principle 4: “Valid scientific conclusions based on p-values and related statistics cannot be drawn without at least knowing how many and which analyses were conducted, and how those analyses (including p-values) were selected for reporting”. ...An unanswered question is how Principle 4 is to operate in a world with ASA II.

The NEJM’s revised guidelines, far from agreeing to use P-values without error probability thresholds, will now be stricter in their use. When no method to adjust for multiplicity of inferences or controlling the Type I error probability is prespecified, the report of secondary endpoints

should be limited to point estimates of treatment effects with 95% confidence intervals. In such cases, the Methods section should note that the widths of the intervals have not been adjusted for multiplicity and that the inferences drawn may not be reproducible. No P values should be reported for these analyses.

Confidence intervals severed from their dualities with tests, from which they were initially developed, lose their error probability guarantees.

Conclusion

The ASA P-value project is lately careering into recommendations on which there has been little balanced discussion and much disagreement. Hardwicke and Ioannidis find that more than half of the respondents deny significance should be excluded from all science, and the 43 papers in

the special issue "Moving to a world beyond 'p < 0.05'" offer a cacophony of competing reforms.

It is hard to resist the missionary zeal of masterful calls: Do you want bad science to thrive? or Do you want to ban significance? (a false dilemma). A question to raise before jumping on the bandwagon: Are they asking the most unbiased questions about the consequences of removing thresholds currently ensconced into hundreds of legal statutes and best practice manuals? This needs to be carefully considered, if the reforms intended to improve credibility of statistics are not to backfire, as they may already be doing.

ASA II is part of a large undertaking; it contains plenty of sagacious advice. Notably the M in ATOM: Modesty.

Be modest by recognizing that different readers may have very different stakes on the results of your analysis, which means you should try to take the role of a neutral judge rather than an advocate for any hypothesis.

ASA II regards its positions "open to debate". An open debate is very much needed.

Here's the [full \(uncorrected\) preprint of my editorial](#).

*Mayo (2018), Mayo and Cox (2006), Mayo and Spanos (2006).

References not linked above

Birnbaum, A. [Statistical Methods in Scientific Inference](#) (letter to the Editor), *Nature* 1970;225(5237):1033.

Box, G. An apology for ecumenism in statistics. In G. E. P. Box, T. Leonard, and D. F. J. Wu (Eds.), *Scientific inference, data analysis, and robustness*. Academic Press, 1983:51-84.

Fisher, RA. *The design of experiments*, Oliver and Boyd, 1947.

Kempthorne, O, Folks, J. Probability, statistics, data analysis. Iowa State University Press 1971.

Mayo (2019). *Statistical Inference as Severe Testing: How to Get Beyond the Statistics Wars* (2018, CUP).

Mayo, D.G. and Cox, D. R. (2006) "[Frequentist Statistics as a Theory of Inductive Inference](#)," *Optimality: The Second Erich L. Lehmann Symposium* (ed. J. Rojo), Lecture Notes-Monograph series, Institute of Mathematical Statistics (IMS), Vol. 49: 77-97.

Mayo, D. G. and Spanos, A. (2006). "[Severe Testing as a Basic Concept in a Neyman-Pearson Philosophy of Induction](#)," *British Journal of Philosophy of Science*, 57: 323-357.

Neyman, J. Tests of statistical hypotheses and their use in studies of natural phenomena. *Communications in Statistics: Theory and Methods* 1976;5(8):737-51.

NEJM Author Guidelines: Retrieved from: <https://www.nejm.org/author-center/new-manuscripts> on July 19, 2019.

Relevant (2019) posts:

[The 2019 ASA Guide to P-values and Significance: Don't say What You don't Mean \(Some Recommendations\)](#)

[The NEJM Issues New Guidelines on Statistical Reporting: Is the ASA P-Value Project Backfiring?](#)