



Revised evidence for statistical standards

In ref. 1, Johnson proposes replacing the usual P = 0.05 standard for significance with the more stringent P = 0.005. This might be good advice in practice, but we remain troubled by Johnson's logic because it seems to dodge the essential nature of any such rule, which is that it expresses a tradeoff between the risks of publishing misleading results and of important results being left unpublished. Ultimately such decisions should depend on costs, benefits, and probabilities of all outcomes.

Johnson's minimax prior is not intended to correspond to any distribution of effect sizes; rather, it represents a worst case scenario under some mathematical assumptions. Minimax and tradeoffs do not play well together (2), and it is hard for us to see how any worst case procedure can supply much guidance on how to balance between two different losses.

For example, in a genomics or shot-in-thedark drug discovery program in which thousands of possibilities are being tried in a search for a few patterns, it could make sense to have a very stringent threshold. But in settings where differences are large, such as the evaluation of teachers (3), it can be a mistake to set aside real differences just because they do not exceed a high level of statistical significance.

Johnson's evidence threshold is chosen relative to a conventional value, namely Jeffreys' target Bayes factor of 1/25 or 1/50, for which we do not see any particular justification except with reference to the tail area probability of 0.025, traditionally associated with statistical significance.

To understand the difficulty of this approach, consider the hypothetical scenario in which R. A. Fisher had chosen P = 0.005 rather than P = 0.05 as a significance threshold. In this alternative history, the discrepancy between P values and Bayes factors remains, and Johnson could have written a paper noting that the accepted 0.005 standard fails to correspond to 200-to-1 evidence against the null. Indeed, a 200:1 evidence in a minimax sense gets processed by his

fixed-point equation of $\gamma = \exp \left| z \sqrt{2 \log(\gamma)} - \right|$

 $\log(\gamma)$ at the value of $\gamma = 0.005$, into z =

 $\sqrt{-2\log(0.005)} = 3.86$, which corresponds to a (one-sided) tail probability of $\Phi(-3.86)$: ~0.0005. Moreover, the proposition approximately divides any small initial *P* level by a factor of $\sqrt{-4\pi \log(p)}$, roughly equal to 10 for the *P*s of interest. Thus, Johnson's recommended threshold of *P* = 0.005 stems from taking 1/20 as a starting point; *P* = 0.005 has no justification on its own (any more than does the *P* = 0.005 threshold derived from the alternative default standard of 1/200).

One might then ask, was Fisher foolish to settle for the P = 0.05 rule that has caused so many problems in later decades? We would argue that the appropriate significance level depends on the scenario and that what

worked well for agricultural experiments in the 1920s might not be so appropriate for many applications in modern biosciences. Thus, Johnson's recommendation to rethink significance thresholds seems like a good idea that needs to include assessments of actual costs, benefits, and probabilities, rather than being based on an abstract calculation.

Andrew Gelman^a and Christian P. Robert^{b,c,1}

^aDepartments of Statistics and Political Science, Columbia University, New York, NY 10027; ^bCentre de Recherche en Mathématiques de la Décision, Université Paris-Dauphine, 75775 Paris Cedex 16, France; and ^cDepartment of Statistics, University of Warwick, Coventry, Coventry CV4 7AL, United Kingdom

Author contributions: A.G. and C.P.R. designed research, performed research, and wrote the paper.

The authors declare no conflict of interest

¹To whom correspondence should be addressed. E-mail: xian@ ceremade.dauphine.fr.

Johnson VE (2013) Revised standards for statistical evidence. Proc Natl Acad Sci USA 110(48):19313–19317.
Berger J (1985) Statistical Decision Theory and Bayesian Analysis

⁽Springer-Verlag, New York), 2nd Ed. 3 Kane TJ (2013) Presumed averageness: The mis-application of classical hypothesis testing in education. The Brown Center

Chalkboard, Brookings Institution. Available at www.brookings.edu/ blogs/brown-center-chalkboard/posts/2013/12/04-classical-hypotesistesting-in-education-kane. Accessed December 4, 2013.