

1 Justify Your Alpha: A Response to “Redefine Statistical Significance”

2

3 Daniel Lakens^{abc1}, Federico G. Adolfini^{bc2}, Casper J. Albers^{ab3}, Farid Anvari^{d4}, Matthew A. J. Apps^{a5},
 4 Shlomo E. Argamon^{ab6}, Thom Baguley^{ab7}, Raymond B. Becker^{ac8}, Stephen D. Benning^{a9}, Daniel E.
 5 Bradford^{a10}, Erin M. Buchanan^{ab11}, Aaron R. Caldwell^{d12}, Ben van Calster^{ab13}, Rickard Carlsson^{d14},
 6 Sau-Chin Chen^{a15}, Bryan Chung^{a16}, Lincoln J. Colling^{a17}, Gary S. Collins^{b18}, Zander Crook^{ab19}, Emily S.
 7 Cross^{d20}, Sameera Daniels^{ab21}, Henrik Danielsson^{a22}, Lisa DeBruine^{a23}, Daniel J. Dunleavy^{ab24}, Brian
 8 D. Earp^{ab25}, Michele I. Feist^{bc26}, Jason D. Ferrell^{ab27}, James G. Field^{ab28}, Nicholas W. Fox^{abc29}, Amanda
 9 Friesen^{d30}, Caio Gomes^{d31}, Monica Gonzalez-Marquez^{abc32}, James A. Grange^{abc33}, Andrew P.
 10 Grieve^{d34}, Robert Guggenberger^{d35}, James Grist^{d36}, Anne-Laura van Harmelen^{ab37}, Fred
 11 Hasselman^{bc38}, Kevin D. Hochard^{d39}, Mark R. Hoffarth^{a40}, Nicholas P. Holmes^{abc41}, Michael Ingre^{ab42},
 12 Peder M. Isager^{b43}, Hanna K. Isotalus^{ab44}, Christer Johansson^{d45}, Konrad Juszczyk^{d46}, David A.
 13 Kenny^{d47}, Ahmed A. Khalil^{abc48}, Barbara Konat^{d49}, Junpeng Lao^{ab50}, Erik Gahner Larsen^{a51}, Gerine M.
 14 A. Lodder^{ab52}, Jiří Lukavský^{d53}, Christopher R. Madan^{d54}, David Manheim^{ab55}, Stephen R. Martin^{abc56},
 15 Andrea E. Martin^{ab57}, Deborah G. Mayo^{d58}, Randy J. McCarthy^{a59}, Kevin McConway^{ab60}, Colin
 16 McFarland^{b61}, Amanda Q. X. Nio^{ab62}, Gustav Nilsson^{ab63}, Cilene Lino de Oliveira^{b64}, Jean-Jacques
 17 Orban de Xivry^{ab65}, Sam Parsons^{bc66}, Gerit Pfuhl^{ab67}, Kimberly A. Quinn^{b68}, John J. Sakon^{a69}, S. Adil
 18 Saribay^{a70}, Iris K. Schneider^{ab71}, Manojkumar Selvaraju^{d72}, Zsuzsika Sjoerds^{b73}, Samuel G. Smith^{b74},
 19 Tim Smits^{a75}, Jeffrey R. Spies^{b76}, Vishnu Sreekumar^{abc77}, Crystal N. Steltenpohl^{abc78}, Neil
 20 Stenhouse^{a79}, Wojciech Świątkowski^{a80}, Miguel A. Vadillo^{a81}, Marcel A. L. M. Van Assen^{ab82}, Matt N.
 21 Williams^{ab83}, Samantha E. Williams^{d84}, Donald R. Williams^{ab85}, Tal Yarkoni^{b86}, Ignazio Ziano^{d87}, Rolf A.
 22 Zwaan^{ab88}

23

24 a) Participated in brainstorming. b) Participated in drafting the commentary. c) Conducted statistical
 25 analyses/data preparation. d) Did not participate in a, b, or c, because the points that they would have
 26 raised had already been incorporated into the commentary, or endorse a sufficiently large part of the
 27 contents as if participation had occurred. Except for the first author, authorship order is alphabetical.

28

29

30

Affiliations

31

32 ¹Human-Technology Interaction, Eindhoven University of Technology, Den Dolech, 5600MB,
 33 Eindhoven, The Netherlands

34 ²Laboratory of Experimental Psychology and Neuroscience (LPEN), Institute of Cognitive and
 35 Translational Neuroscience (INCYT), INECO Foundation, Favaloro University, Pacheco de Melo
 36 1860, Buenos Aires, Argentina

37 ²National Scientific and Technical Research Council (CONICET), Godoy Cruz 2290, Buenos Aires,
 38 Argentina

39 ³Heymans Institute for Psychological Research, University of Groningen, Grote Kruisstraat 2/1,
 40 9712TS Groningen, The Netherlands

- 1 ⁴College of Education, Psychology & Social Work, Flinders University, Adelaide, GPO Box 2100,
2 Adelaide, SA, 5001, Australia
- 3 ⁵Department of Experimental Psychology, University of Oxford, New Radcliffe House, Oxford, OX2
4 6GG, UK
- 5 ⁶Department of Computer Science, Illinois Institute of Technology, Chicago, IL, 10 W. 31st Street,
6 Chicago, IL 60645, USA
- 7 ⁷Department of Psychology, Nottingham Trent University, Nottingham, 50 Shakespeare Street,
8 Nottingham, NG1 4FQ, UK
- 9 ⁸Faculty of Linguistics and Literature, Bielefeld University, Bielefeld, Universitätsstraße 25, 33615
10 Bielefeld, Germany
- 11 ⁹Psychology, University of Nevada, Las Vegas, Las Vegas, 4505 S. Maryland Pkwy., Box 455030, Las
12 Vegas, NV 89154-5030, USA
- 13 ¹⁰Psychology, University of Wisconsin-Madison, Madison, 1202 West Johnson St. Madison WI. 53706,
14 USA
- 15 ¹¹Psychology, Missouri State University, 901 S. National Ave, Springfield, MO, 65897, USA
- 16 ¹²Health, Human Performance, and Recreation, University of Arkansas, Fayetteville, 155 Stadium
17 Drive, HPER 321, Fayetteville, AR, 72701, USA
- 18 ¹³Department of Development and Regeneration, KU Leuven, Leuven, Herestraat 49 box 805, 3000
19 Leuven, Belgium, Belgium
- 20 ¹³Department of Medical Statistics and Bioinformatics, Leiden University Medical Center, Postbus
21 9600, 2300 RC, Leiden, The Netherlands
- 22 ¹⁴Department of Psychology, Linnaeus University, Kalmar, Stagneliusgatan 14, 392 34, Kalmar,
23 Sweden
- 24 ¹⁵Department of Human Development and Psychology, Tzu-Chi University, No. 67, Jieren St., Hualien
25 City, Hualien County, 97074, Taiwan
- 26 ¹⁶Department of Surgery, University of British Columbia, Victoria, #301 - 1625 Oak Bay Ave, Victoria
27 BC Canada, V8R 1B1 , Canada
- 28 ¹⁷Department of Psychology, University of Cambridge, Cambridge CB2 3EB, UK
- 29 ¹⁸Centre for Statistics in Medicine, University of Oxford, Windmill Road, Oxford, OX3 7LD, UK
- 30 ¹⁹Department of Psychology, The University of Edinburgh, 7 George Square, Edinburgh, EH8 9JZ, UK
- 31 ²⁰School of Psychology, Bangor University, Bangor, Adeilad Brigantia, Bangor, Gwynedd, LL57 2AS,
32 UK
- 33 ²¹Ramsey Decision Theoretics, 4849 Connecticut Ave. NW #132, Washington, DC 20008, USA
- 34 ²²Department of Behavioural Sciences and Learning, Linköping University, SE-581 83, Linköping,
35 Sweden
- 36 ²³Institute of Neuroscience and Psychology, University of Glasgow, Glasgow, 58 Hillhead Street, UK
- 37 ²⁴College of Social Work, Florida State University, 296 Champions Way, University Center C,
38 Tallahassee, FL, 32304, USA
- 39 ²⁵Departments of Psychology and Philosophy, Yale University, 2 Hillhouse Ave, New Haven CT
40 06511, USA

- 1 ²⁶Department of English, University of Louisiana at Lafayette, P. O. Box 43719, Lafayette LA 70504,
2 USA
- 3 ²⁷Department of Psychology, St. Edward's University, 3001 S. Congress, Austin, TX 78704, USA
- 4 ²⁷Department of Psychology, University of Texas at Austin, 108 E. Dean Keeton Stop A8000, Austin,
5 TX 78712-1043, USA
- 6 ²⁸Department of Management, West Virginia University, 1602 University Avenue, Morgantown, WV
7 26506, USA
- 8 ²⁹Department of Psychology, Rutgers University, New Brunswick, 53 Avenue E, Piscataway NJ
9 08854, USA
- 10 ³⁰Department of Political Science, Indiana University Purdue University, Indianapolis, Indianapolis,
11 425 University Blvd CA417, Indianapolis, IN 46202, USA
- 12 ³¹Booking.com, Herengracht 597, 1017 CE Amsterdam, The Netherlands
- 13 ³²Department of English, American and Romance Studies, RWTH - Aachen University, Aachen,
14 Kármánstraße 17/19, 52062 Aachen, Germany
- 15 ³³School of Psychology, Keele University, Keele, Staffordshire, ST5 5BG, UK
- 16 ³⁴Centre of Excellence for Statistical Innovation, UCB Celltech, 208 Bath Road, Slough, Berkshire SL1
17 3WE, UK
- 18 ³⁵Translational Neurosurgery, Eberhard Karls University Tübingen, Tübingen, Germany
- 19 ³⁵University Tübingen, International Centre for Ethics in Sciences and Humanities, Germany
- 20 ³⁶Department of Radiology, University of Cambridge, Box 218, Cambridge Biomedical Campus, CB2
21 0QQ, UK
- 22 ³⁷Department of Psychiatry, University of Cambridge, Cambridge, 18b Trumpington Road, CB2 8AH,
23 UK
- 24 ³⁸Behavioural Science Institute, Radboud University Nijmegen, Montessorilaan 3, 6525 HR, Nijmegen,
25 The Netherlands
- 26 ³⁹Department of Psychology, University of Chester, Chester, Department of Psychology, University of
27 Chester, Chester, CH1 4BJ, UK
- 28 ⁴⁰Department of Psychology, New York University, 4 Washington Place, New York, NY 10003, USA
- 29 ⁴¹School of Psychology, University of Nottingham, Nottingham, University Park, NG7 2RD, UK
- 30 ⁴²None, Independent, Stockholm, Skåpvägen 5, 12245 ENSKEDE, Sweden
- 31 ⁴³Department of Clinical and Experimental Medicine, University of Linköping, 581 83 Linköping,,
32 Sweden
- 33 ⁴⁴School of Clinical Sciences, University of Bristol, Bristol, Level 2 academic offices, L&R Building,
34 Southmead Hospital, BS10 5NB, UK
- 35 ⁴⁵Occupational Orthopaedics and Research, Sahlgrenska University Hospital, 413 45 Gothenburg,
36 Sweden
- 37 ⁴⁶The Faculty of Modern Languages and Literatures, Institute of Linguistics, Psycholinguistics
38 Department, Adam Mickiewicz University, Al. Niepodległości 4, 61-874, Poznań, Poland
- 39 ⁴⁷Department of Psychological Sciences, University of Connecticut, Storrs, CT, Department of
40 Psychological Sciences, U-1020, Storrs, CT 06269-1020, USA

- 1 ⁴⁸Center for Stroke Research Berlin, Charité - Universitätsmedizin Berlin, Hindenburgdamm 30, 12200
2 Berlin, Germany
- 3 ⁴⁸Max Planck Institute for Human Cognitive and Brain Sciences, Stephanstraße 1a, 04103 Leipzig,
4 Germany
- 5 ⁴⁸Berlin School of Mind and Brain, Humboldt-Universität zu Berlin, Luisenstraße 56, 10115 Berlin,
6 Germany
- 7 ⁴⁰Social Sciences, Adam Mickiewicz University, Poznań, Szamarzewskiego 89, 60-568 Poznan,
8 Poland
- 9 ⁵⁰Department of Psychology, University of Fribourg, Faucigny 2, 1700 Fribourg, Switzerland
- 10 ⁵¹School of Politics and International Relations, University of Kent, Canterbury CT2 7NX, UK
- 11 ⁵²Department of Sociology / ICS, University of Groningen, Grote Rozenstraat 31, 9712 TG Groningen,
12 The Netherlands
- 13 ⁵³Institute of Psychology, Czech Academy of Sciences, Hybernská 8, 11000 Prague, Czech Republic
- 14 ⁵⁴School of Psychology, University of Nottingham, Nottingham, NG7 2RD, UK
- 15 ⁵⁵Pardee RAND Graduate School, RAND Corporation, 1200 S Hayes St, Arlington, VA 22202, USA
- 16 ⁵⁶Psychology and Neuroscience, Baylor University, Waco, One Bear Place 97310, Waco TX, USA
- 17 ⁵⁷Psychology of Language Department, Max Planck Institute for Psycholinguistics, Nijmegen,
18 Wundtlaan 1, 6525XD, The Netherlands
- 19 ⁵⁷Department of Psychology, School of Philosophy, Psychology, and Language Sciences, University
20 of Edinburgh, 7 George Square, EH8 9JZ Edinburgh, UK
- 21 ⁵⁸Dept of Philosophy, Major Williams Hall, Virginia Tech, Blacksburg, VA, US
- 22 ⁵⁹Center for the Study of Family Violence and Sexual Assault, Northern Illinois University, DeKalb, IL,
23 125 President's BLVD., DeKalb, IL 60115, USA
- 24 ⁶⁰School of Mathematics and Statistics, The Open University, Milton Keynes, Walton Hall, Milton
25 Keynes MK7 6AA, UK
- 26 ⁶¹Skyscanner, 15 Laurison Place, Edinburgh, EH3 9EN, UK
- 27 ⁶²School of Biomedical Engineering and Imaging Sciences, King's College London, London, UK
- 28 ⁶³Stress Research Institute, Stockholm University, Stockholm, Frescati Hagväg 16A, SE-10691
29 Stockholm, Sweden
- 30 ⁶³Department of Clinical Neuroscience, Karolinska Institutet, Nobels väg 9, SE-17177 Stockholm,
31 Sweden
- 32 ⁶³Department of Psychology, Stanford University, 450 Serra Mall, Stanford, CA 94305, USA
- 33 ⁶⁴Laboratory of Behavioral Neurobiology, Department of Physiological Sciences, Federal University of
34 Santa Catarina, Florianópolis, Campus Universitário Trindade, 88040900, Brazil
- 35 ⁶⁵Department of Kinesiology, KU Leuven, Leuven, Tervuursevest 101 box 1501, B-3001 Leuven,
36 Belgium
- 37 ⁶⁶Department of Experimental Psychology, University of Oxford, Oxford, UK
- 38 ⁶⁷Department of Psychology, UiT The Arctic University of Norway, Tromsø, Norway
- 39 ⁶⁸Department of Psychology, DePaul University, Chicago, 2219 N Kenmore Ave, Chicago, IL 60657,
40 USA

- 1 ⁶⁹Center for Neural Science, New York University, 4 Washington PI Room 809 New York, NY 10003,
2 USA
- 3 ⁷⁰Department of Psychology, Boğaziçi University, Bebek, 34342, Istanbul, Turkey
- 4 ⁷¹Psychology, University of Cologne, Cologne, Herbert-Lewin-St. 2, 50931, Cologne, Germany
- 5 ⁷²Saudi Human Genome Program, King Abdulaziz City for Science and Technology (KACST);
6 Integrated Gulf Biosystems, Riyadh, Saudi Arabia
- 7 ⁷³Cognitive Psychology Unit, Institute of Psychology, Leiden University, Wassenaarseweg 52, 2333
8 AK Leiden, The Netherlands
- 9 ⁷³Leiden Institute for Brain and Cognition, Leiden University, Leiden, The Netherlands
- 10 ⁷⁴Leeds Institute of Health Sciences, University of Leeds, Leeds, LS2 9NL, UK
- 11 ⁷⁵Institute for Media Studies, KU Leuven, Leuven, Belgium
- 12 ⁷⁶Center for Open Science, 210 Ridge McIntire Rd Suite 500, Charlottesville, VA 22903, USA
- 13 ⁷⁶Department of Engineering and Society, University of Virginia, Thornton Hall, P.O. Box 400259,
14 Charlottesville, VA 22904, USA
- 15 ⁷⁷Surgical Neurology Branch, National Institute of Neurological Disorders and Stroke, National
16 Institutes of Health, Bethesda, MD 20892, USA
- 17 ⁷⁸Department of Psychology, University of Southern Indiana, 8600 University Boulevard, Evansville,
18 Indiana, USA
- 19 ⁷⁹Life Sciences Communication, University of Wisconsin-Madison, Madison, Wisconsin, 1545
20 Observatory Drive, Madison, WI 53706, USA
- 21 ⁸⁰Department of Social Psychology, Institute of Psychology, University of Lausanne, Quartier UNIL-
22 Mouline, Bâtiment Géopolis, CH-1015 Lausanne, Switzerland
- 23 ⁸¹Departamento de Psicología Básica, Universidad Autónoma de Madrid, c/ Ivan Pavlov 6, 28049
24 Madrid, Spain
- 25 ⁸²Department of Methodology and Statistics, Tilburg University, Warandelaan 2, 5000 LE Tilburg, The
26 Netherlands
- 27 ⁸²Department of Sociology, Utrecht University, Padualaan 14, 3584 CH, Utrecht, The Netherlands
- 28 ⁸³School of Psychology, Massey University, Auckland, Private Bag 102904, North Shore, Auckland,
29 0745, New Zealand
- 30 ⁸⁴Psychology, Saint Louis University, St. Louis, MO, 3700 Lindell Blvd, St. Louis, MO 63108, USA
- 31 ⁸⁵Psychology, University of California, Davis, Davis, One Shields Ave, Davis, CA 95616, USA
- 32 ⁸⁶Department of Psychology, University of Texas at Austin, 108 E. Dean Keeton Stop A8000, Austin,
33 TX 78712-1043, USA
- 34 ⁸⁷Marketing Department, Ghent University, Tweekerkenstraat 2, 9000 Ghent, Belgium
- 35 ⁸⁸Department of Psychology, Education, and Child Studies, Erasmus University Rotterdam,
36 Rotterdam, Burgemeester Oudlaan 50, 3000 DR, Rotterdam, The Netherlands

37

38 **Author Note:** We'd like to thank Dale Barr, Felix Cheung, David Colquhoun, Hans IJzerman,
39 Harvey Motulsky, and Richard Morey for helpful discussions while drafting this commentary.

40

1 **Funding Statement:** Daniel Lakens was supported by NWO VIDI 452-17-013. Federico G.
2 Adolphi was supported by CONICET. Matthew Apps was funded by a Biotechnology and
3 Biological Sciences Research Council AFL Fellowship (BB/M013596/1). Gary Collins was
4 supported by the NIHR Biomedical Research Centre, Oxford. Zander Crook was supported
5 by the Economic and Social Research Council [grant number C106891X]. Emily S. Cross
6 was supported by the European Research Council (ERC-2015-StG-677270). Lisa DeBruine
7 is supported by the European Research Council (ERC-2014-CoG-647910 KINSHIP). Anne-
8 Laura van Harmelen is funded by a Royal Society Dorothy Hodgkin Fellowship (DH150176).
9 Mark R. Hoffarth was supported by the National Science Foundation under grant SBE
10 SPRF-FR 1714446. Junpeng Lao was supported by the SNSF grant 100014_156490/1.
11 Cilene Lino de Oliveira was supported by AvH, Capes, CNPq. Andrea E. Martin was
12 supported by the Economic and Social Research Council of the United Kingdom [grant
13 number ES/K009095/1]. Jean-Jacques Orban de Xivry is supported by an internal grant from
14 the KU Leuven (STG/14/054) and by the Fonds voor Wetenschappelijk Onderzoek
15 (1519916N). Sam Parsons was supported by the European Research Council (FP7/2007–
16 2013; ERC grant agreement no; 324176). Gerine Lodder was funded by NWO VICI 453-14-
17 016. Samuel Smith is supported by a Cancer Research UK Fellowship (C42785/A17965).
18 Vishnu Sreekumar was supported by the NINDS Intramural Research Program (IRP). Miguel
19 A. Vadillo was supported by Grant 2016-T1/SOC-1395 from Comunidad de Madrid. Tal
20 Yarkoni was supported by NIH award R01MH109682.

21

22 **Abstract:** In response to recommendations to redefine statistical significance to $p \leq .005$, we
23 propose that researchers should transparently report and justify all choices they make when
24 designing a study, including the alpha level.

25

1 **Justify Your Alpha: A Response to “Redefine Statistical Significance”**

2
3 *“Tests should only be regarded as tools which must be used with discretion and*
4 *understanding, and not as instruments which in themselves give the final verdict.”*

5 Neyman & Pearson, 1928, p. 58.

6
7 Renewed concerns about the non-replication of scientific findings have prompted
8 widespread debates about its underlying causes and possible solutions. As an actionable
9 step toward improving standards of evidence for new discoveries, 72 researchers proposed
10 changing the conventional threshold that defines “statistical significance” (i.e., the alpha
11 level) from $p \leq .05$ to $p \leq .005$ for all novel claims with relatively low prior odds (Benjamin et
12 al., 2017). They argued that this change will “immediately improve the reproducibility of
13 scientific research in many fields” (Benjamin et al., 2017, p. 5).

14
15 Benjamin et al. (2017) provided two arguments against the current threshold for statistical
16 significance of .05. First, a p -value of .05 provides only weak evidence for the alternative
17 hypothesis. Second, under certain assumptions, a p -value threshold of .05 leads to a high
18 false positive report probability (FPRP; the probability that a significant finding is a false
19 positive, Wacholder et al., 2004; also referred to as the false positive rate, or false positive
20 risk, Benjamin et al., 2017; Colquhoun, 2017). The authors claim that lowering the threshold
21 for statistical significance to .005 will increase evidential strength for novel discoveries and
22 reduce the FPRP.

23
24 We share the concerns raised by Benjamin et al. (2017) regarding the apparent non-
25 replicability¹ of many scientific studies and appreciate their attempt to provide a concrete,
26 easy-to-implement suggestion to improve science. We further agree that the current default
27 alpha level of .05 is arbitrary and may result in weak evidence for the alternative hypothesis.
28 However, we do not think that redefining the threshold for statistical significance to the lower,
29 but equally arbitrary threshold of $p \leq .005$ is advisable. In this commentary, we argue that (1)
30 there is insufficient evidence that the current standard for statistical significance is in fact a
31 “leading cause of non-reproducibility” (Benjamin et al., 2017, p. 5), (2) the arguments in favor
32 of a blanket default of $p \leq .005$ are not strong enough to warrant the immediate and
33 widespread implementation of such a policy, and (3) a lower significance threshold will likely
34 have positive and negative consequences, both of which should be carefully evaluated

¹ We use ‘replicability’ to refer to the question of whether a conclusion that is sufficiently similar to an earlier study could be drawn from data obtained from a new study, and ‘reproducibility’ to refer to getting the same results when re-analysing the same data (Peng, 2009).

1 before any large-scale changes are proposed. We conclude with an alternative suggestion,
2 whereby researchers justify their choice for an alpha level before collecting the data, instead
3 of adopting a new uniform standard.

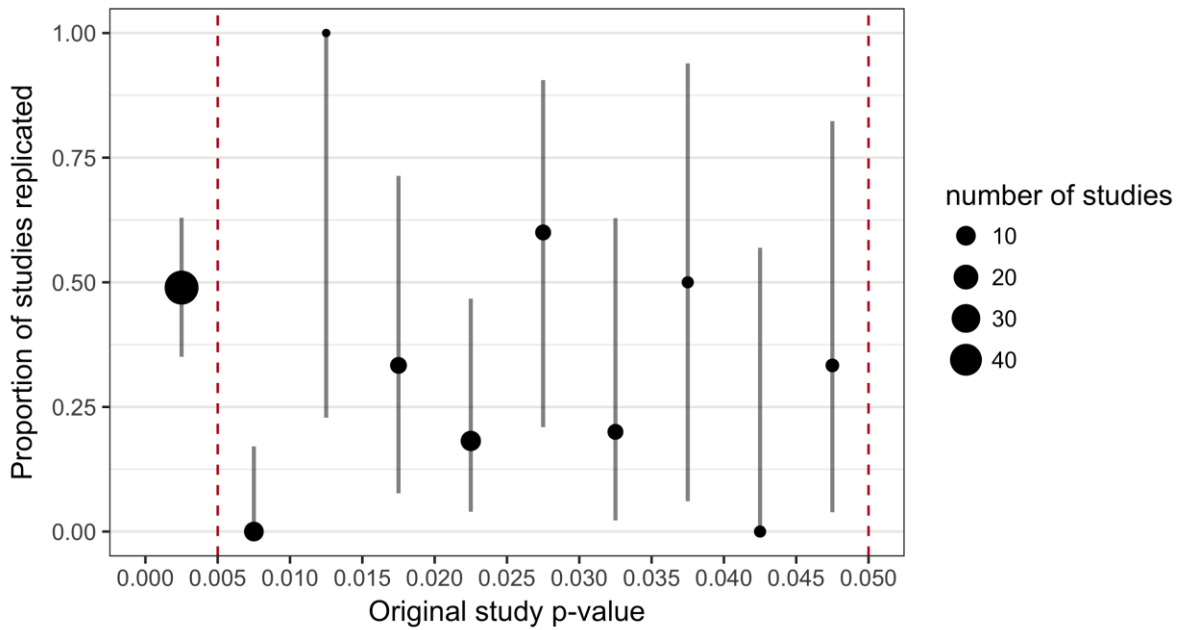
5 ***Lack of evidence that $p \leq .005$ improves replicability***

7 One of the main claims made by Benjamin et al. (2017) is that the expected proportion of
8 studies that can be replicated will be considerably higher for studies that observe $p \leq .005$
9 than for studies that observe $.005 < p \leq .05$, due to a lower FPRP. All else being equal, we
10 agree with Benjamin et al. (2017) that improvement in replicability is *in theory* related to the
11 FPRP, and that lower alpha levels will reduce false positive results in the literature. However,
12 it is difficult to predict how much the FPRP will change *in practice*, because quantifying the
13 FPRP requires accurate estimates of several unknowns, such as the prior odds that the
14 examined hypotheses are true, the true power of any performed experiments, and the
15 (change in) actual behavior of researchers should the newly proposed threshold be put in
16 place.

18 An analysis of the results of the Reproducibility Project: Psychology (RP:P; Open Science
19 Collaboration, 2015) shows that 49% (23 out of 47) of the original findings with p -values
20 below $.005$ yielded $p \leq .05$ in the replication study, whereas only 24% (11 out of 45) of the
21 original studies with $.005 < p \leq .05$ yielded $p \leq .05$ in the replication study ($\chi^2(1) = 5.92$, $p =$
22 $.015$, $BF_{10} = 6.84$). Benjamin et al. (2017, p. 9) presented this analysis as empirical evidence
23 of the “potential gains in reproducibility that would accrue from the new threshold.” However,
24 as they acknowledged, their obtained p -value of $.015$ is only “suggestive” of such a
25 conclusion, according to their own proposal. Moreover, there is considerable variation in
26 replication rates across p -values (see Figure 1), with few observations in bins of size $.005$ for
27 $.005 < p \leq .05$. In addition, the lower replication rate for p -values just below $.05$ is likely
28 confounded by p -hacking (the practice of flexibly analysing data until the p -value passes the
29 ‘significance’ threshold) in the original study. This implies that at least some of the
30 differences in replication rates between studies with $.005 < p \leq .05$ compared to studies with
31 $p \leq .005$ are not due to the level of evidence *per se*, but rather due to other mechanisms
32 (e.g., flexibility during data analysis). Indeed, depending on the degree of flexibility exploited
33 by researchers, such p -hacking can be used to overcome any inferential threshold.

35 Even with a $p \leq .005$ threshold, only 49% of studies replicated successfully. Furthermore,
36 only 11 out of 30 studies (37%) with $.0001 < p \leq .005$ replicated at $\alpha = .05$. By contrast, a
37 *prima facie* more satisfactory replication success rate of 71% was obtained only for $p <$

1 .0001 (12 out of 17 studies). This suggests that a relatively small number of studies with p -
 2 values much lower than .005 were largely responsible for the 49% replication rate for studies
 3 with $p \leq .005$. Further analysis is needed, therefore, to explain the low replication rate of
 4 studies with $p \leq .005$ before this alpha level is recommended as a new significance threshold
 5 for novel discoveries across scientific disciplines.



6
 7 *Figure 1.* The proportion of studies (Open Science Collaboration, 2015) that replicated at $\alpha =$
 8 .05 (with a bin width of 0.005). Window start and end positions are plotted on the horizontal
 9 axis. The error bars denote 95% Jeffreys confidence intervals. R code to reproduce Figure 1
 10 is available from <https://github.com/VishnuSreekumar/Alpha005>

11

12 ***Weak justifications for the new $p \leq .005$ threshold***

13

14 Even though p -values close to .05 never provide strong ‘evidence’ against the null
 15 hypothesis on their own (Wasserstein & Lazar, 2016), the argument that p -values provide
 16 weak evidence based on Bayes factors has been called into question (Casella & Berger,
 17 1987; Greenland et al., 2016; Senn, 2001). Redefining the alpha level as a function of the
 18 strength of relative evidence measured by the Bayes factor is undesirable, given that the
 19 marginal likelihood is very sensitive to different (somewhat arbitrary) choices for the models
 20 that are compared (Gelman et al., 2013). Benjamin et al. (2017) stated that p -values of .005
 21 imply Bayes factors between 14 and 26, but the level of evidence depends on the model
 22 priors and the choice of hypotheses tested, and different modelling assumptions would imply
 23 a different p -value threshold. The Bayesian analysis that underlies the recommendation
 24 actually overstates the evidence against the null from the perspective of error statistics. It

1 would, with high probability, deem an alternative highly probable, even if it's false (Mayo,
2 1997, 2018). Finally, Benjamin et al. (2017) provided no rationale for why the new p -value
3 threshold *should* align with equally arbitrary Bayes factor thresholds representing
4 'substantial' or 'strong' evidence. Indeed, it has been argued that such classifications of
5 Bayes factors themselves introduce arbitrary meaning to a continuous measure (e.g., Morey,
6 2015). We (even those of us prepared to use likelihoods and Bayesian approaches in lieu of
7 p -values when interpreting results) caution against the idea that the alpha level at which an
8 error rate is controlled should be based on the amount of relative evidence indicated by a
9 Bayes factor. Extending Morey, Wagenmakers, and Rouder (2016), who argued against the
10 frequentist calibration of Bayes factors, we argue against the necessity of a Bayesian
11 calibration of error rates.

12

13 The second argument Benjamin et al. (2017) provided for $p \leq .005$ is that the FPRP can be
14 high with $\alpha = .05$. To calculate the FPRP one needs to define the alpha level, the power of
15 the tests that examine true effects, and the ratio of true to false hypotheses tested (the prior
16 odds). The FPRP is only problematic when a high proportion of examined hypotheses are
17 false, and thus Benjamin et al. (2017, p. 10) stated that their "recommendation applies to
18 disciplines with prior odds broadly in the range depicted in Figure 2." Their Figure 2 displays
19 FPRPs for scenarios where many examined hypotheses are false, with ratios of true to false
20 hypotheses (i.e., prior odds) of 1 to 5, 1 to 10, and 1 to 40. Benjamin et al. (2017)
21 recommended $p \leq .005$ because this threshold reduces the *minimum* FPRP to less than 5%,
22 assuming 1 to 10 prior odds of examining a true hypothesis (the true FPRP might still be
23 substantially higher in studies with very low power). This estimate of prior odds is based on
24 data from the RP:P (Open Science Collaboration, 2015) using an analysis that modelled
25 publication bias for 73 studies (Johnson et al., 2017; see also Ingre, 2016, for a more
26 conservative estimate). Without stating the reference class for the 'base-rate of true nulls'
27 (i.e., does this refer to all hypotheses in science, in a discipline, or by a single researcher?),
28 the concept of 'prior odds that H1 is true' has little meaning in practice. The modelling effort
29 by Johnson et al. (2017) ignored practices that inflate error rates (e.g., p -hacking) and thus
30 likely does not provide an accurate estimate of bias, given the prevalence of such practices
31 (Fiedler & Schwarz, 2016; John et al., 2012). An estimate of the prior probability that a
32 hypothesis is true, similar to that of Johnson et al. (2017), was derived from 92 participants'
33 subjective ratings of the prior probability that the alternative hypothesis was true for 44
34 studies included in the RP:P (Dreber et al., 2015). As Dreber et al. (2015, p. 15345) noted,
35 "This relatively low average prior may reflect [the fact] that top psychology journals focus on
36 publishing surprising findings, i.e., positive findings on relatively unlikely hypotheses." These
37 observations imply that there are not sufficient representative data to accurately estimate the

1 prior odds that researchers examine a true hypothesis, and thus, there is currently no strong
2 argument based on FPRP to redefine statistical significance to $p \leq .005$.

3

4 ***Ways in which a threshold of $p \leq .005$ might harm scientific practice***

5

6 Benjamin et al. (2017) acknowledged that lowering the p -value threshold will not ameliorate
7 other practices that negatively impact the replicability of research findings (such as p -
8 hacking, publication bias, and low power). Yet, they did not address ways in which a $p \leq .005$
9 threshold might harm scientific practice. Chief among our concerns are (1) a reduction in the
10 number of replication studies that can be conducted if such a threshold is adopted, (2) a
11 concomitant reduction in generalisability and breadth of research findings due to a likely
12 increased reliance on convenience samples, and (3) exacerbation of an already exaggerated
13 focus on single p -values.

14

15 *Risk of fewer replication studies.* Replication studies are central to generating reliable
16 scientific knowledge, especially when conclusions are largely based on p -values. As Fisher
17 (1926, p. 85) noted: "A scientific fact should be regarded as experimentally established only
18 if a properly designed experiment *rarely fails* to give this level of significance." Replication
19 studies are at the heart of scientific progress. In the field of medicine, for example, the FDA
20 requires two independent pre-registered clinical trials, both significant with $p \leq .05$, before
21 issuing marketing approval for new drugs (for a discussion, see Senn, 2007, p. 188).

22 Researchers have limited resources, and when studies require larger sample sizes scientists
23 will have to decide what research they will invest in. Achieving 80% power with $\alpha = .005$,
24 compared to $\alpha = .05$, will require a 70% larger sample size in a between-subjects design with
25 a two-sided test (and an 88% larger sample size for one-sided tests). This means that
26 researchers can complete almost two studies each powered at $\alpha = .05$ (e.g., one novel study
27 and one replication study), or only one study powered at $\alpha = .005$. Therefore, at a time when
28 replication studies are rare, lowering the alpha level to .005 might reduce the number of
29 replication studies. Indeed, general recommendations for evidence thresholds need to
30 carefully balance statistical and non-statistical considerations (e.g., the value of evidence per
31 novel study vs. the value of independent replications).

32

33 *Risk of reduced generalisability and breadth.* All things equal, larger sample sizes increase
34 the informational value of studies, but requiring larger sample sizes across all scientific
35 disciplines would potentially compound problems with over-reliance on convenience samples
36 (such as undergraduate students or Mechanical Turk workers). Lowering the significance
37 threshold could adversely affect the type of breadth of research questions examined if it is

1 done without (1) increased funding, (2) a reward system that values large-scale
2 collaboration, or (3) clear recommendations for how to evaluate research with lower
3 evidential value due to sample size constraints. Achieving a lower p -value in studies with
4 unique populations (e.g., people with rare genetic variants, people diagnosed with post-
5 traumatic stress disorder) or in studies with time- or otherwise resource-intensive data
6 collection (e.g., longitudinal studies) requires exponentially more effort than increasing the
7 amount of evidence in studies that use undergraduate students or Mechanical Turk workers.
8 Thus, researchers may become less motivated, or even tacitly discouraged, to study the
9 former populations or collect those types of data. Hence, lowering the alpha threshold may
10 indirectly reduce the generalisability and breadth of findings (Peterson & Merunka, 2014).

11

12 *Risk of exaggerating the focus on single p -values.* If anything, an excessive focus on p -value
13 thresholds has the potential to mask or even discourage opportunities for more fruitful
14 changes in scientific practice and education. Many researchers have come to recognise p -
15 hacking, low power, and publication bias as more important reasons for non-replication.
16 Benjamin et al. (2017) acknowledged that changing the threshold could be considered a
17 distraction from other solutions, and yet their proposal risks reinforcing the idea that relying
18 only on p -values is a sufficient, if imperfect, way to evaluate findings. The proposed $p \leq .005$
19 threshold is not intended as a publication threshold. However, given the long history of
20 misuse of statistical recommendations, there is a substantial risk that redefining $p \leq .005$ as
21 'statistically significant' will increase publication bias, which, in turn, would bias effect size
22 estimates upwards to an even greater extent (Lane & Dunlap, 1978). As such, Benjamin et
23 al.'s recommendation could divert attention from the burgeoning movement towards a more
24 cumulative evaluation of findings, where the converging results of multiple studies are taken
25 into account when addressing specific research questions. Examples of such approaches
26 are: multiple replications (both registered and multi-lab; see, e.g., Hagger et al., 2016),
27 continuously updating meta-analyses (Braver et al., 2014), p -curve analysis (Simonsohn et
28 al., 2014), and pre-registration of studies.

29

30 ***No one alpha to rule them all***

31

32 Benjamin et al. (2017) recommended that only p -values lower than .005 should be called
33 'statistically significant' and that studies should generally be designed with $\alpha = .005$. Our
34 recommendation is similarly twofold. First, when describing results, we recommend that the
35 label 'statistically significant' simply no longer be used. Instead, researchers should provide
36 a more meaningful interpretation (Eysenck, 1960). While p -values can inspire statements
37 about the probability of data (e.g., 'the observed difference in the data was surprisingly large,

1 assuming the null hypothesis is true'), they should not be treated as indices that, on their
2 own, signify evidence for a theory.

3

4 Second, when designing studies, we propose that authors transparently specify their design
5 choices. These include (where applicable) the alpha level, the null and alternative models,
6 assumed prior odds, statistical power for a specified effect size of interest, the sample size,
7 and/or the desired accuracy of estimation. Without imposing a single value on any of these
8 design parameters, we ask authors to *justify their choices* before the data are collected.

9 Fellow researchers can evaluate these decisions on their merits and discuss how
10 appropriate they are for a specific research question, and whether the conclusions follow
11 from the study design. Ideally, this evaluation process occurs prior to data collection when
12 reviewing a Registered Report submission (Chambers, Dienes, McIntosh, Rotshtein, &
13 Willmes, 2015). Providing researchers (and reviewers) with accessible information on ways
14 to justify (and evaluate) these design choices, tailored to specific research areas, would
15 improve current research practices.

16

17 The optimal alpha level will sometimes be lower and sometimes be higher than the current
18 convention of .05 (see Field, Tyre, Jonzén, Rhodes, & Possingham, 2004; Grieves, 2015;
19 Mudge, Baker, Edge, & Houlahan, 2012; Pericchi & Pereira, 2016). Some fields, such as
20 genomics and physics, have lowered the alpha level. However, in genomics the overall false
21 positive rate is still controlled at 5%; the lower alpha level is only used to correct for multiple
22 comparisons (Storey & Tibshirani, 2003). In physics, a five sigma threshold ($p \leq 2.87 \times 10^{-7}$)
23 is required to publish an article with 'discovery of' in the title, with less stringent alpha levels
24 being used for article titles with 'evidence for' or 'measurement of' (Franklin, 2014). In
25 physics researchers have also argued against a blanket rule, and instead setting the alpha
26 level based on factors such as how surprising the result would be and how much practical or
27 theoretical impact the discovery would have (Lyons, 2013). In non-human animal research,
28 minimising the number of animals used needs to be directly balanced against the probability
29 of false positives; other trade-offs may be relevant in other areas. Thus, a broadly applied p
30 $\leq .005$ threshold will rarely be optimal.

31

32 Benjamin et al. (2017, p. 5) stated that a "critical mass of researchers" now endorse the
33 standard of a $p \leq .005$ threshold for "statistical significance." However, the presence of a
34 critical mass can only be identified *after* a norm or practice has been widely adopted, not
35 *before*. Even if a $p \leq .005$ threshold was widely endorsed, this would only reinforce the
36 flawed idea that a single alpha level is universally applicable. Ideally, the decision of where
37 to set the alpha level for a study should be based on statistical decision theory, where costs

1 and benefits are compared against a utility function (Neyman & Pearson, 1933; Skipper,
2 Guenther, & Nass, 1967). Such an analysis can be expected to differ based on the type of
3 study being conducted: for example, analysis of a large existing dataset versus primary data
4 collection relying on hard-to-obtain samples. Science is necessarily diverse, and it is up to
5 scientists within specific fields to justify the alpha level they decide to use. To quote Fisher
6 (1956, p. 42): "...no scientific worker has a fixed level of significance at which, from year to
7 year, and in all circumstances, he rejects hypotheses; he rather gives his mind to each
8 particular case in the light of his evidence and his ideas."
9

10 **Conclusion**

11
12 It is laudable that Benjamin et al. (2017) suggested a concrete step designed to immediately
13 improve science. However, it is not clear that lowering the significance threshold to $p \leq .005$
14 will in practice amount to an improvement in replicability that is worth the potential costs.
15 Instead of simple heuristics and an arbitrary blanket threshold, research should be guided by
16 principles of rigorous science (Casadevall & Fang, 2016; LeBel, Vanpaemel, McCarthy,
17 Earp, & Elson, 2017; Meehl, 1990). These principles include not only sound statistical
18 analyses, but also experimental redundancy (e.g., replication, validation, and generalisation),
19 avoidance of logical traps, intellectual honesty, research workflow transparency, and full
20 accounting for potential sources of error. Single studies, regardless of their p -value, are
21 never enough to conclude that there is strong evidence for a *theory*. We need to train
22 researchers to recognise what cumulative evidence looks like and work towards an unbiased
23 scientific literature.
24

25 Although we agree with Benjamin et al. (2017) that the relatively high rate of non-replication
26 in the scientific literature is a cause for concern, we do not believe that redefining statistical
27 significance is a desirable solution: (1) there is not enough evidence that a blanket threshold
28 of $p \leq .005$ will improve replication sufficiently to be worth the additional cost in data
29 collection, (2) the justifications given for the new threshold are not strong enough to warrant
30 the widespread implementation of such a policy, and (3) there are realistic concerns that a p
31 $\leq .005$ threshold will have negative consequences for science, which should be carefully
32 examined before a change in practice is instituted. Instead of a narrower focus on p -value
33 thresholds, we call for a broader mandate whereby all *justifications* of key choices in
34 research design and statistical practice are pre-registered whenever possible, fully
35 accessible, and transparently evaluated.
36
37

References

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37

- Benjamin, D. J., Berger, J., Johannesson, M., Nosek, B. A., Wagenmakers, E.-J., Berk, R., ... Johnson, V. (2017, July 22). Redefine statistical significance. <https://doi.org/10.17605/OSF.IO/MKY9J>
- Braver, S. L., Thoemmes, F. J., & Rosenthal, R. (2014). Continuously cumulating meta-analysis and replicability. *Perspectives on Psychological Science*, 9(3), 333-342. <https://doi.org/10.1177/1745691614529796>
- Casella, G., & Berger, R. L. (1987). Testing Precise Hypotheses: Comment. *Statistical Science*, 2(3), 344–347. <https://doi.org/10.1214/ss/1177013243>
- Casadevall, A., & Fang, F. C. (2016). Rigorous Science: a How-To Guide. *mBio*, 7(6), e01902-16. <https://doi.org/10.1128/mBio.01902-16>
- Chambers, C.D., Dienes, Z., McIntosh, R.D., Rotshtein, P., & Willmes, K. (2015). Registered Reports: Realigning incentives in scientific publishing. *Cortex*, 66, A1-2. <https://doi.org/10.1016/j.cortex.2015.03.022>
- Colquhoun, D. (2017). The reproducibility of research and the misinterpretation of p -values. *bioRxiv*, 144337. <https://doi.org/10.1101/144337>
- Dreber, A., Pfeiffer, T., Almenberg, J., Isaksson, S., Wilson, B., Chen, Y., ... Johannesson, M. (2015). Using prediction markets to estimate the reproducibility of scientific research. *Proceedings of the National Academy of Sciences*, 112(50), 15343–15347. <https://doi.org/10.1073/pnas.1516179112>
- Eysenck, H. J. (1960). The concept of statistical significance and the controversy about one-tailed tests. *Psychological Review*, 67(4), 269-271. <http://dx.doi.org/10.1037/h0048412>
- Fiedler, K., & Schwarz, N. (2016). Questionable Research Practices Revisited. *Social Psychological and Personality Science*, 7(1), 45–52. <http://doi.org/10.1177/1948550615612150>
- Field, S. A., Tyre, A. J., Jonzen, N., Rhodes, J. R., & Possingham, H. P. (2004). Minimizing the cost of environmental management decisions by optimizing statistical thresholds. *Ecology Letters*, 7(8), 669-675. <https://doi.org/10.1111/j.1461-0248.2004.00625.x>
- Fisher, R. A. (1926). The arrangement of field experiments. *Journal of the Ministry of Agriculture of Great Britain*, 33, 503-513.
- Fisher R. A. (1956). *Statistical Methods and Scientific Inferences*. Hafner: New York.
- Franklin, A. (2014). *Shifting standards: Experiments in particle physics in the twentieth century*. University of Pittsburgh Press.
- 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

- 1 guide to misinterpretations. *European Journal of Epidemiology*, 31(4), 337–350.
2 <https://doi.org/10.1007/s10654-016-0149-3>
- 3 Grieve, A. P. (2015). How to test hypotheses if you must. *Pharmaceutical Statistics*, 14(2),
4 139–150. <https://doi.org/10.1002/pst.1667>
- 5 Hagger, M. S., Chatzisarantis, N. L. D., Alberts, H., Anggono, C. O., Batailler, C., Birt, A. R.,
6 . . . Zwieneberg, M. (2016). A multilab preregistered replication of the ego-depletion
7 effect. *Perspectives on Psychological Science*, 11, 546–573.
8 <https://doi.org/10.1177/17456916166652873>
- 9 Ingre, M. (2016). Recent reproducibility estimates indicate that negative evidence is
10 observed over 30 times before publication. *arXiv preprint arXiv:1605.06414*.
11 <https://arxiv.org/abs/1605.06414>
- 12 John, L. K., Loewenstein, G., & Prelec, D. (2012). Measuring the prevalence of questionable
13 research practices with incentives for truth-telling. *Psychological Science*, 23(5),
14 524–532. <https://doi.org/10.2139/ssrn.1996631>
- 15 Johnson, V. E., Payne, R. D., Wang, T., Asher, A., & Mandal, S. (2017). On the
16 reproducibility of psychological science. *Journal of the American Statistical*
17 *Association*, 112(517), 1–10. <https://doi.org/10.1080/01621459.2016.1240079>
- 18 Koole, S. L., & Lakens, D. (2012). Rewarding replications: A sure and simple way to improve
19 psychological science. *Perspectives on Psychological Science*, 7(6), 608-614.
20 <https://doi.org/10.1177/1745691612462586>
- 21 Lakens, D. (2015). On the challenges of drawing conclusions from p-values just below 0.05.
22 *PeerJ*, 3, e1142. <https://doi.org/10.7717/peerj.1142>
- 23 Lane, D. M., & Dunlap, W. P. (1978). Estimating effect size: Bias resulting from the
24 significance criterion in editorial decisions. *British Journal of Mathematical and*
25 *Statistical Psychology*, 31(2), 107-112. [https://doi.org/10.1111/j.2044-](https://doi.org/10.1111/j.2044-8317.1978.tb00578.x)
26 [8317.1978.tb00578.x](https://doi.org/10.1111/j.2044-8317.1978.tb00578.x)
- 27 LeBel, E. P., Vanpaemel, W., McCarthy, R. J., Earp, B. D., & Elson, M. (2017). A Unified
28 Framework to Quantify the Trustworthiness of Empirical Research.
29 <https://doi.org/10.17605/OSF.IO/UWMR8>
- 30 Lyons, L. (2013). Discovering the Significance of 5 sigma. *arXiv preprint arXiv:1310.1284*.
- 31 Meehl, P. E. (1990). Appraising and amending theories: The strategy of Lakatosian defense
32 and two principles that warrant it. *Psychological Inquiry*, 1(2), 108–141.
33 https://doi.org/10.1207/s15327965pli0102_1
- 34 Mayo, D. (1997). Error statistics and learning from error: Making a virtue of necessity.
35 *Philosophy of Science*, Vol. 64, Part II: Symposia Papers, S195-S212.
- 36 Mayo, D. (2018). *Statistical Inference as Severe Testing: How to Get Beyond the Statistics*
37 *Wars*. Cambridge University Press.

- 1 Morey, R. (2015). On verbal categories for the interpretation of Bayes Factors. *BayesFactor*.
2 <https://bayesfactor.blogspot.nl/2015/01/on-verbal-categories-for-interpretation.html>
- 3 Morey, R. D., Wagenmakers, E.-J., & Rouder, J. N. (2016). Calibrated Bayes factors
4 should not be used: A reply to Hoijtink, van Kooten, and Hulsker. *Multivariate Behavioral*
5 *Research*, 51(1), 11–19. <http://dx.doi.org/10.1080/00273171.2015.1052710>
- 6 Mudge, J. F., Baker, L. F., Edge, C. B., & Houlahan, J. E. (2012). Setting an Optimal α That
7 Minimizes Errors in Null Hypothesis Significance Tests. *PLOS ONE*, 7(2), e32734.
8 <https://doi.org/10.1371/journal.pone.0032734>
- 9 Neyman, J., & Pearson, E. S. (1928). On the use and interpretation of certain test criteria for
10 purposes of statistical inference: Part I. *Biometrika*, 175-240.
11 <https://doi.org/10.2307/2331945>
- 12 Neyman, J., & Pearson, E. S. (1933). On the Problem of the Most Efficient Tests of
13 Statistical Hypotheses. *Philosophical Transactions of the Royal Society of London A:*
14 *Mathematical, Physical and Engineering Sciences*, 231(694–706), 289–337.
15 <https://doi.org/10.1098/rsta.1933.0009>
- 16 Nosek, B. A., Spies, J. R., & Motyl, M. (2012). Scientific Utopia II: Restructuring Incentives
17 and Practices to Promote Truth Over Publishability. *Perspectives on Psychological*
18 *Science*, 7(6), 615-631. <http://doi.org/10.1177/1745691612459058>
- 19 Open Science Collaboration. (2015). Estimating the reproducibility of psychological science.
20 *Science*, 349(6251), aac4716. <https://doi.org/10.1126/science.aac4716>
- 21 Peng, R. D. (2009). Reproducible research and biostatistics. *Biostatistics*, 10(3), 405–408.
22 <https://doi.org/10.1093/biostatistics/kxp014>
- 23 Pericchi, L., & Pereira, C. (2016). Adaptive significance levels using optimal decision rules:
24 Balancing by weighting the error probabilities. *Brazilian Journal of Probability and*
25 *Statistics*, 30(1), 70–90. <https://doi.org/10.1214/14-BJPS257>
- 26 Peterson, R. A., & Merunka, D. R. (2014). Convenience samples of college students and
27 research reproducibility. *Journal of Business Research*, 67(5), 1035-1041.
28 <https://doi.org/10.1016/j.jbusres.2013.08.010>
- 29 Senn, S. (2001) Two cheers for p-values? *Journal of Epidemiology and Biostatistics*, 6, 193-
30 204. <https://doi.org/10.1080/135952201753172953>
- 31 Senn, S. (2007). *Statistical issues in drug development* (2nd ed). Chichester, England;
32 Hoboken, NJ: John Wiley & Sons.
- 33 Simonsohn, U., Nelson, L. D., & Simmons, J. P. (2014). P-curve and effect size: Correcting
34 for publication bias using only significant results. *Perspectives on Psychological*
35 *Science*, 9, 666-681. <https://doi.org/10.1177/1745691614553988>

- 1 Skipper, J. K., Guenther, A. L., & Nass, G. (1967). The sacredness of .05: A note concerning
2 the uses of statistical levels of significance in social science. *The American*
3 *Sociologist*, 2(1), 16–18.
- 4 Storey, J. D., & Tibshirani, R. (2003). Statistical significance for genomewide studies.
5 *Proceedings of the National Academy of Sciences*, 100(16), 9440–9445.
6 <https://doi.org/10.1073/pnas.1530509100>
- 7 Wacholder, S., Chanock, S., Garcia-Closas, M., El Ghormli, L., & Rothman, N. (2004).
8 Assessing the probability that a positive report is false: An approach for molecular
9 epidemiology studies. *Journal of the National Cancer Institute*, 96, 434-442.
10 <https://doi.org/10.1093/jnci/djh075>
- 11 Wasserstein, R. L., & Lazar, N. A. (2016). The ASA's Statement on p-Values: Context,
12 Process, and Purpose. *The American Statistician*, 70(2), 129–133.
13 <https://doi.org/10.1080/00031305.2016.1154108>