



Subject Areas:

Cognition and decision making,
Psychology and cognitive
neuroscience, Statistics

Keywords:

reasoning, statistics, statistical
cognition, decision making,
significance testing

Author for correspondence:

Richard D. Morey

e-mail:

moreyr@cardiff.ac.uk

Use of significance test logic by scientists in a novel reasoning task

Richard D. Morey¹, Rink Hoekstra²

¹School of Psychology, Cardiff University

²Faculty of Behavioural and Social Sciences,
University of Groningen

Although statistical significance testing is one of the most widely-used techniques across science, previous research has suggested that scientists have a poor understanding of how it works. If scientists misunderstand one of their primary inferential tools the implications are dramatic: potentially unchecked, unjustified conclusions and wasted resources. Scientists' apparent difficulties with significance testing have led to calls for its abandonment or increased reliance on alternative tools, which would represent a substantial, untested, shift in scientific practice. However, if scientists' understanding of significance testing is truly as poor as thought, one could argue such drastic action is required. We show using a novel experimental method that scientists do, in fact, understand the logic of significance testing and can use it effectively. This suggests that scientists may not be as statistically-challenged as often believed, and that reforms should take this into account.

For most of the past century, the dominant method of statistical inference has been statistical significance testing (SST). In a significance test, the statistical evidence in the form of a test statistic is compared to what would be expected under a particular hypothesis (often called the “null” hypothesis). If it would be surprising to observe evidence as strong as what was observed under this hypothesis, the evidence is strong enough to call the assumed hypothesis into doubt, at least tentatively [see also 1,2]. The rarity of evidence as strong as what was observed under the assumed hypothesis—the so-called p value—is the typical way that results of significance tests are reported. The key feature of SST for our purposes is the assessment of evidence by means of comparing a result to a “null” distribution.

Despite the use of SST in a majority of research projects across fields, there is debate over whether scientists understand SST and can use it competently. Methodologists and statistical cognition researchers point to evidence from questionnaires and vignette studies to argue that researchers do not, in fact, grasp the core logic of SST. In one highly influential study of research psychologists, Oakes [3] presented six statements about a hypothetical significance test result to be categorized as true or false (e.g., “[The p value provides] the probability of the null hypothesis being true”). Despite all of these statements being false, 97% of the research psychologists categorized at least one as true. Oakes argues that this shows that the participants have an “[un]sound understanding of the logic of the significance test” (p. 82).

Oakes’ basic method and results have been replicated and extended with various groups, showing that students [4], instructors [5], and statisticians [6] all misinterpret SST results. Moreover, these misinterpretations are difficult to eliminate even through targeted interventions [7]. As a result, many have argued that use of SST should be discontinued or dramatically reduced, and may even contribute to wide-spread replication problems in the sciences [3,8–11].

The interpretation of studies of researchers’ understanding of SST are limited, however, by their methodology. A typical study presents a vignette describing research results. Statistical results are offered to the participants (e.g., a t statistic and p value), who are then asked to explicitly give or endorse various interpretations that are thought to represent participants’ understanding of SST. However, there are reasons to be cautious of drawing strong conclusions from these studies, including the abstract nature of such vignettes, the lack of investment researchers have in the fictional research, and their disconnection from research activity (e.g., experimentation and replication). It is unclear how well vignette studies (including ones by the present authors [12]) tap understanding of the core logic of SST or, say, familiarity with the technical terminology used to present statistical results. Conceptual understanding and fluency with common representations are both important, but are distinct.

A second major piece of evidence for misunderstandings of SST logic is reasoning errors in published papers. [13–16]. Like evidence from vignette studies, however, these errors are difficult to interpret as misunderstandings of SST logic *per se*. To avoid conflating problems with statistical reasoning with motivated scientific reasoning, one must test scientists’ statistical reasoning where these incentives are not in play.

Using a simulated experimental task, we tested researchers’ and trainees’ understanding and application of SST. Our key innovation allowing us to focus on SST was to prevent—or make difficult—the use of alternative strategies. If participants have poor understanding of SST, one would predict they would perform poorly and would not be sensitive to SST-related information. In fact, scientists performed well, explicitly reported using SST strategies, and were sensitive to SST-related information. Our results strongly suggest that common methods for assessing scientists’ competence may miss important aspects of their statistical knowledge, and hence that the case for abandoning significance testing may be overstated.

1. Testing reasoning by withholding information

In tests of perception—e.g., of colorblindness [17]—it is common to eliminate one cue (e.g., brightness) in order to assess the ability to use another (e.g., color). If color is the only useful cue for

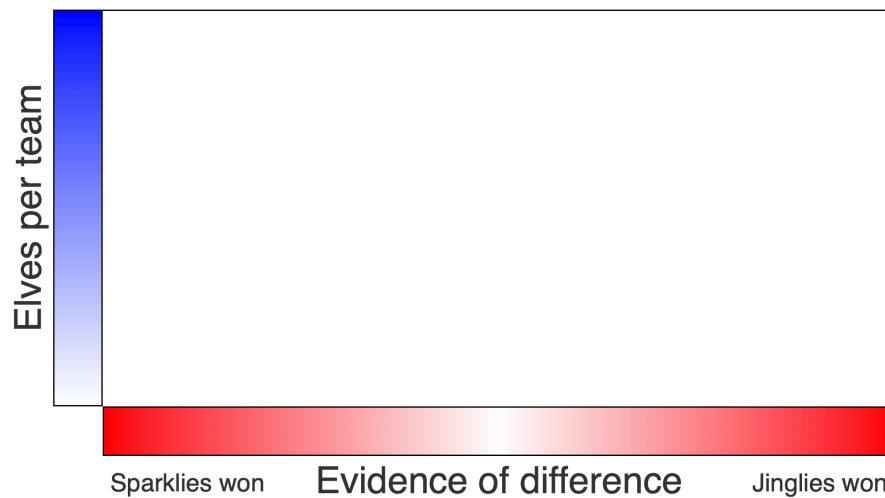


Figure 1: The interface on which samples were reported. The x -axis monotonically (but nonlinearly) related to the strength of the statistical evidence (z statistic) favoring one group; the y -axis is monotonically (but nonlinearly) related to the sample size. Underlying numerical values of the statistical evidence and sample sizes were unknown to the participant.

reading a number, deficits in color vision make the number difficult to read. We adopt a similar strategy to test statistical reasoning: we eliminate numerical information from statistical results to test scientists' ability to interpret results with reference to a null sampling distribution (i.e., SST logic). Without numerical information, many other strategies and heuristics, such as confidence intervals or Bayesian inference, are difficult or impossible to apply.¹

Participants were scientists or trainees recruited via social media. Our statistical reasoning task required them to perform a series of experiments to judge which of two groups of "Christmas elves" — "Jinglies" or "Sparklies" — could make more of a particular toy, where the true effect size was a randomly-assigned group difference between 0 and 1 standard deviations (δ). Because the study was run around the Christmas holiday season, it was hoped that the theme would make the task more enjoyable. The numerical information for an experiment, including sample size and the test statistic, was translated into color and location and displayed as a point on a two-dimensional visual interface (Figure 1). Participants could change the sample size per group for each experiment (increasing the time required to return a result), but did not know its numerical value. Importantly, the meaning of the colors and locations was unknown to the participants, aside from the monotone relationship with the sample size and statistical evidence.

Participants were also offered the opportunity to sample "random shuffle reports" that were described as the results of experiments with random assignment of elves to groups: that is, the result of experiments in which the null hypothesis was true. These results took no time to return. Participants were not told how to use these samples, only that they might use them.

Our experiment was constructed such that the only way to assess the evidence in the data was by comparison to a null sampling distribution: either the one provided by the random shuffle reports, or a simpler null that assumes that the evidence will favor one team or the other with 50% probability. Thus, the information afforded only the information in a p value, but it was not described as such, and participants had to *discover* for themselves how to use the information.

After sampling as many "experiments" and "random shuffle reports" as they liked, participants could report whether they believed Jinglies or Sparklies were better, that they could

¹ A formal statistical explanation showing that the task is difficult or impossible to perform using non-SST logic is given in Section 3 of [Supplement A](#).

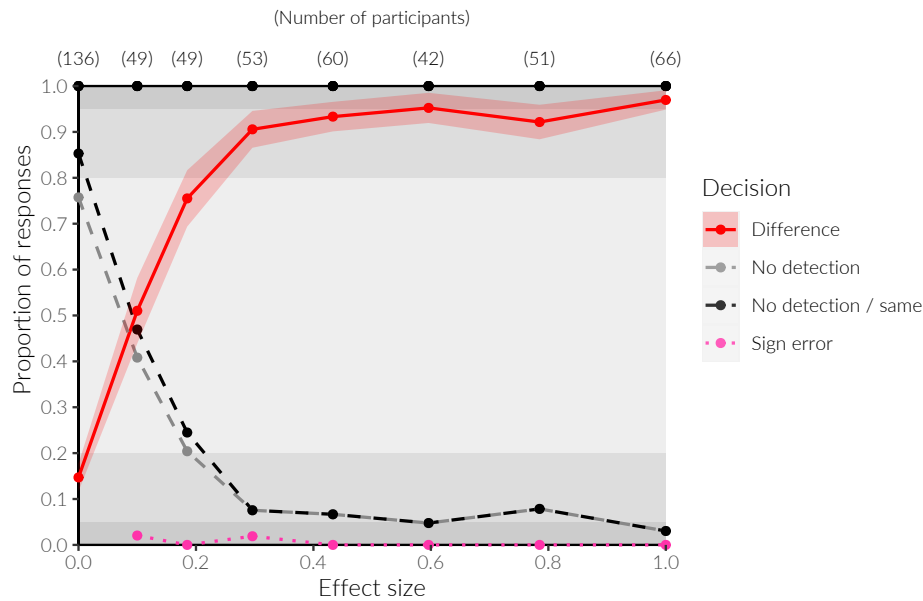


Figure 2: Correct decision and error rates as a function of true effect size. The red solid line represents the rate at which participants *correctly* determined the direction of the effect, except for at an effect size of 0, where it represents an incorrect claim of an effect. The shaded region around the solid line is one standard error wide. The difference between the black and the gray dashed lines shows the proportion of participants who decided that the two groups were the same (as opposed to failing to detect a difference).

not detect a difference, that there was no difference, or that they were bored and wanted to stop. Following their decision they were asked several open-ended questions about their strategy. Our key questions are whether participants can effectively find the “truth”, whether they report strategies consistent with SST, and whether their behaviour shows evidence of SST use.

Here, we report the results of 506 scientists or trainees who completed the statistical reasoning task. We provide an interactive app for exploring participants’ responses at <https://richarddmores.shinyapps.io/explore/>.

2. Success rates identifying effect sign

Decision error rates as a function of true effect size are shown in Figure 2.

Of the 136 participants for whom the null hypothesis was true (i.e. $\delta = 0$), 20 participants (15%) incorrectly indicated an effect. This is larger than the typically-accepted 5% false positive rate in many sciences; however, participants were performing a novel task with no recourse to numbers or statistical software. Those who did not indicate an effect when $\delta = 0$ tended to indicate that they *did not detect* an effect (103; 76%), which is the correct conclusion from the SST perspective. The other 13 (10%) indicated that the groups were the same, which under SST is a fallacy.

When there was a true effect ($|\delta| > 0$), correct decisions increased as a function of effect size, plateauing at about 95%. Of the 370 participants for whom $|\delta| > 0$, only 2 (0.5%) indicated the incorrect team [a sign, or Type S, error; 18]. For larger effect sizes, participants never incorrectly indicated that the two groups were the same.

Another way one can evaluate the participants’ responses is whether they reflect the information in the display at the time the decision is made, taking into account all points. To roughly quantify the evidence for a difference for each participant, we computed p values from the Wilcoxon rank-sum test between the shuffle reports and the fictitious experimental results as

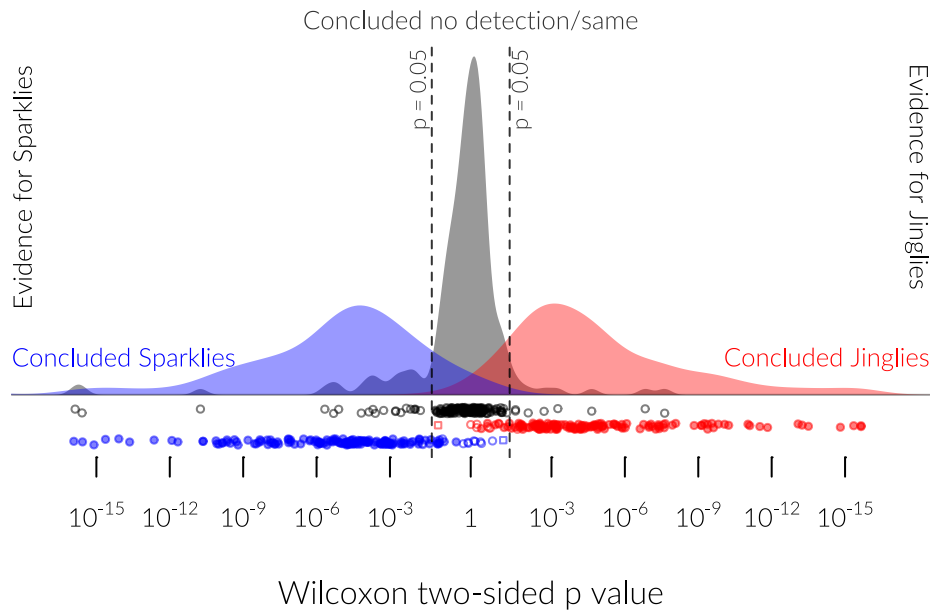


Figure 3: Statistical evidence underlying participants' decisions. The Wilcoxon rank-sum test p value (x axis) between the participants' null samples and experimental samples is used as a rough index of evidential strength. Kernel density estimates for the evidence are shown for three relevant conclusions. Each point at the bottom represents a single participant. Filled circles show correct decisions; hollow circles, incorrect decisions. The two hollow squares show sign errors.

Table 1: Frequencies of self-reported strategies.

	Strong	Only weak	Neither	Total	No shuffles	Missing
Count	362	69	75	506	28	29
%	71.54%	13.64%	14.82%	100%	5.53%	5.73%

they stood when the participant made their decision. Small Wilcoxon p values suggest a difference between the shuffle reports and the experiments.

Figure 3 shows the distribution of Wilcoxon p values (arranged by the direction of the decision). Kernel density estimates show the distributions of p values when participants indicated that Sparklies were faster, no detection/same, or that Jingles were faster. With few exceptions, participants' conclusions appear reasonable given the information in the display, though a few participants appear to ignore clear evidence of an effect.

3. Self-Reported SST Strategies

After they reported their decision regarding which team they believed was faster, we asked participants three questions about how they performed the task: what was the most salient information for their decision, what was their general strategy, and whether/how they used the shuffle reports.

We coded their responses according to whether they indicated comparing to the shuffle reports or using them to assess sampling variability (which we term "strong" significance testing strategies), assessing asymmetry in the display (a "weak" significance testing strategy, because it ignores information), and whether they explicitly deny using the shuffle reports.

As Table 1 shows that a large majority of participants (362, 71.54%) indicated using strong significance testing strategies. We should be cautious in directly interpreting this high number alone, however, because participants were told in the instructions that the shuffle reports could be used for assessing sampling variability. We did this to make clear what the shuffle reports were, but without explaining how to use them. To some extent, then, the text responses may reflect the instructions. However, the data strongly suggest a deeper understanding; first, among the responses were richer, lucid descriptions of SST logic, such as:

“[t]he random [shuffles] showed quite often such ‘strong evidence’, even at high sample sizes. That should not happen when the evidence is really strong, so probably the end of the scale was not [so] strong evidence... The random [shuffles] helped me to judge how common misleading evidence in that order of magnitude is, and after 5 samples from the real experiment I concluded that this result is probably not misleading evidence.”

Secondly—and most importantly—the instructions did not tell the participants *how* they should use the shuffles reports, yet many participants gave detailed accounts. As we show next, participants’ decisions were sensitive to a key manipulation of the shuffles reports in the manner that one would expect if they were using SST logic.

4. Sensitivity to SST-Relevant Information

In addition to a random effect size, participants were randomly assigned to one of two transformations of the location/color test statistic from an underlying z statistic. Of particular interest was how the transformation affected responding for the same visual deviation from the center. In one condition (“wide”), the null sampling distribution was visually about twice as wide as in the other (“narrow”).

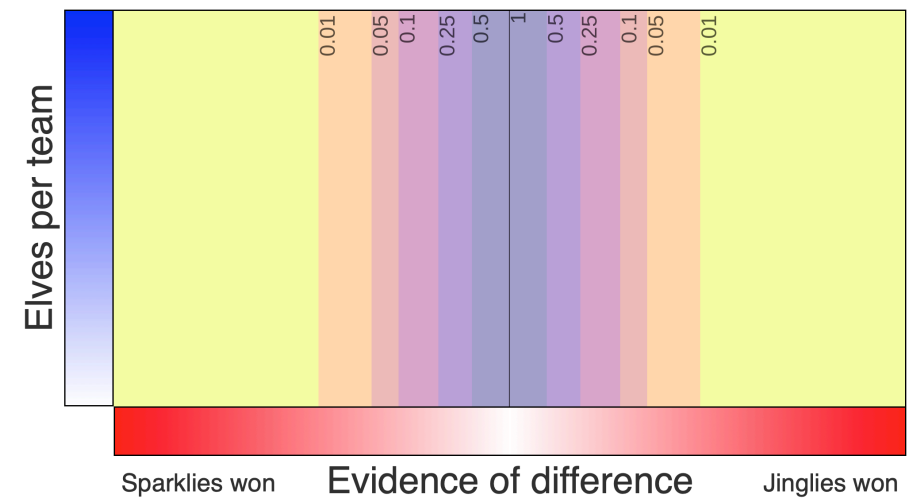
The effects of the manipulation are shown in Figure 4. The two experimental conditions used different arbitrary monotone mappings from the underlying Z -statistic to the visual space. Intuitively, this would be like deciding to use Z^3 instead of Z in all Z tests; one would need to adjust the significance criteria to account for the cubing (e.g., use $|1.96^3| = 7.53$ instead of $|1.96|$ for a $\alpha = 0.05$ level test), but the underlying test remains the same. The manipulation changes only the visual impression of the sampling distributions, allowing us to see how sensitive their responses are to the null sampling distribution as represented by the random shuffle reports.

If participants were using the shuffle reports to interpret the data, as would be predicted if they were using SST logic, the transformation should affect their interpretation of the visual evidence: a visually-extreme point should be more discounted against the sampling distribution that is wider. When we break down responses by the *visual* extremeness of the evidence, responses in two conditions should appear different; when we break down responses by *statistical* extremeness (i.e., the p value) responses in the two conditions should appear very similar, because the visual manipulation is irrelevant, given the p value.

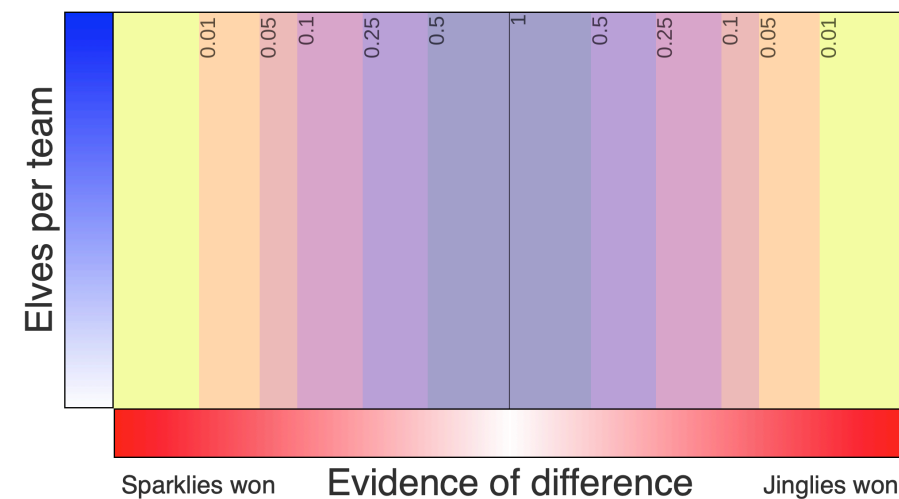
Figure 5 (top) shows responses (no detect/same or Jingles/Sparklies) as a function of the most extreme experiment sampled (x axis) and the transformation. There was a strong effect of the transformation consistent with use of the null sampling distribution; participants randomly assigned to the “narrow” evidence transformation responded “Jingles/Sparklies” for much less visually extreme evidence (sequential LRT: $\chi^2_2 = 35.492, p < .001$).

A logistic regression relating responses to the visual extremeness of the evidence and the transformation provides predicted probabilities of responding “Jingles/Sparklies” when the visual evidence corresponded to $p = 0.05$ for the null sampling distribution. In both the wide and the narrow conditions, the predicted probability of a “Jingles/Sparklies” response at the critical value was about 22%, despite that in the wide transformation condition this point was about twice as visually extreme.

Applying the same analysis to the responses corrected for their respective sampling distributions (Figure 5, bottom) almost completely eliminates the effect of experimental condition,



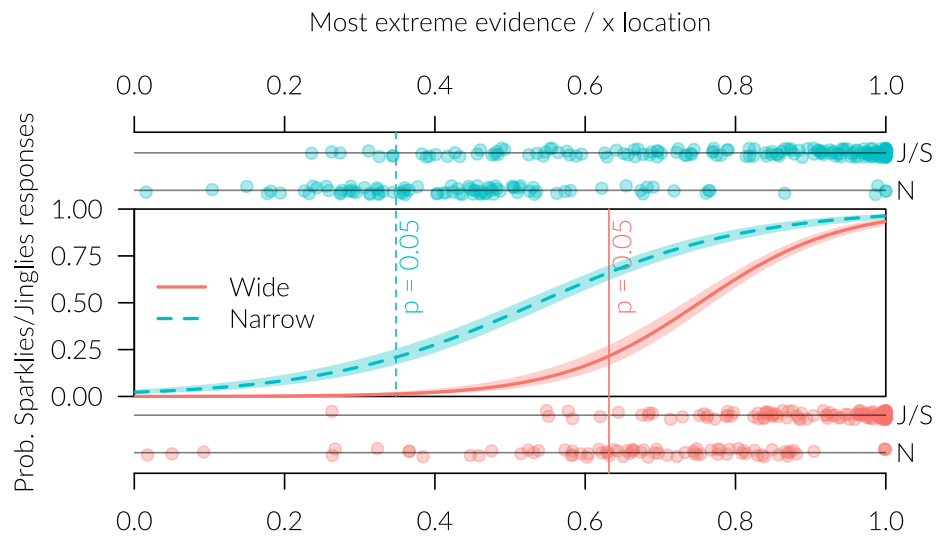
(A) Implied p values for the “narrow” condition.



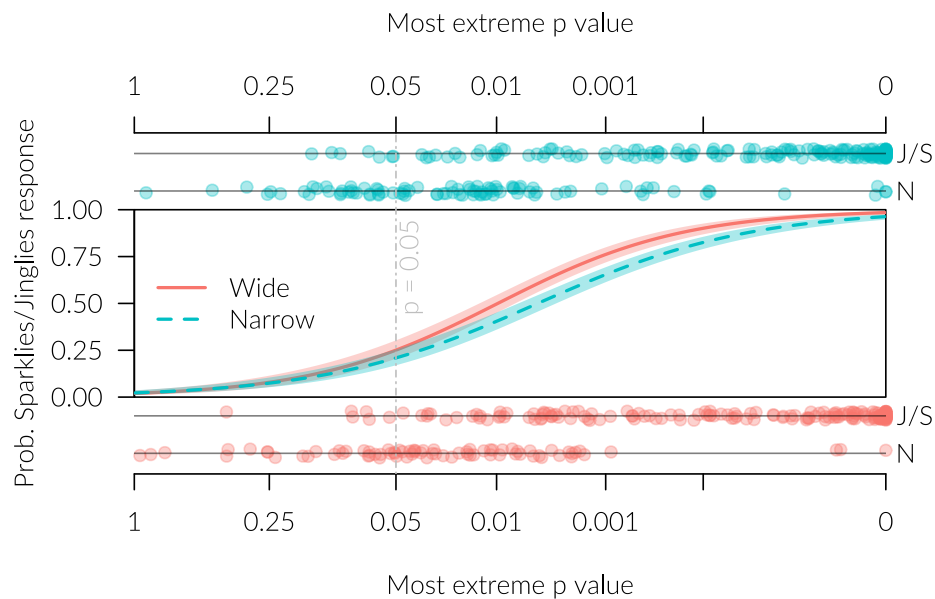
(B) Implied p values for the “wide” condition.

Figure 4: Selected two-tailed p values from the null sampling distributions of the narrow (top) and wide (bottom) conditions, projected onto the experimental interface. Participants were not shown the p values; they had to intuit them by using samples from the null distribution (“random shuffle reports”).

as would be expected if most participants were using the sampling distributions to calibrate (sequential LRT: $\chi^2_2 = 3.505, p = 0.173$). It is noteworthy that when they responses are aligned by sampling distribution, the wide condition appears to slightly dominate; this is consistent with a small number of participants incorrectly using the non-diagnostic visual extremeness to perform the task. If more people had been fooled by the irrelevant visual manipulation of the sampling distribution, we would expect this effect to be substantially larger.



(A) Predicted response probabilities relative to visual extremity. Vertical lines show the critical 0.05 for the corresponding null sampling distribution.



(B) Predicted response probabilities relative to the null sampling distributions (implicit p values).

Figure 5: The effect of the evidence transformation manipulation on responding. Points on top (narrow scale; $q = 7$) and bottom (wide scale; $q = 3$) represent participants' decisions as a function of the most extreme experiment sampled. See the methods details for the interpretation of q . "N" indicates a "no detect" or "same" response; "J/S" indicates a response in favor of a difference between the groups. Curves show predicted probability by a logistic regression fit with standard errors.

5. Discussion

Although it has previously been suggested that scientists have dramatic misunderstandings of SST logic, scientists and trainees in our experiment demonstrate both understanding and the ability to use it to come to the correct conclusion in a simulated statistical task. Moreover, they report strategies consistent with SST, and the signature of SST reasoning can be seen in their responses. Because we removed numerical effect size and sample size information — making strategies other than pure significance testing difficult or impossible to apply — our results are strong evidence that scientists *can* successfully deploy SST logic. It is still an open question what causes SST to be misunderstood so often, but it seems comprehension of its underlying logic can be ruled out for many scientists.

There are some major differences between our methods and the methods of previous demonstrations that scientists fail when reasoning about SST. We have already mentioned the fact that previous demonstrations typically use survey vignettes, as opposed to our more engaged experimental method. A second difference is our sample which is a substantially larger and more diverse convenience sample than previous demonstrations. It is possible that the recent so-called “replication crisis” [19] has raised awareness among scientists of deficits of statistical reasoning. If this increase in awareness accounts for our more optimistic findings, this would be good news for statistical educators.

Our findings echo other work showing that human reasoning can, under some conditions, be better than previously understood. [20,21]. Suggestions that SST be discontinued due to scientists’ apparent misunderstandings may have been hasty. Of course, there may be other reasons why some favor abandoning SST, but our work shows that given the opportunity, scientists successfully deploy basic SST logic. In spite of scientists’ real-life statistical behaviour often resembling a “ritual” [22], when we eliminate the ritual — no p value, or any other familiar number, was offered — they think statistically, very often arriving at the correct conclusion about the sign of the effect. Effective methodological reform in the sciences may be easier than abandoning significance testing, if we can harness scientists’ already-existing statistical competence.

6. Methods

(a) Participants

Participants were recruited via social media platforms such as Twitter and Facebook. All participants gave informed consent. Data inclusion criteria included sampling at least one shuffle report and experimental result, working in a scientific field, having at least some University education in science, and that it was their first time participating. Details are given in [Supplement B](#).

After applying all inclusion criteria, 506 participants remained for analysis.

(b) Experimental Design and Procedure

Each participant was randomly assigned to one of eight true effect sizes (from $\delta = 0$ to $\delta = 1$) and one of two evidence powers (“wide” $q = 3$ or “narrow” $q = 7$; see “Evidence Distributions” below). The probability of being assigned $\delta = 0$ was 25%, while the remaining effect sizes were equally probable at 11%. The probability of assignment to either evidence power was 50%. Details are given in Table 1.1 in [Supplement A](#).

After offering informed consent, participants read the cover story and instructions. During the instructions, the participant was introduced to the task through sampling random shuffle reports. After a brief recap of the instructions, participants performed the main task — sampling either random shuffles or experiments — until they made a decision about which, if either, elf

group was faster. They were then asked several open-ended questions about their strategy, some informational questions (results in [Supplement B](#)) and debriefed.

Qualtrics' duration estimate indicated that the median time spent on the experiment was 21 minutes.

(c) Evidence distributions

The evidence/horizontal (x) location test statistic presented to the participant was derived from a transformed Z statistic:

$$Z \sim \text{Normal}(\delta\sqrt{n/2}, 1)$$

where δ is a true effect size (randomly assigned to each participant, from 0 to 1) and n is the selected but unknown sample size (from 10 to 200 participants per group). Z then transformed to the $(-1,1)$ space:

$$x = \text{sgn}(Z) \left[1 - \left(1 - F_{\chi_1^2}(Z^2) \right)^{\frac{1}{q}} \right], \quad -1 \leq x \leq 1.$$

where $F_{\chi_1^2}$ is the cumulative distribution function of a χ_1^2 random variable, and $q \in \{3, 7\}$ was randomly assigned for each participant. $x = -1$ represented the left edge of the interface, $x = 0$ the middle, and $x = 1$ the right edge. The setting of q determined how spread out the test statistic was on the display. This arbitrary transformation was done to ensure that the test statistic's distribution was unfamiliar to the participant. See [Supplement A](#) for more details, including graphical depictions of the evidence distributions.

(d) Coding of open-ended strategy questions

We determined the coding scheme and independently categorized the first 20 participant, discussing the source of disagreements. After categorizing the remaining participants, some disagreements were resolved through mutual agreement, and a discussion between the authors was had over what caused the disagreements. The remainder of the disagreements were re-coded separately, and a final round of discussion resolved the remaining disagreements. The coding of participants' responses is described in detail in [Supplement B](#).

- **Funding:** This research was not supported by external funding.
- Compiled 2020-06-04 18:09:53 (Europe/London) under R version 4.0.0 (2020-04-24).

Ethics. This research project was evaluated by the Cardiff University School of Psychology (application number EC.18.12.11.5526G). It was found to be within the ethical guidelines for experiments with human participants. All participants gave informed consent prior to their participation.

Data Accessibility. Data and relevant code for this research work are stored in GitHub: https://github.com/richarddmorey/Morey_Hoekstra_StatCognition and have been archived within the Zenodo repository: <https://doi.org/10.5281/zenodo.3877106>

Authors' Contributions. RDM conceptualized and designed the study in consultation with RH. RDM analysed the data and created the materials and figures. The manuscript was written by RDM and RH.

Competing Interests. The authors declare no conflicts of interest.

References

1. Dempster AP. 1964 On the Difficulties Inherent in Fisher's Fiducial Argument. *Journal of the American Statistical Association* **59**, 56–66.
2. Greenland S. 2019 Valid P -values behave exactly as they should: some misleading criticisms of P -values and their resolution with S -values. *The American Statistician* **73**, 106–114. Publisher: Taylor & Francis.

- 259 3. Oakes M. 1986 *Statistical inference: A commentary for the social and behavioral sciences*. Chichester:
260 Wiley.
- 261 4. Falk R, Greenbaum CW. 1995 Significance Tests Die Hard: The Amazing Persistence of a
262 Probabilistic Misconception. *Theory & Psychology* **5**, 75–98.
- 263 5. Haller H, Krauss S. 2002 Misinterpretations of Significance: A Problem Students Share with
264 Their Teachers?. *Methods of Psychological Research Online* **7**.
- 265 6. Lecoutre MP, Poitevineau J, Lecoutre B. 2003 Even statisticians are not immune to
266 misinterpretations of Null Hypothesis Tests. *International Journal of Psychology* **38**, 37–45.
- 267 7. Kalinowski P, Fidler F, Cumming G. 2008 Overcoming the Inverse Probability Fallacy: A
268 Comparison of Two Teaching Interventions. *Methodology* **4**, 152–158.
- 269 8. Carver R. 1978 The Case Against Statistical Significance Testing. *Harvard Educational Review*
270 **48**, 378–399.
- 271 9. Fidler F. 2006 Should Psychology abandon *p* values and teach CIs instead? Evidence-based
272 reforms in statistics education. In *Proceedings of the 7th International Conference on Teaching*
273 *Statistics*.
- 274 10. The B. 2011 Significance testing - are we ready yet to abandon its use?. *Current Medical Research*
275 *and Opinion* **27**, 2087–2090. PMID: 21916530.
- 276 11. Wasserstein RL, Lazar NA. 2016 The ASA's Statement on p-Values: Context, Process, and
277 Purpose. *The American Statistician* **70**, 129–133.
- 278 12. Hoekstra R, Morey RD, Rouder JN, Wagenmakers EJ. 2014 Robust Misinterpretation of
279 Confidence Intervals. *Psychonomic Bulletin & Review* **21**, 1157–1164.
- 280 13. Gelman A, Stern H. 2006 The difference between “significant” and “not significant” is not
281 itself statistically significant. *The American Statistician* **60**, 328–331.
- 282 14. Hoekstra R, Finch S, Kiers HAL, Johnson A. 2006 Probability as certainty: Dichotomous
283 thinking and the misuse of *p* values. *Psychonomic Bulletin & Review* **13**, 1033–1037.
- 284 15. Nieuwenhuis S, Forstmann BU, Wagenmakers EJ. 2011 Erroneous analyses of interactions in
285 neuroscience: A problem of significance. *Nature Neuroscience* **14**, 1105–1107.
- 286 16. Weisburd D, Lum CM, Yang SM. 2003 When can we Conclude that Treatments or Programs
287 “Don’t work”?. *The Annals of the American Academy of Political and Social Science* **587**, 31–48.
- 288 17. Ishihara S. 1972 *Tests for colour-blindness*. Tokyo: Kanehara Shuppan Co. Ltd 24 plate edition
289 edition.
- 290 18. Gelman A, Carlin J. 2014 Beyond Power Calculations: Assessing Type S (Sign) and Type M
291 (Magnitude) Errors. *Perspectives on Psychological Science* **9**, 641–651. PMID: 26186114.
- 292 19. Baker M. 2016 1,500 scientists lift the lid on reproducibility. *Nature News* **533**, 452. Section:
293 News Feature.
- 294 20. Cosmides L, Tooby J. 1992 Cognitive Adaptations for Social Exchange. In *The Adapted Mind:*
295 *Evolutionary psychology and the generation of culture*, pp. 163–228. New York: Oxford University
296 Press.
- 297 21. Gigerenzer G, Hoffrage U. 1995 How to improve Bayesian reasoning without instruction:
298 frequency formats. *Psychological Review* **102**, 684–704.
- 299 22. Gigerenzer G, Krauss S, Vitouch O. 2004 The null ritual: What you always wanted to know
300 about significance testing but were afraid to ask. In Kaplan D, editor, *The Sage handbook of*
301 *quantitative methodology for the social sciences*, . Thousand Oaks, CA: Sage.