# Moving to a world beyond p < 0.05: where we came from and the road map

Alice Richardson

National Centre for Epidemiology & Population Health
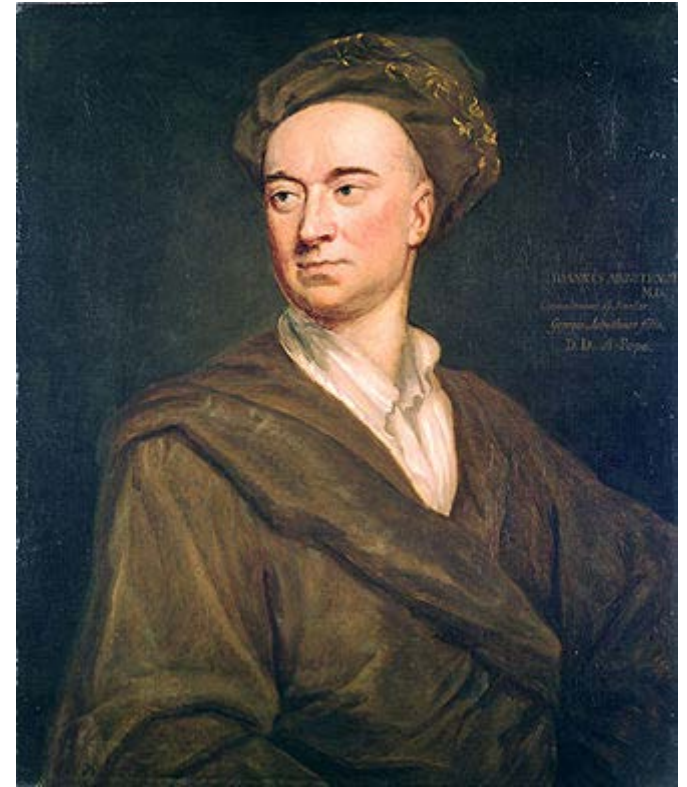
20 June 2019

# What is a p value?

- The probability
- Of seeing a test statistic
- At least as extreme as the one you got
- Given the null hypothesis is true

# History

- Arbuthnot (1710) "An argument for divine providence taken from the constant regularity observed in the births of both sexes"
- Calculated a p value form a binomial distribution

# History

- Laplace (1827) measuring the effect of the moon on atmospheric tides
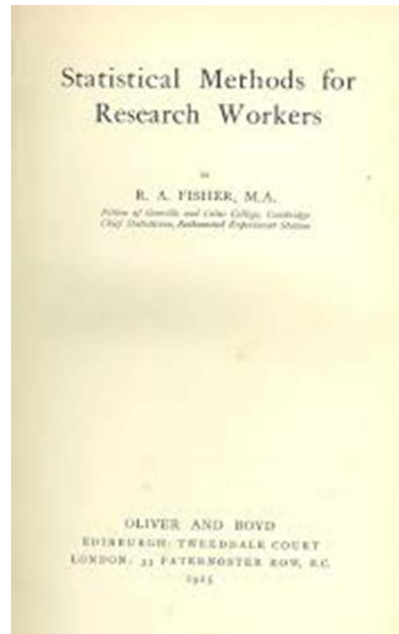
- Calculated a p value from a Normal distribution

# History

- Poisson (1837) modelling jury verdicts (7/12 was sufficient to convict)

- Calculated confidence intervals from Normal distributions

# History

- Fisher (1922) introduced the term "p-value"

Statistical Methods for Research Workers

R. A. FISHER, M.A.

*Fellow of Gonville and Caius College, Cambridge
Chief Statistician, Rothamsted Experimental Station*

OLIVER AND BOYD
EDINBURGH: TWEEDDALE COURT
LONDON: 33 PATERNOSTER ROW, E.C.
1925

- Describing the standard normal distribution: "The value for which **P=0.05**, or 1 in 20, is 1.96 or nearly 2; it is convenient to take this point as a limit in judging whether a deviation ought to be considered **significant** or not. Deviations exceeding twice the standard deviation are thus formally regarded as **significant**. Using this criterion we should be led to follow up a false indication only once in 22 trials, even if the statistics were the only guide available."

- "If one in twenty does not seem high enough odds, we may, if we prefer it, draw the line at one in fifty … or one in a hundred …. Personally, the writer prefers to set a low standard of significance at the 5 per cent point, and ignore entirely all results which fail to reach this level. A scientific fact should be regarded as experimentally established only if a properly designed experiment rarely fails to give this level of significance."

# History

- Neyman & Pearson (1933) introduced the concept of hypothesis tests
- Neyman (1937) introduced the term "confidence interval"
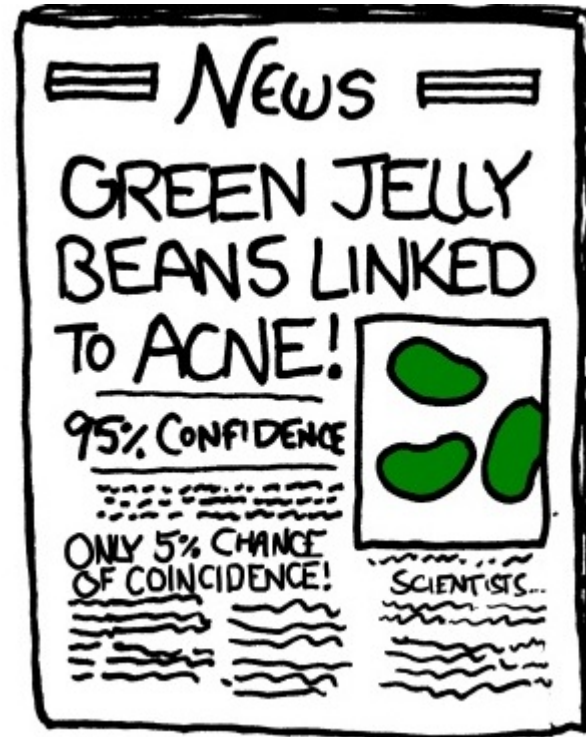
# Poking fun at p values

- Jokes about the threshold

- teetering on the brink of significance (p=0.06)
- not significant in the narrow sense of the word (p=0.29)
- partial significance (p>0.09)
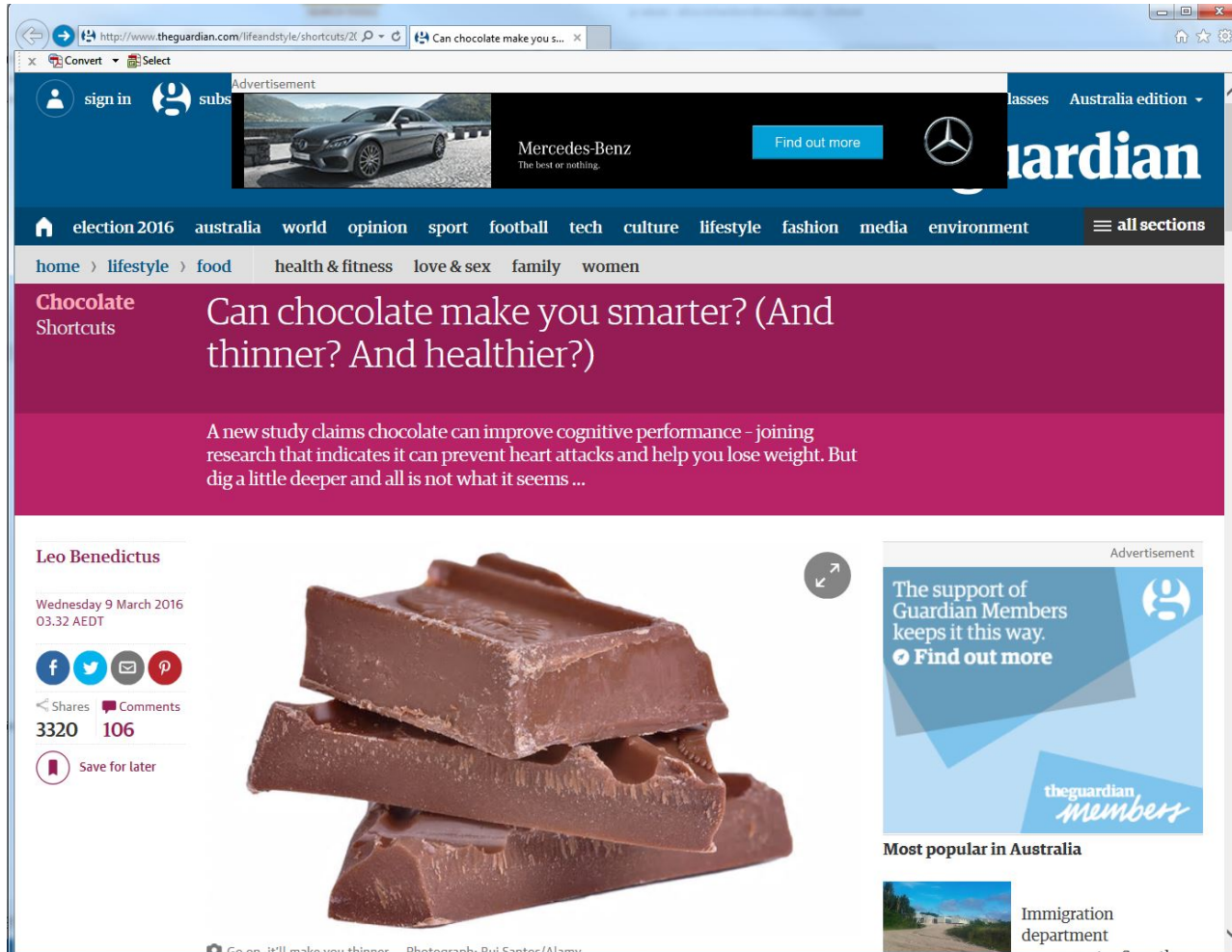- significant at the .07 level

• Jokes about multiple comparisons

- Jokes about the terminology

# But statisticians can poke fun too …

# • BASP bans the p-value

Routledge
Taylor & Francis Group

## Editorial

David Trafimow and Michael Marks
*New Mexico State University*

The *Basic and Applied Social Psychology* (BASP) 2014 Editorial emphasized that the null hypothesis significance testing procedure (NHSTP) is invalid, and thus authors would be not required to perform it (Trafimow, 2014). However, to allow authors a grace period, the Editorial stopped short of actually banning the NHSTP. The purpose of the present Editorial is to announce that the grace period is over. From now on, BASP is banning the NHSTP.

With the banning of the NHSTP from BASP, what are the implications for authors? The following are anticipated questions and their corresponding answers.

**Question 1.** *Will manuscripts with p-values be desk rejected automatically?*

**Answer to Question 1.** No. If manuscripts pass the preliminary inspection, they will be sent out for review. But prior to publication, authors will have to remove all vestiges of the NHSTP (*p*-values, *t*-values, *F*-values, statements about "significant" differences or lack thereof, and so on).

**Question 2.** *What about other types of inferential statistics such as confidence intervals or Bayesian methods?*

**Answer to Question 2.** Confidence intervals suffer from an inverse inference problem that is not very different from that suffered by the NHSTP. In the NHSTP, the problem is in traversing the distance from the probability of the finding, given the null hypothesis, to the probability of the null hypothesis, given the finding. Regarding confidence intervals, the problem is that, for example, a 95% confidence interval does not indicate that the parameter of interest has a 95% probability of being within the interval. Rather, it means merely that if an infinite number of samples were taken and confidence intervals computed, 95% of the confidence intervals would capture the population parameter. Analogous to how the NHSTP fails to provide the probability of the null hypothesis, which is needed to provide

a strong case for rejecting it, confidence intervals do not provide a strong case for concluding that the population parameter of interest is likely to be within the stated interval. Therefore, confidence intervals also are banned from BASP.

Bayesian procedures are more interesting. The usual problem with Bayesian procedures is that they depend on some sort of Laplacian assumption to generate numbers where none exist. The Laplacian assumption is that when in a state of ignorance, the researcher should assign an equal probability to each possibility. The problems are well documented (Chihara, 1994; Fisher, 1973; Glymour, 1980; Popper, 1983; Suppes, 1994; Trafimow, 2003, 2005, 2006). However, there have been Bayesian proposals that at least somewhat circumvent the Laplacian assumption, and there might even be cases where there are strong grounds for assuming that the numbers really are there (see Fisher, 1973, for an example). Consequently, with respect to Bayesian procedures, we reserve the right to make case-by-case judgments, and thus Bayesian procedures are neither required nor banned from BASP.

**Question 3.** *Are any inferential statistical procedures required?*

**Answer to Question 3.** No, because the state of the art remains uncertain. However, BASP will require strong descriptive statistics, including effect sizes. We also encourage the presentation of frequency or distributional data when this is feasible. Finally, we encourage the use of larger sample sizes than is typical in much psychology research, because as the sample size increases, descriptive statistics become increasingly stable and sampling error is less of a problem. However, we will stop short of requiring particular sample sizes, because it is possible to imagine circumstances where more typical sample sizes might be justifiable.

We conclude with one last thought. Some might view the NHSTP ban as indicating that it will be easier to publish in BASP, or that less rigorous manuscripts will be acceptable. This is not so. On the contrary, we believe

Correspondence should be sent to David Trafimow, Department of Psychology, MSC 3452, New Mexico State University, P.O. Box 30001, Las Cruces, NM 88003-8001. E-mail: dtrafimo@nmsu.edu

# ASA statement of six principles

# • ASA editorial introducing 40+ papers
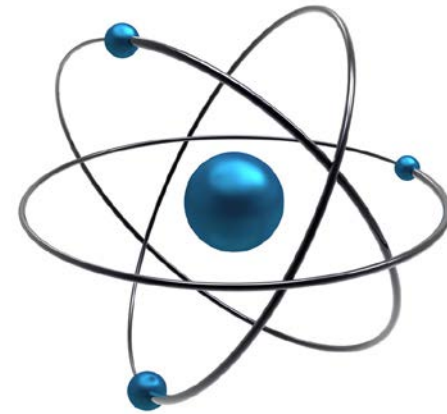
**EDITORIAL**

OPEN ACCESS     Check for updates

## Moving to a World Beyond "$p < 0.05$"

Some of you exploring this special issue of *The American Statistician* might be wondering if it's a scolding from pedantic statisticians lecturing you about what *not* to do with $p$-values, without offering any real ideas of what *to do* about the very hard problem of separating signal from noise in data and making decisions under uncertainty. Fear not. In this issue, thanks to 43 innovative and thought-provoking papers from forward-looking statisticians, help is on the way.

special issue of *The American Statistician*. Authors were explicitly instructed to develop papers for the variety of audiences interested in these topics. If you use statistics in research, business, or policymaking but are not a statistician, these articles were indeed written with YOU in mind. And if you are a statistician, there is still much here for you as well.

The papers in this issue propose many new ideas, ideas that in our determination as editors merited publication to enable broader consideration and debate. The ideas in this editorial are

- ATOM(IC)
- Accept uncertainty; be
- Thoughtful
- Open, and
- Modest

- Institutional Change

# Resources

- https://www.howresearchers.com/episodes/episode-2/

- https://i2insights.org/2019/04/30/replacing-p-values/