ON THE USE AND ABUSE OF P-VALUES

Giorgio Bedogni

www.giorgiobedogni.it

http://www.giorgiobedogni.it/tutorials/metodologia/pmetodologia.html



The mechanical, ritualistic application of statistics is contributing to a crisis in science. Education, software and peer review have encouraged poor practice – and it is time for statisticians to fight back. By **Philip B. Stark** and **Andrea Saltelli**

40 | SIGNIFICANCE | August 2018

«The problem is one of cargo-cult statistics – the ritual miming of statistics rather than conscientious practice»

CARGO-CULT SCIENCE



Cargo Cult Science

by RICHARD P. FEYNMAN

Some remarks on science, pseudoscience, and learning how to not fool yourself. Caltech's 1974 commencement address.







http://calteches.library.caltech.edu/51/2/CargoCult.htm

«In our experience many applications of statistics are cargo-cult statistics: practitioners go through the notions of fitting models, computing p-values or confidence intervals»

«They invoke statistical terms and procedures <u>as incantations</u>, with scant understanding of the assumptions or relevance of the calculations, or even the meaning of the terminology»

«Statistical software enables and promotes cargo-cult statistics»

 «The <u>misuse of p-values</u>, hypothesis tests, and confidence intervals might be deemed frequentist cargo-cult statistics. There is also Bayesian cargo-cult statistics»

On the use and abuse of p-values



The ASA's Statement on p-Values: Context, Process, and Purpose

Ronald L. Wasserstein & Nicole A. Lazar

To cite this article: Ronald L. Wasserstein & Nicole A. Lazar (2016) The ASA's Statement on p-Values: Context, Process, and Purpose, The American Statistician, 70:2, 129-133, DOI: <u>10.1080/00031305.2016.1154108</u>

To link to this article: http://dx.doi.org/10.1080/00031305.2016.1154108

ASA position statement

- "The ASA Board was also stimulated by <u>highly visible discussions</u> over the last few years"
- "The ASA has <u>not</u> previously taken positions on specific matters of statistical practice"

ASA position statement

 "That turned out to be relatively easy to do [to find points of agreement], but <u>it was just as easy</u> to find points of intense disagreement"

ASA position statement – What is a p-value?

 "Informally, a p-value is the probability <u>under a specified statistical</u> <u>model</u> that a statistical summary of the data (e.g., the sample mean difference between two compared groups) would be equal to or more extreme than its observed value"

 "P-values can indicate how <u>incompatible</u> the data are with a specified statistical model. A p-value provides one approach to summarizing the incompatibility between a particular set of data and a <u>proposed</u> <u>model</u> for the data"

- "The most common context is a model, constructed under a set of assumptions, together with a so-called <u>null hypothesis</u>"
- "Often the null hypothesis postulates the absence of an effect, such as no difference between two groups, or the absence of a relationship between a factor and an outcome"

- "The smaller the p-value, the greater the statistical incompatibility of the data with the null hypothesis, <u>if the underlying assumptions</u> used to calculate the p-value <u>hold</u>"
- "This incompatibility can be interpreted as casting doubt on or providing evidence against the null hypothesis <u>or the underlying</u> <u>assumptions</u>"

- "P-values do <u>not</u> measure the probability that the studied hypothesis is true, or the probability that the data were produced by random chance alone"
- "Researchers often wish to turn a p-value into a statement about the truth of a null hypothesis, or about the probability that random chance produced the observed data"
- "The p-value is neither. It is a statement about data in relation to a specified hypothetical explanation, and is not a statement about the explanation itself"

- "Scientific conclusions and business or policy decisions should <u>not</u> be based only on whether a p-value passes a specific threshold"
- "Practices that reduce data analysis or scientific inference to mechanical 'bright-line' rules (such as 'p<0.05') for justifying scientific claims or conclusions can lead to erroneous beliefs and poor decision making"

- "A conclusion does not immediately become "true" on one side of the divide and "false" on the other"
- "Researchers should bring <u>many contextual factors</u> into play to derive scientific inferences, including the design of a study, the quality of the measurements, the external evidence for the phenomenon under study, and the validity of assumptions that underlie the data analysis"

- "Pragmatic considerations often require binary, "yes-no" decisions, but this does <u>not</u> mean that p-values alone can ensure that a decision is correct or incorrect"
- "The widespread use of "statistical significance" (generally interpreted as 'p<0.05') as a license for making a claim of a scientific finding (or implied truth) leads to considerable <u>distortion</u> of the scientific process"

- "Proper inference requires full reporting and transparency. P-values and related analyses should <u>not</u> be reported selectively"
- "Conducting multiple analyses of the data and reporting only those with certain p-values (typically those passing a significance threshold) renders the reported p-values <u>essentially uninterpretable</u>"

 "Cherry-picking promising findings, also known by such terms as data dredging, significance chasing, significance questing, selective inference, and 'p-hacking', leads to a <u>spurious excess</u> of statistically significant results in the published literature and should be vigorously avoided"

- "One need not formally carry out multiple statistical tests for this problem to arise"
- "<u>Whenever a researcher chooses</u> what to present based on statistical results, valid interpretation of those results is severely compromised if the reader is not informed of the choice and its basis"
- "<u>Researchers should disclose</u> the number of hypotheses explored during the study, all data collection decisions, all statistical analyses conducted, and all p-values computed"

 "Valid scientific conclusions based on p-values and related statistics <u>cannot be drawn</u> without at least knowing how many and which analyses were conducted, and how those analyses (including pvalues) were selected for reporting"

- "A p-value, or statistical significance, does <u>not</u> measure the size of an effect or the importance of a result"
- "Statistical significance is <u>not</u> equivalent to scientific, human, or economic significance"
- "Smaller p-values do <u>not</u> necessarily imply the presence of larger or more important effects, and larger p-values do not imply a lack of importance or even lack of effect"

- "Any effect, no matter how tiny, can produce a small p-value if the sample size or measurement precision is high enough, and large effects may produce unimpressive p-values if the sample size is small or measurements are imprecise"
- "Similarly, identical estimated effects will have different p-values if the precision of the estimates differs"

- "By itself, a p-value does <u>not</u> provide a good measure of evidence regarding a model or hypothesis"
- "Researchers should recognize that a p-value without context or other evidence provides limited information"

- "For example, a p-value near 0.05 taken by itself offers only weak evidence against the null hypothesis"
- "Likewise, a relatively large p-value does not imply evidence in favor of the null hypothesis; <u>many other hypotheses may be equally or</u> <u>more consistent with the observed data</u>"
- "For these reasons, data analysis should not end with the calculation of a p-value when other approaches are appropriate and feasible"

ASA position statement – other methods

• "In view of the prevalent misuses of and misconceptions concerning p-values, some statisticians prefer to supplement or even replace p-values with other approaches"

ASA position statement – other methods

 "These include methods that emphasize estimation over testing, such as confidence, credibility, or prediction intervals; Bayesian methods; alternative measures of evidence, such as likelihood ratios or Bayes Factors; and other approaches such as decision-theoretic modeling and false discovery rates"

ASA position statement – other methods

• "All these measures and approaches <u>rely on further assumptions</u>, but they may more directly address the size of an effect (and its associated uncertainty) or whether the hypothesis is correct"

ASA position statement – conclusion

• No single index should be a substitute for scientific reasoning.

You are invited to read this wonderful paper...

Eur J Epidemiol (2016) 31:337–350 DOI 10.1007/s10654-016-0149-3





Statistical tests, *P* values, confidence intervals, and power: a guide to misinterpretations

Sander Greenland¹ · Stephen J. Senn² · Kenneth J. Rothman³ · John B. Carlin⁴ · Charles Poole⁵ · Steven N. Goodman⁶ · Douglas G. Altman⁷

Caveat

 "Analysis interpretation depends on <u>contextual judgements</u> about how reality is mapped onto the model and how the formal analysis results are mapped back to reality"

Greenland S. Eur J Epidemiol. 2017;32:3-20.

ASA position statement - 1 year after

The ASA's *p*-value statement, one year on

Its aim was to stop the misuse of statistical significance testing. But **Robert Matthews** argues that little has changed in the 12 months since the ASA's intervention

38 | SIGNIFICANCE | April 2017

ASA position statement - 1 year after -Robert Matthews

- "Yet a year on, it is not clear that the ASA's statement has had any substantive effect at all"
- "The commentaries accompanying the statement show that even at the time some participants feared there would be no lasting impact"
- "Goodman's frustration with this resonates with mine: <u>exactly how</u> [are] scientists supposed to do that?"

Matthews R et al. Significance. 2017;14:38.

ASA position statement – 1 year after – – Ron Wasserstein*

- "The statement did not go as far as it should go, but it went as far as it could go"
- "To address the statement's shortcomings the ASA is convening a Symposium on Statistical Inference this October"

*Executive Director, American Statistical Association

Matthews R et al. Significance. 2017;14:38.

ASA position statement – 1 year after – – David Spigelhalter*

- "If by this he [Matthews] means [to ban] all use of p-values than I must disagree"
- "Fortunately the Neyman-Pearson's decision-theoretic idea of "accepting" the null has just been consigned to the overwflowing dustbin of inappropriate scientific ideas"

Matthews R et al. Significance. 2017;14:38.

*President, Royal Statistical Association

ASA position statement – 1 year after – – David Spigelhalter

 "While it would be wonderful if every analysis would be informed by someone skilled in statistical methodology, <u>whether a nominal</u> <u>statistician or not</u>, the rise of data science means that even more practitioners will be without a full professional training and <u>continue</u> <u>to do their *t*-ing and *p*-ing</u>"

Matthews R et al. Significance. 2017;14:38.

ASA position statement – 1 year after



ASA position statement – 2 years after?



In gentle praise of significance tests – Sir David Cox





«Skin in the game keeps human hubris in check»

Taleb N. Skin in the game. Penguin; 2017.

«Academia has a tendency, when unchecked (from lack of skin in the game), to evolve into a ritualistic self-referential publishing game»

Taleb N. Skin in the game. Penguin; 2017.

«Indeed, journals may reject papers that use more reliable or more rigorous methods the discipline is accustomed to, simply because the methods are unfamiliar»

 «I will not be proselitizing for a given statistical school, so you can relax. Frankly, <u>they all have shortcomings</u>, insofar as one can even glean a clear statement of a given statistical "school"».

Mayo D. Statistical inference as severe testing: how to get beyond the statistics wars. Cambridge University Press; 2018.

How do I (= Giorgio) put my skin the game?

 «If they've made a mistake, correct them gently and show them where they went wrong. If you can't do that, then the blame lies with you. Or no one»

Marcus Aurelius, Ad se ipsum (X 4)

http://www.perseus.tufts.edu/hopper/text?doc=urn:cts:greekLit:tlg0562.t lg001.perseus-grc1:10.4.1

How do I (= Giorgio) put my skin the game?

«But if on my part I fail to produce yourself as my one witness to confirm what I say, I consider I have achieved nothing of any account»

(Socrates in) Plato, Gorgias

http://data.perseus.org/citations/urn:cts:greekLit:tlg0059.tlg023.perseusgrc1:472b

Thank you for your attention!

I invite you to join Qeios:

www.qeios.com/invitation-to-join/researcher/314

