Log in     Register

The online home for the publications of the American Statistical Association

✓ Free access

◀) **Listen** | ▶

Editorial

# Editorial

Ronald L. Wasserstein ✉ & Nicole A. Lazar

❝ Download citation     ⬈ https://doi.org/10.1080/00031305.2016.1154108     ◉ Check for updates

📄 **Full Article**     🖻 **Figures & data**     📕 **References**     ➕ **Supplemental**     ❝ **Citations**

📊 **Metrics**     🖶 **Reprints & Permissions**     📄 **PDF**     ➕ AddThis

In this article     ☰

Taylor & Francis
Taylor & Francis Group

EDITORIAL

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

In February 2014, George Cobb, Professor Emeritus of Mathematics and Statistics at Mount Holyoke College, posed these questions to an ASA discussion forum:

Q: Why do so many colleges and grad schools teach $p = 0.05$?
A: Because that's still what the scientific community and journal editors use.
Q: Why do so many people still use $p = 0.05$?
A: Because that's what they were taught in college or grad school.

Cobb's concern was a long-worrisome circularity in the sociology of science based on the use of bright lines such as $p < 0.05$: "We teach it because it's what we do; we do it because it's what we teach." This concern was brought to the attention of the ASA Board.

The ASA Board was also stimulated by highly visible discussions over the last few years. For example, ScienceNews (Siegfried 2010) wrote: "It's science's dirtiest secret: The 'scientific method' of testing hypotheses by statistical analysis stands on a flimsy foundation." A November 2013, article in Phys.org Science News Wire (2013) cited "numerous deep flaws" in null hypothesis significance testing. A ScienceNews article (Siegfried 2014) on February 7, 2014, said "statistical techniques for testing hypotheses ... have more flaws than Facebook's privacy policies." A week later, statistician and "Simply Statistics" blogger Jeff Leek responded. "The problem is not that people use P-values poorly," Leek wrote, "it is that the vast majority of data analysis is not performed by people properly trained to perform data analysis" (Leek 2014). That same week, statistician and science writer Regina Nuzzo published an article in Nature entitled "Scientific Method: Statistical Errors" (Nuzzo 2014). That article is now one of the most highly viewed Nature articles, as reported by altmetric.com (http://www.altmetric.com/details/2115792#score).

Of course, it was not simply a matter of responding to some articles in print. The statistical community has been deeply concerned about issues of reproducibility and replicability of scientific conclusions. Without getting into definitions and distinctions of these terms, we observe that much confusion and even doubt about the validity of science is arising. Such doubt can lead to radical choices, such as the one taken by the editors of Basic and Applied Social Psychology, who decided to ban p-values (null hypothesis significance testing) (Trafimow and Marks 2015). Misunderstanding or misuse of statistical inference is only one cause of the "reproducibility crisis" (Peng 2015), but to our community, it is an important one.

When the ASA Board decided to take up the challenge of developing a policy statement on p-values and statistical significance, it did so recognizing this was not a lightly taken step. The ASA has not previously taken positions on specific matters of statistical practice. The closest the association has come to this is a statement on the use of value-added models (VAM) for educational assessment (Morganstein and Wasserstein 2014) and a statement on risk-limiting post-election audits (American Statistical Association 2010). However, these were truly policy-related statements. The VAM statement addressed a key educational policy issue, acknowledging the complexity of the issues involved, citing limitations of VAMs as effective performance models, and urging that they be developed and interpreted with the involvement of statisticians. The statement on election auditing was also in response to a major but specific policy issue (close elections in 2008), and said that statistically based election audits should become a routine part of election processes.

By contrast, the Board envisioned that the ASA statement on p-values and statistical significance would shed light on an aspect of our field that is too often misunderstood and misused in the broader research community, and, in the process, provides the community a service. The intended audience would be researchers, practitioners, and science writers who are not primarily statisticians. Thus, this statement would be quite different from anything previously attempted.

The Board tasked Wasserstein with assembling a group of experts representing a wide variety of points of view. On behalf of the Board, he reached out to more than two dozen such people, all of whom said they would be happy to be involved. Several expressed doubt about whether agreement could be reached, but those who did said, in effect, that if there was going to be a discussion, they wanted to be involved.

Over the course of many months, group members discussed what format the statement should take, tried to more concretely visualize the audience for the statement, and began to find points of agreement. That turned out to be relatively easy to do, but it was just as easy to find points of intense disagreement.

The time came for the group to sit down together to hash out these points, and so in October 2015, 20 members of the group met at the ASA Office in Alexandria, Virginia. The 2-day meeting was facilitated by Regina Nuzzo, and by the end of the meeting, a good set of points around which the statement could be built was developed.

The next 3 months saw multiple drafts of the statement, reviewed by group members, by Board members (in a lengthy discussion at the November 2015 ASA Board meeting), and by members of the target audience. Finally, on January 29, 2016, the Executive Committee of the ASA approved the statement.

The statement development process was lengthier and more controversial than anticipated. For example, there was considerable discussion about how best to address the issue of multiple potential comparisons (Gelman and Loken 2014). We debated at some length the issues behind the words "a p-value near 0.05 taken by itself offers only weak evidence against the null

---

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

Ronald L. Wasserstein[a]* & Nicole A. Lazar[a]

In this article ≡

In February 2014, George Cobb, Professor Emeritus of Mathematics and Statistics at Mount Holyoke College, posed these questions to an ASA discussion forum:

> Q: Why do so many colleges and grad schools teach $p$ = 0.05?
>
> A: Because that's still what the scientific community and journal editors use.
>
> Q: Why do so many people still use $p$ = 0.05?
>
> A: Because that's what they were taught in college or grad school.

Cobb's concern was a long-worrisome circularity in the sociology of science based on the use of bright lines such as $p < 0.05$: "We teach it because it's what we do; we do it because it's what we teach." This concern was brought to the attention of the ASA Board.

The ASA Board was also stimulated by highly visible discussions over the last few years. For example, ScienceNews (Siegfried ) wrote: "It's science's dirtiest secret: The 'scientific method' of testing hypotheses by statistical analysis stands on a flimsy foundation." A November 2013, article in Phys.org Science News Wire () cited "numerous deep flaws" in null hypothesis significance testing. A ScienceNews article (Siegfried ) on February 7, 2014, said "statistical techniques for testing hypotheses... have more flaws than Facebook's privacy policies." A week later, statistician and "Simply Statistics" blogger Jeff Leek responded. "The problem is not that people use P-values poorly," Leek wrote, "it is that the vast majority of data analysis is not performed by people properly trained to perform data analysis" (Leek ). That same week, statistician and science writer Regina Nuzzo published an article in *Nature* entitled "Scientific Method: Statistical Errors" (Nuzzo ). That article is now one of the most highly viewed *Nature* articles, as reported by altmetric.com (*http://www.altmetric.com/details/2115792#score*).

In this article          ≡

*replicability* of scientific conclusions. Without getting into definitions and distinctions of these terms, we observe that much confusion and even doubt about the validity of science is arising. Such doubt can lead to radical choices, such as the one taken by the editors of *Basic and Applied Social Psychology*, who decided to ban *p*-values (null hypothesis significance testing) (Trafimow and Marks ). Misunderstanding or misuse of statistical inference is only one cause of the "reproducibility crisis" (Peng ), but to our community, it is an important one.

When the ASA Board decided to take up the challenge of developing a policy statement on *p*-values and statistical significance, it did so recognizing this was not a lightly taken step. The ASA has not previously taken positions on specific matters of statistical practice. The closest the association has come to this is a statement on the use of value-added models (VAM) for educational assessment (Morganstein and Wasserstein ) and a statement on risk-limiting post-election audits (American Statistical Association ). However, these were truly policy-related statements. The VAM statement addressed a key educational policy issue, acknowledging the complexity of the issues involved, citing limitations of VAMs as effective performance models, and urging that they be developed and interpreted with the involvement of statisticians. The statement on election auditing was also in response to a major but specific policy issue (close elections in 2008), and said that statistically based election audits should become a routine part of election processes.

By contrast, the Board envisioned that the ASA statement on *p*-values and statistical significance would shed light on an aspect of our field that is too often misunderstood and misused in the broader research community, and, in the process, provides the community a service. The intended audience would be researchers, practitioners, and science writers who are not primarily statisticians. Thus, this statement would be quite different from anything previously attempted.

In this article          ≡

two dozen such people, all of whom said they would be happy to be involved. Several expressed doubt about whether agreement could be reached, but those who did said, in effect, that if there was going to be a discussion, they wanted to be involved.

Over the course of many months, group members discussed what format the statement should take, tried to more concretely visualize the audience for the statement, and began to find points of agreement. That turned out to be relatively easy to do, but it was just as easy to find points of intense disagreement.

The time came for the group to sit down together to hash out these points, and so in October 2015, 20 members of the group met at the ASA Office in Alexandria, Virginia. The 2-day meeting was facilitated by Regina Nuzzo, and by the end of the meeting, a good set of points around which the statement could be built was developed.

The next 3 months saw multiple drafts of the statement, reviewed by group members, by Board members (in a lengthy discussion at the November 2015 ASA Board meeting), and by members of the target audience. Finally, on January 29, 2016, the Executive Committee of the ASA approved the statement.

The statement development process was lengthier and more controversial than anticipated. For example, there was considerable discussion about how best to address the issue of multiple *potential* comparisons (Gelman and Loken ). We debated at some length the issues behind the words "a *p*-value near 0.05 taken by itself offers only weak evidence against the null hypothesis" (Johnson ). There were differing perspectives about how to characterize various alternatives to the *p*-value and in how much detail to address them. To keep the statement reasonably simple, we did not address alternative hypotheses, error types, or power (among other things), and not everyone agreed with that approach.

As the end of the statement development process neared, Wasserstein contacted

In this article          ≡

good platform to reach a broad and general statistical readership. Together, we decided that the addition of an online discussion would heighten the interest level for the *TAS* audience, giving an opportunity to reflect the aforementioned controversy.

To that end, a group of discussants was contacted to provide comments on the statement. You can read their statements in the online supplement, and a guide to those statements appears at the end of this editorial. We thank Naomi Altman, Douglas Altman, Daniel J. Benjamin, Yoav Benjamini, Jim Berger, Don Berry, John Carlin, George Cobb, Andrew Gelman, Steve Goodman, Sander Greenland, John Ioannidis, Joseph Horowitz, Valen Johnson, Michael Lavine, Michael Lew, Rod Little, Deborah Mayo, Michele Millar, Charles Poole, Ken Rothman, Stephen Senn, Dalene Stangl, Philip Stark and Steve Ziliak for sharing their insightful perspectives.

Of special note is the following article, which is a significant contribution to the literature about *p*-values and statistical significance.

> Greenland, S., Senn, S.J., Rothman, K.J., Carlin, J.B., Poole, C., Goodman, S.N. and Altman, D.G.: ``Statistical Tests, *P*-values, Confidence Intervals, and Power: A Guide to Misinterpretations."

Though there was disagreement on exactly what the statement should say, there was high agreement that the ASA should be speaking out about these matters.

Let us be clear. Nothing in the ASA statement is new. Statisticians and others have been sounding the alarm about these matters for decades, to little avail. We hoped that a statement from the world's largest professional association of statisticians would open a fresh discussion and draw renewed and vigorous attention to changing the practice of science with regards to the use of statistical inference.

In this article　　　≡

Many of the participants in the development of the ASA statement contributed commentary about the statement or matters related to it. Their comments are posted as online supplements to the statement. We provide here a list of the supplemental articles.

## Supplemental Material to the ASA Statement on *P*-Values and Statistical Significance

*Altman, Naomi*: Ideas from multiple testing of high dimensional data provide insights about reproducibility and false discovery rates of hypothesis supported by *p*-values

*Benjamin, Daniel J, and Berger, James O:* A simple alternative to *p*-values

*Benjamini, Yoav:* It's not the *p*-values' fault

*Berry, Donald A:* *P*-values are not what they're cracked up to be

*Carlin, John B:* Comment: Is reform possible without a paradigm shift?

*Cobb, George:* ASA statement on p-values: Two consequences we can hope for

*Gelman, Andrew:* The problems with *p*-values are not just with *p*-values

*Goodman, Steven N:* The next questions: Who, what, when, where, and why?

*Greenland, Sander:* The ASA guidelines and null bias in current teaching and practice

*Ioannidis, John P.A.:* Fit-for-purpose inferential methods: abandoning/changing

In this article          ≡

*Johnson, Valen E.:* Comments on the "ASA Statement on Statistical Significance and *P*-values" and marginally significant *p*-values

*Lavine, Michael, and Horowitz, Joseph:* Comment

*Lew, Michael J:* Three inferential questions, two types of *P*-value

*Little, Roderick J:* Discussion

*Mayo, Deborah G:* Don't throw out the error control baby with the bad statistics bathwater

*Millar, Michele:* ASA statement on *p*-values: some implications for education

*Rothman, Kenneth J:* Disengaging from statistical significance

*Senn, Stephen:* Are *P*-Values the Problem?

*Stangl, Dalene:* Comment

*Stark, P.B.:* The value of *p*-values

*Ziliak, Stephen T:* The significance of the ASA statement on statistical significance and *p*-values

---

Ronald L.Wasserstein and Nicole A. Lazar

ron@amstat.org

*American Statistical Association, 732 NorthWashington Street,*

*Alexandria, VA 22314-1943.*

---

In this article     ☰

# and *P*-values

Ronald L. Wasserstein[a]

## 1. Introduction

Increased quantification of scientific research and a proliferation of large, complex datasets in recent years have expanded the scope of applications of statistical methods. This has created new avenues for scientific progress, but it also brings concerns about conclusions drawn from research data. The validity of scientific conclusions, including their reproducibility, depends on more than the statistical methods themselves. Appropriately chosen techniques, properly conducted analyses and correct interpretation of statistical results also play a key role in ensuring that conclusions are sound and that uncertainty surrounding them is represented properly.

Underpinning many published scientific conclusions is the concept of "statistical significance," typically assessed with an index called the *p*-value. While the *p*-value can be a useful statistical measure, it is commonly misused and misinterpreted. This has led to some scientific journals discouraging the use of *p*-values, and some scientists and statisticians recommending their abandonment, with some arguments essentially unchanged since *p*-values were first introduced.

In this context, the American Statistical Association (ASA) believes that the scientific community could benefit from a formal statement clarifying several widely agreed upon principles underlying the proper use and interpretation of the *p*-value. The issues touched on here affect not only research, but research funding, journal practices, career advancement, scientific education, public policy, journalism, and law. This statement does not seek to resolve all the issues relating to sound statistical practice, nor to settle foundational controversies. Rather, the statement articulates in

In this article            ☰

statistical community.

## 2. What is a *p*-Value?

Informally, a *p*-value is the probability under a specified statistical model 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.

## 3. Principles

1. **_P_-values can indicate how incompatible 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 proposed model for the data. The most common context is a model, constructed under a set of assumptions, together with a so-called "null hypothesis." 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, if the underlying assumptions used to calculate the *p*-value hold. This incompatibility can be interpreted as casting doubt on or providing evidence against the null hypothesis or the underlying assumptions.

2. **_P_-values do not 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

In this article         ≡

a specified hypothetical explanation, and is not a statement about the explanation itself.

3. **Scientific conclusions and business or policy decisions should not 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 many contextual factors 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 not 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 distortion of the scientific process.

4. **Proper inference requires full reporting and transparency**

   *P*-values and related analyses should not 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 essentially uninterpretable. Cherry-picking promising findings, also known by such terms as data dredging, significance chasing, significance questing, selective inference, and "*p*-hacking," leads to a spurious excess 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: Whenever a researcher chooses what to present based on statistical results, valid interpretation of those results is severely compromised if the

In this article    ≡

decisions, all statistical analyses conducted, and all *p*-values computed. Valid scientific conclusions based on *p*-values and related statistics cannot be drawn without at least knowing how many and which analyses were conducted, and how those analyses (including *p*-values) were selected for reporting.

5. **A *p*-value, or statistical significance, does not measure the size of an effect or the importance of a result.**

   Statistical significance is not equivalent to scientific, human, or economic significance. Smaller *p*-values do not 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.

6. **By itself, a *p*-value does not 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; many other hypotheses may be equally or more consistent with the observed data. For these reasons, data analysis should not end with the calculation of a *p*-value when other approaches are appropriate and feasible.

---

# 4. Other Approaches

In view of the prevalent misuses of and misconceptions concerning *p*-values, some

In this article          ≡

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

## 5. Conclusion

Good statistical practice, as an essential component of good scientific practice, emphasizes principles of good study design and conduct, a variety of numerical and graphical summaries of data, understanding of the phenomenon under study, interpretation of results in context, complete reporting and proper logical and quantitative understanding of what data summaries mean. No single index should substitute for scientific reasoning.

## Acknowledgments

In this article        ≡

Edited by Ronald L.Wasserstein, Executive Director

*On behalf of the American Statistical Association*

*Board of Directors*

# Supplemental material

In this article          ≡

Process, and Purpose

Showing 1/23: utas_a_1154108_sm5368.pdf

# Statistical Tests, *P*-values, Confidence Intervals, and Power: A Guide to Misinterpretations

Sander GREENLAND, Stephen J. SENN, Kenneth J. ROTHMAN, John B. CARLIN, Charles POOLE, Steven N. GOODMAN, and Douglas G. ALTMAN

Misinterpretation and abuse of statistical tests, confidence intervals, and statistical power have been decried for decades, yet remain rampant. A key problem is that there are no interpretations of these concepts that are at once simple, intuitive, correct, and foolproof. Instead, correct use and interpretation of these statistics requires an attention to detail which seems to tax the patience of working scientists. This high cognitive demand has led to an epidemic of shortcut definitions and interpretations that are simply wrong, sometimes disastrously so—and yet these misinterpretations dominate much of the scientific literature.

In light of this problem, we provide definitions and a discussion of basic statistics that are more general and critical than typically found in traditional introductory expositions. Our goal is to provide a resource for instructors, researchers, and consumers of statistics whose knowledge of statistical theory and technique may be limited but who wish to avoid and spot misinterpretations. We emphasize how violation of often unstated analysis protocols (such as selecting analyses for presentation based on the *P*-values they produce) can lead to small *P*-values even if the declared test hypothesis is correct, and can lead to large *P*-values even if that hypothesis is incorrect. We then provide an explanatory list of 25 misinterpretations of *P*-values, confidence intervals, and power. We conclude with guidelines for improving statistical interpretation and reporting.

## Introduction

Misinterpretation and abuse of statistical tests has been decried for decades, yet remains so rampant that some scientific journals discourage use of "statistical significance" (classifying results as "significant" or not based on a *P*-value) (Lang et al. 1998). One journal now bans all statistical tests and mathematically related procedures such as confidence intervals (Trafimow and Marks 2015), which has led to considerable discussion and debate about the merits of such bans (e.g., Ashworth 2015; Flanagan 2015).

Despite such bans, we expect that the statistical methods at issue will be with us for many years to come. We thus think it imperative that basic teaching as well as general understanding of these methods be improved. Toward that end, we attempt to explain the meaning of significance tests, confidence intervals, and statistical power in a more general and critical way than is traditionally done, and then review 25 common misconceptions in light of our explanations. We also discuss a few more subtle but nonetheless pervasive problems, explaining why it is important to examine and synthesize all results relating to a scientific question, rather than focus on individual findings. We further explain why statistical tests should never constitute the sole input to inferences or decisions about associations or effects. Among the many reasons are that, in most scientific set-

**00_GREENLAND-ETAL.PDF**

## Related Research Data

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

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

In this article      ≡

The ASA's statement on p-values: context, process, and purpose

Linking provided by Schole**plorer

## People also read

Article

### Quality Control for Scientific Research: Addressing Reproducibility, Responsiveness, and Relevance  ›

Douglas W. Hubbard et al.

The American Statistician
Volume 73, 2019 - Issue sup1

**Published online:** 20 Mar 2019

Article

### The *p*-value Function and Statistical Inference  ›

D. A. S. Fraser

In this article          ≡

Article

## Statistical Inference Enables Bad Science; Statistical Thinking Enables Good Science  ›

Christopher Tong

The American Statistician
Volume 73, 2019 - Issue sup1

**Published online:** 20 Mar 2019

Article

## An Introduction to Second-Generation $p$-Values  ›

Jeffrey D. Blume et al.

The American Statistician
Volume 73, 2019 - Issue sup1

**Published online:** 20 Mar 2019

Article

## Assessing Statistical Results: Magnitude, Precision, and Model Uncertainty  ›

In this article     ☰

Article

# The False Positive Risk: A Proposal Concerning What to Do About $p$-Values   ›

David Colquhoun

In this article     ☰