¹ Covered in P: The consequences of null hypothesis significance testing on point and interval

² estimates

³ Raphael T. Gerraty[1], Matti Vuorre[1]

⁴ [1] Columbia University, Department of Psychology

⁵ Author note

⁶ The authors declare no conflicts of interest.

⁷ Correspondence concerning this article should be addressed to Raphael T. Gerraty,

⁸ Columbia University, Psychology Department, 406 Schermerhorn, 1190 Amsterdam Avenue,

⁹ New York, NY 10027. E-mail: https://github.com/neurostorm

Abstract

Confidence intervals (CI) do not allow post-data inference on parameter values.

*Keywords:* confidence interval, NHST


Word count: Short and sweet.

15 Covered in P: The consequences of null hypothesis significance testing on point and interval

16                                                  estimates

## Introduction

18      Scientists, psychological or otherwise, routinely use null hypothesis significance testing

19 procedures (NHSTP) to move from data to conclusions—a practice thats applicability has

20 been debated since its inception. Recently, concerns about the replicability and reliability of

21 empirical findings (Collaboration, 2015) have underlined the concerns about NHSTP as *the*

22 valid form of statistical inference (Cumming, 2014; Gelman & Loken, 2014).

23      One response to the growing concerns regarding the reliability of NHSTP has been an

24 appeal to effect size and interval estimation in addition—or as replacement—to NHSTP test

25 statistics (Cumming, 2014). For example, many journals in psychology and neuroscience now

26 ask authors to include *confidence intervals* (CI) with their test statistics. These confidence

27 intervals are intended, by practitioners, to communicate unbiased estimates of likely ranges

28 of parameter values, although the common computations simply result in ranges of

29 parameter values that would not be rejected by the very statistical test the CI is supposed to

30 supplement. This approach is problematic from at least two perspectives:

31  - CIs do not support probability statements about parameters

32  - CIs are designed (Neyman, 1957) to be a procedure of generating a set of intervals

33     that, as a whole set, contain the true parameter value X% of the time

34     – This property is severely compromised by the current practice of using CIs as a

35        post-data inferential tool, as we show in this paper

## A Confidence Interval

37      While the first of these interpretations is inaccurately held by many researchs in

38 psychology and other fields (citing stuff),

<sup>39</sup> In this paper, we report an unappealing property of confidence intervals. Because the

<sup>40</sup> claim of confidence intervals is to have a coverage proportion of the true parameter value

<sup>41</sup> equal to the nominal value (usually 95%), it is crucial that this claim is substantiated in its

<sup>42</sup> long-run property for a CI to be what it claims to be (not a *confidence* interval.) We show

<sup>43</sup> that using confidence intervals *in addition* to P values leads to an undesirable distortion of

<sup>44</sup> the coverage proportion. This is a direct result of the more general problem with interpreting

<sup>45</sup> any particular *obtained* CI in terms of frequency coverage, but we feel the specific case of

<sup>46</sup> significance thresholding on this interpretation is worth describing, given the per

<sup>47</sup> **P stains the nominal coverage proportion**

<sup>48</sup> <div align="center">**Methods**</div>

<sup>49</sup> We performed a simulation study. . .

<sup>50</sup> <div align="center">**Results**</div>

<sup>51</sup> [1] 0.6457245

<sup>52</sup> <div align="center">**Discussion**</div>

<sup>53</sup> Here we show a pervasive bias in the paramaters and intervals passing a null

<sup>54</sup> hypothesis significance threshold. We don't know if this result is well known to statisticians,

<sup>55</sup> but from the perspective of practitioners, we found it suprising. This paper was motivated in

<sup>56</sup> part by the discussions with colleagues who were equally suprised by the biases induced by

<sup>57</sup> hypothesis testing, especially on interval estimation. The "significance filter" has been

<sup>58</sup> discussed previously (Gelman, 2011), but to our knowledge there have been no discussions of

<sup>59</sup> the effect of this filter on the frequency properties of confidence intervals.

<sup>60</sup> We note that, while the issues discussed in this paper are related to questionable

<sup>61</sup> research practices as well as known issues in null hypothesis testing such as alpha inflation

<sup>62</sup> due to multiple comparisons, the biased point estimates and interval coverage for significant

results we discuss here are present in expectation even for single tests. Thus this bias will be
more severe for significant results which have been filtered through such processes, but. . .

## References

Collaboration, O. S. (2015). Estimating the reproducibility of psychological. *Science*,
    *349*(6251), aac4716. doi:10.1126/science.aac4716

Cumming, G. (2014). The new statistics why and how. *Psychological Science*, *25*(1), 7–29.
    doi:10.1177/0956797613504966

Gelman, A. (2011, October 9). *The statistical significance filter. Statistical modeling, causal
    inference, and social science.* Retrieved December 9, 2015, from
    http://andrewgelman.com/2011/09/10/the-statistical-significance-filter/

Gelman, A., & Loken, E. (2014). The statistical crisis in science. *American Scientist*, *102*(6),
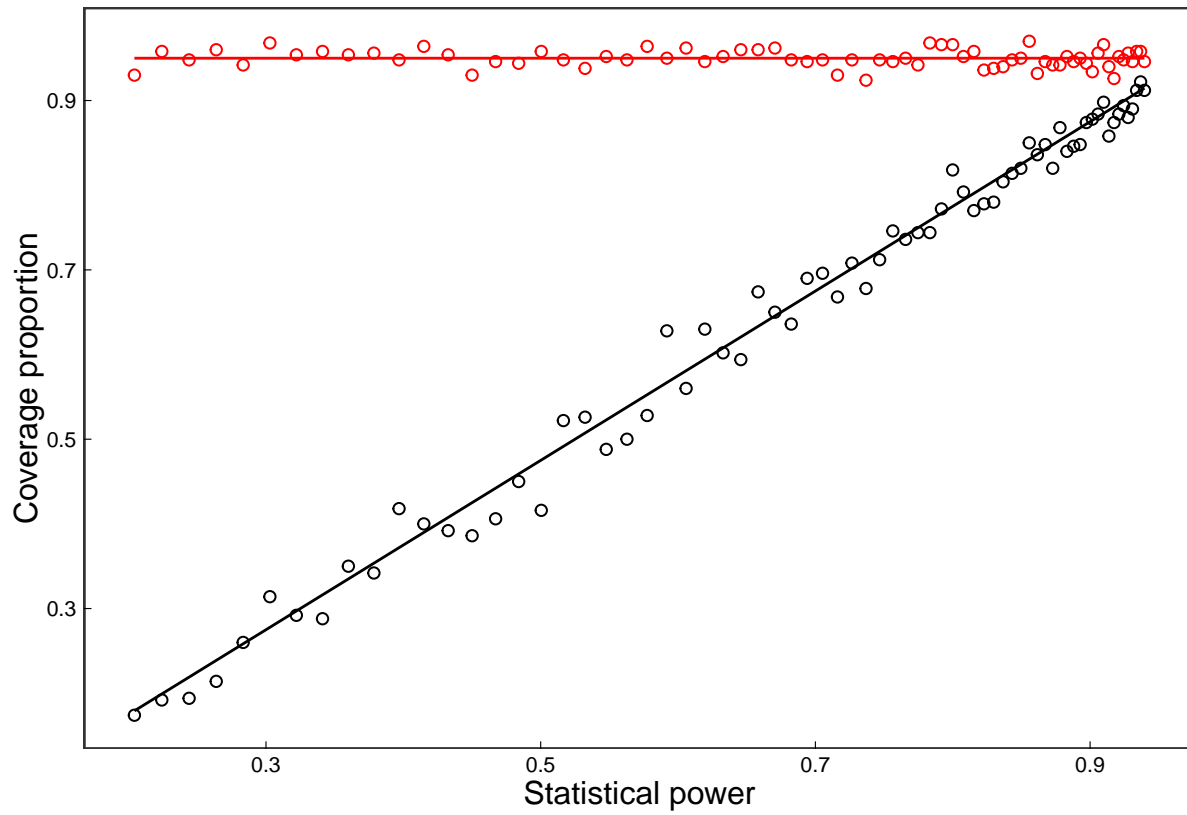    460. doi:10.1511/2014.111.460

*Figure 1*. Coverage proportion of confidence intervals versus statistical power. The red line represents the nominal coverage proportion (95%), black line is the actual coverage proportion when CIs are conditioned on the result being significant.