

The manuscript “On Data ‘Janitor Work’ in Political Science: The Case of Thermostatic Support for Democracy” is an important replication study. I recommend publishing it with minor revisions.

PSRM is providing an innovative contribution to scientific research by explicitly publishing replication studies as an article format. This follows a wider trend in political science over the last decade where academic standards and best practices within the discipline and within science in general are reassessed. This is mainly done and pushed forward by the open science movement. This movement highlights the importance of replication and open availability of research materials such as full replication materials (data, source code, documentation).

The article on the “janitor work” adds to this debate and provides two contributions. First, it highlights the relevance of data-entry-errors and provides a general assessment of its implications for political science research. Second, it replicates Claassen (2020) an important study on how changes in democracy affect democratic support among the public. The results of the replication demonstrate that the effects the Claassen (2020) study establishes do not hold when data-entry-errors are accounted for.

In my view, the replication analysis is convincing, valuable and should be published. Nevertheless, I believe that the manuscript would benefit from presenting the main argument about the relevance of data-entry-errors within a more general open science perspective. I provide details and recommendations in the sections below.

First and foremost, the authors must be congratulated for the resource intense and time-consuming replication of Claassen (2020) by creating anew the dataset used in the original study based on the primary data sources (surveys). The authors demonstrate that the coding criteria that Claassen (2020) defines are not applied coherently when creating the compiled dataset used in the original analysis. The authors recreate the compiled dataset based on the original data sources (surveys). They demonstrate that a reanalysis with the recreated new dataset does not replicate the main findings of the original study. Claassen (2020) established a finding that public support responds thermostatically to changes in democracy. This finding is in difference to other studies that demonstrate growing democratic experience generates robust public support for democracy. Hence, it is an important contribution to establish that the findings of Claassen (2020) do not hold after data-entry-errors are considered.

I am a little more skeptical about the presentation of the relevance of data-entry-errors. I agree with the argument that a full perspective on replication needs to start from the original data source(s). Hence the creation of compiled data sets according to established coding rules needs to be replicated as well. I have two suggestions on improving the presentation of the data-entry-error argument.

First, I think that the manuscript would benefit from establishing the relevance of “data janitor” work more broadly. Currently, the relevance is established based on two references (Bachard and Pace 2011, Lohr 2014) one of them a newspaper article. I am convinced that a literature search would lead to references that are more relevant and recent. The argument about the relevance of “data janitor” is for example made repeatedly by Hadley Wickham in his publications on “tidy” data.

Second, I am convinced that the argument about the relevance of data-entry-errors and “data janitor” work can be better linked to the emerging literature in political science that advocates using open science principles. Open science advocates argue that the entire elements of research work should be made publicly available (e.g. all data, software, documentation, etc.). The author’s suggestion to “automate data entry” (p. 9) in the discussion relates to this argument.

In this study and in Claassen (2020) the creation of the compiled data set is based on existing surveys and could be fully automated. Hence, the compilation of the dataset can be established by well documented and tested computer code where the selected variables and the recoding thresholds are documented in a data file. For supporters of open science this establishment of high-quality documentation of all research elements that lead to an article is a key element (e.g., Christensen 2019: Transparent and reproducible social science research, ch. 11). Hence, I would encourage the authors to make this link to the open science literature more explicit and to present it earlier in the manuscript (see also Special issue PS: Political Science & Politics: “Opening Political Science”).