

Report on MS-2344: **A Replication Plan for “Does Social Media Reduce Corruption?”**
(Information Economics and Policy, 2017)

Summary: The present manuscript discusses the replication plan using a recent paper titled “Does Social Media Reduce Corruption?” published in *Information Economics and Policy* in 2017. The manuscript discusses why replication is important, how to properly replicate a paper, and what are the problems associated with replication.

I did learn quite a bit about the replication process and why it is important. I think the author has done an excellent job of discussing why a replication is important and the importance of replicating the paper in correct way. It is also careful to note that a replicator can also be fallible and may have his/her own biases. The author is also careful to note that a replication may fail to produce the exact results reported in the paper and provides reason for this. He then goes on to note that this does not necessarily mean that author(s) of the original paper (in this instance, Jha and Sarangi) may have altered the result intentionally to obtain the desired conclusions.

In the strictest sense, the replication of a paper must be done using the same dataset and the methodology as described in the paper that is being replicated. If the results do not hold after that, it casts a doubt on the validity of the published results. However, in a less restrictive sense, replication can also involve some robustness checks. These checks, in my opinion, should be performed using the same methodology and variables, but may involve updated dataset. Note, however, that in this case if the replication results do not agree with the published results, it does not mean an error on part of the authors of the paper being replicated but simply that the result (for example using the Jha and Sarangi (2017) {*JS henceforth*} paper as an example, the relationship between social media and corruption) may not persist over time for several reasons including the possibility of a dynamic relationship.

Keeping such an approach in mind, my main comments are as follows:

1. The first important issue the manuscript raises is about data availability and data updates. This issue is especially important for cross-country regressions because several data sources such as the World Bank update certain data multiple times. An example would be GDP per capita. Since the paper being replicated here is a cross-country study, this is an important issue. In case of the paper that is being replicated in this manuscript, the dependent variable, Control of Corruption Index, is also revised over time. As a result, it is important, and I agree with the author, that first the authors of the paper being replicated (JS in this case) should be contacted for the data to ensure that correct version of the data is being used. Since it may take many years for a paper to be published in a journal – papers are first circulated as working papers, then are under review for several months, and for most papers a few more months under revision – the data are more likely than not to be revised during this time period. Hence, recreating the entire data is unlikely to reproduce the exact dataset used by the authors of the study that is being replicated. I think the manuscript should elaborate on this issue before noting that the new data should be collected anyways and be used to replicate the original paper (which is important as well).

2. The author has been very careful to find at least one country, “Argentina”, that was not included in the original JS analysis, which, according to the author, should reasonably have been included, particularly since the JS paper does not provide any reasons for this omission.

This raises an important question – would the inclusion of this country significantly alter the published results? While JS may have had their reasons for doing this, any meaningful replication in the absence of a valid reason, should include this country to obtain new results. The replicator should possibly also solicit from the original authors (JS in this case) possible reasons for the omission. The author could then present both sets of results if they are significantly different and offer explanations as to why this might be the case. This would significantly add value to the exercise. Thus, I would like the manuscript to discuss some of these issues.

3. Regarding the methodology, the replicator must use the same methodology as used in the paper being replicated. The replicator however must verify that the paper being replicated does indeed use the methodology it claims to use. This subtle difference needs to be clearly spelt out in the paper.
4. The author mentions that “it might seem reasonable to wonder whether changes in corruption are correlated with changes in corruption”. In my opinion (and the author also concurs), this is not replicating the paper, rather investigating this relationship using a different (extended) dataset and a different methodology. This is no longer just replicating the paper but reinvestigating the issue, which, though important, is not the issue at hand. In my opinion, this fact needs to be either clearly explained or this part should be removed from the paper. This is clearly not replication, but something that may be worth exploring.
5. In a liberal replication plan, in case of a cross-country study such as the present study being replicated, one can also check the robustness of the results/conclusions using a different year’s data: do the conclusions stay the same or do they change? This robustness check would ensure that published results are not driven by any year-specific events. The manuscript currently does not discuss this issue, and may consider adding a discussion of this, subject to data availability.
6. Another point I like is the fact that the author identifies that the paper being replicated does not include “Protestant share” in the regressions unlike several papers that are cited therein. The author however does not discuss the likely reasons of such an omission. Does the published version of the paper mention a reason why it was not included? Were there other variables included to replace/proxy this variable? Was it because the questions of the paper being replicated was different from those that are cited and hence this variable (that is, protestant share) may not be crucial? Finally, if the data on this variable are available, one can check the robustness of the results by including this variable as a regressor. I would like to see a discussion on this.