

The Reproducibility of Economics Research: A Case Study

Hautahi Kingi* Flavio Stanchi Lars Vilhuber Sylvérie Herbert†

March 29, 2019

Preliminary - do not cite or quote without contacting authors.

Abstract

Published reproductions or replications of economics research are rare. However, recent years have seen increased recognition of the important role of replication in the scientific endeavor. We describe and present the results of a large reproduction exercise in which we assess the reproducibility of research articles published in the American Economic Journal: Applied Economics over the last decade. 69 of 162 eligible replication attempts successfully replicated the article's analysis 42.6%. A further 68 (42%) were at least partially successful. A total of 98 out of 303 (32.3%) relied on confidential or proprietary data, and were thus not reproducible by this project. We also conduct several bibliometric analyses of reproducible vs. non-reproducible articles.

JEL codes: B41; C80; C81; C87; C88

Keywords: Replication

*IMPAQ International. hautahikingi@gmail.com

†Labor Dynamics Institute, Cornell University. Stanchi: fs379@cornell.edu, Vilhuber: lars.vilhuber@cornell.edu, Herbert: sh2258@cornell.edu

1 Introduction

Replication, reproduction, and falsification of published articles is an important part of the scientific endeavor, which helps to make science “robust and reliable” (**Bollen2015-vb**). The ability (or lack thereof) to reproduce an article has been discussed in economics for at least thirty years (**Dewald1986**; **Vinod2005**; **AndersonEtAl2005**; **King95**; **BurmanEtAl2010**; **Duvendack2017**; **Hamermesh2017**; **Sukhtankar2017**; **Hoeffler2017**; **Coffman2017**; **Chang2017**; **Berry2017**; **Anderson2017**).

Though not unheard of (**Hoeffler2017a**; **Chang2017**; **ChangLi2015**; **camerer2016**), actual published reproductions or replications are rare (**BellMiller2013b**; **Duvendack2017**). For example, **MuellerLanger18** found that just 0.1% of the 126,505 articles published between 1974 and 2014 in the top 50 economics journals were replications. **sukhtankar17** found that, of the 1,138 empirical development economics articles published between 2000 and 2015 in the “top 10”¹ economics journals, just 6.2% were replicated in a published or working paper. The paucity of replications in economics is, in part, because it is often difficult to find the materials required to conduct reproducibility or replication exercises (**Dewald1986**; **McCullough2006**; **McCullough03**). Despite a long standing explicit recognition of the importance of replication in economics (**Frisch1933**), it has been suggested that “there is no tradition of replication in economics” (**McCullough2006**).²

The scientific community is actively engaged in identifying ways to support replication (**Bollen2015-vb**; **NAP25116**). More and more journals are adopting “data and code availability” policies (the American Economic Association (**AEA**) has had one since 2005), though some doubt their effectiveness (**Stodden2018**; **Hoeffler2017**). **Duvendack2015** finds that 27 of the 333 economics journals listed in the Thomson Reuters *Web of Science* as of September 2013 regularly publish data and code for empirical articles, and 10 of those journals explicitly state that they publish replication studies. While that number seems low, it is higher than it was a decade earlier. More recently, the Journal of the American Statistical Association (**JASA**) has moved towards much more stringent replication requirements (**Fuentes2016-wz**), and the **AEA** in 2017/2018 appointed as Data Editor the last author in lexicographic order of this article (**10.1257/pandp.108.745**).³

A variety of replication concepts are used (**Bollen2015-vb**; **Hamermesh2017**; **Clemens2015**). In this article, we adopt the definitions articulated by **Bollen2015-vb**, among others. *Reproducibility* refers to “the ability [...] to duplicate the results of a prior study using the same materials and procedures as were used by the original investigator,” and is related to the “narrow” sense of replication of **Pesaran2003**. Use of the “same procedures” may imply using the same computer code or re-implementing the statistical procedures in a different software package. **Hamermesh2007** calls this “pure replication”, which **Christensen2018** argue is the “basic standard [that] should be expected of all published economics research, and hope this expectation is universal among researchers.” *Replicability* refers to “the ability of a researcher to duplicate the results of a prior study if the same procedures are followed but new data are collected” (**Pesaran2003**), while *generalizability* refers to the extension of the scientific findings to other populations, contexts, and time

¹These were the traditional “top 5” and the American Economic Journal: Applied Economics (**AEJ:AE**), the American Economic Journal: Economic Policy (**AEJ:EP**), the Economic Journal (**EJ**), the Journal of the European Economic Association (**JEEA**) and the Review of Economics and Statistics (**ReStat**).

²Though **Hamermesh2017**; **Hamermesh2007** disagree.

³Much of the research reported in this article was started before the appointment.

frames, perhaps using different methods. Because there is a grey zone between these last two definitions, we will generally refer to either context as “replicability”, which **Hamermesh2017** calls “scientific replication.” In this text, we will use the terms as defined above when the distinction is material. However, we may refer to the overall concept of redoing the analysis as “replicability” (in part because it conjugates better).

In this article, we set out to assess how well a particular journal’s “data availability” policy, combined with light enforcement, yields *reproducible* articles. Our protocol is set with a relatively high bar: can undergraduates, armed only with the information provided by authors on the journal website, successfully reproduce the tables and figures presented by the author in the article? Unlike **Dewald1986** and **McCullough03**, who requested data and programs from the original authors, we did not attempt any contact with authors to clarify issues that arose. While our replicators were instructed to do their best to fix any bugs or inconsistencies that they encountered, they were limited both by time and training.

We conducted this experiment over several summers and during the 2018 Fall semester, using the [AEJ:AE](#) as our source of articles, which we chose primarily for two reasons. First, because of the empirical nature of its articles and its policy of publishing papers “only if the data used in the analysis are clearly and precisely documented and are readily available to any researcher for purposes of replication,” we expect that nearly all articles have some empirical component. Second, while other journals may also have theoretical or more complex empirical papers, using a variety of software, we wanted articles to be reproducible by the undergraduate student armed with knowledge of Stata and Matlab only. However, nothing in the methodology used in this paper is specific to the [AEJ:AE](#), and we are planning to expand to other journals. Part of the motivation here is also to test the feasibility of “pre-publication verification” similar to what is done at the American Journal of Political Science ([AJPS](#)) (**JacobyShouldJournalsBe2017**; **Christian2018**)

We start by describing, in Section 2, our methods of selection, analysis, and reproduction. We present our results in Section 3. A preliminary discussion closes the paper.

2 Description of Reproduction Procedure

Our replication exercise was conducted over the course of four summers from 2014 to 2018 as well as the fall semester of 2018. This section describes the procedure used to conduct the replication exercise, which was split into 3 parts. First, an “assessor” evaluated the task of replication as to the availability of the required components and its difficulty, and recorded a selection of article characteristics. A “replicator” then conducted the actual replication. Often but not always, the same person would be assessor and replicator. Finally, upon completion of the replication attempt, replicators filled out a guided report with questions about the exercise, such as whether or not the main results of the article could be replicated and, if not, the main barriers impeding a successful replication. To complement the information collected during the replication exercise, we also obtained descriptive article and author information such as citations and h-indices. Each of these steps are discussed in turn in the following subsections.

2.1 Initial Assessment

We first assessed each article. An assessor filled out a questionnaire (see Appendix B), gathering descriptive information, and providing an initial assessment of the expected level of ‘replicability’ of an article. Each assessor was provided the Digital Object Identifier (DOI) of the article,⁴ and then verified the following elements:

- The presence of one or more downloadable datasets (including DOI or URL, if any), an online appendix, the programs used by the authors and documentation on how to run them (the “Readme”);⁵
- The clarity and completeness of the documentation and of the program metadata;
- The presence of clear references to the original data provenance and a description of how to construct the initial datasets;
- Data availability (e.g., restricted access data, private data, public use data, etc.)

Although some of the responses to the questionnaire could have been captured via web-scraping tools, it is not possible to assess the completeness of the supplemental data without inspection by the assessor. While many supplemental data packages contained some content, they often did not contain all the data or programs. Sometimes, the data provenance might be described in an online appendix, while the instructions for the programs might be enclosed in the supplemental “data” package. Thus, a “clerical” review of each article’s webpage, and some careful reading of the actual article and online appendix were the only way to collect all the information requested. We will return to the aspect of machine-readability (machine-reproducibility) in the concluding discussion.

Based on the initial objective enumeration of the characteristics of the article and on subjective evaluation by the assessor of the complexity of the task described in the “Readme” document, the assessor was asked to provide a subjective rating of the replication difficulty, from 1 (easiest) to 5 (most difficult), based on a set of heuristics (see Table 1).⁶ Assessors also recorded the programming languages and separately the data storage formats contained within each archive. Archives can and do contain programs and data in multiple formats. While all articles had some supplementary data, not all articles were accompanied by the datasets necessary for replication. Assessors recorded whether the articles were accompanied by a dataset and, if not, the (apparent) reason why data was not provided. This included an assessment of whether the data were confidential or proprietary, for which we provided some guidance.

⁴DOIs are a managed identifier space built on top of the Handle System (**Handle**), a technology for distributed, persistent, and unique naming for digital objects. Virtually all academic publishers assign DOIs at the article level in all of their publications. In addition, DOIs are increasingly used to identify data (**PollardWilkinson2010**). In particular, each DOI provides a persistent identifier (**DOI2012**) for a digital object: an article or data artifact.

⁵In theory, the author might have provided a DOI to a third-party data archive for some or all of the content. Ideally, each component – the article, the online appendix, the data and the programs – would have a separate DOI. In the case of the **AEJ:AE**, only the article itself is assigned a DOI. Supplemental data, programs, and online appendices are linked from the landing page associated with the article’s DOI.

⁶A score of 1 assigned to an article does not imply that its ‘replicability’ cannot be improved. For example, the programs provided by authors might lack a complete header or DOI for each database, but overall the article appears to be easily replicable.

Table 1: Criteria for Assessment of Difficulty of Replication

Rating	Description
1	The article possesses all desired features that ensure replicability. Datasets are provided and their use is public. The documentation and program metadata are clear and complete. Negligible changes might be required to run the programs (e.g. path redirection).
2	The article possesses most desired features that ensure replicability. Datasets are provided and their use is public. The documentation and program metadata are present but the programs might need some changes to run cleanly.
3	The article replication may present some difficulties. Datasets are provided and their use is public, but the documentation is incomplete or unclear. Substantial changes might be needed to run the programs.
4	The article replication may present substantial difficulties and/or additional steps are required to recover the datasets used by the authors. Datasets may not be provided but their use is public or available on request.
5	The article is not replicable. The dataset are not provided and their access is private or restricted. The programs are not provided. Documentation is absent or incomprehensible.

2.2 Reproducing the Empirical Analyses

Once the article was assessed and determined to be amenable for replication, a replicator was assigned based on characteristics of the article, and in particular the type of software required and the assessed difficulty of the replication task. Most often, the initial assessor self-assigned themselves as the replicator although some replications were conducted by the supervising team or distributed to a replicator based on their familiarity with a particular programming language. The replicator was instructed to document any changes made to the author-provided programs using a version control system (VCS).⁷ The materials for each article were downloaded, and used to populate a article-specific repository. If multiple replicators worked on the same article, they would work in separate subdirectories of the same repository. In addition to recording changes in the VCS, replicators were asked to record information in two additional files. They recorded one or more Uniform Record Locators (URLs) of materials obtained in ‘*SRC.txt*’. And, in addition to any VCS commit messages, they were asked to provide a high-level summary of steps undertaken for the replication and results obtained in ‘*REPLICATION.txt*’.

Once the VCS area was populated and all data files were downloaded, each replicator was instructed to read the author-provided “Readme”, and attempt to run the programs in the author-provided archive. Replicators were told to keep modifications to the absolute minimum, starting with the adjustment of path structures to the replication system (Figure 1). Whenever a program required more extensive changes to run, the replicator would do so to the best of their ability. Replicators were free to use any computer they had access to for the replication, unless the author materials specifically mandated a particular operating system. This was quite rare, but did happen a few times. We note that the journal does not require authors to specify the required software, operating system (OS), and versions thereof, and most articles were silent on the topic. In addition to replicators’ laptops, our team had access to university-provided Windows remote desktops and a Linux cluster, and were unlikely to be constrained by computing resources. Both Windows and Linux systems had access to the latest version of the computing software available at the time of replication (Stata

⁷Until August 2018, the team used an restricted-access Subversion repository. Since September 2018, the team uses a restricted-access Git repository.

13 and 14, SAS 9.4, and SPSS 23 at the time of the replication). In general, we assume that versions of OS and software are different from the version originally used by the authors, given the considerable time lag between submission, publication, and the time of the replication exercise. In some cases, programs needed small modifications to run cleanly due to software version discrepancies. If successful, these modifications were treated as ‘negligible’ and did not affect the score of the article. If unsuccessful, however, the article was classified as non- or partially reproducible. We discuss the necessary code changes required for successful replication along with our results in Section 3.⁸

It is important to note that the articles considered were relatively recent, and trying to replicate papers even just a few years older might present more difficulties related to differences in software version. To assess the authenticity of the results, we would ideally use the same software version used by the authors of an article, but such software is often difficult - and, in some cases, impossible - to find or run. Most authors did not provide software version information and, to the best of our knowledge, the journals do not attempt to capture this information from authors. However, based on time-lag to publication, and the age of the articles, we expect that multiple versions of each software lie between when the authors ran their programs and when our team ran the programs. For instance, Matlab updates their software distribution twice a year - for an article published in 2010, it is likely that the version of Matlab used by the author was released in 2008 or 2009, at least six years before we replicated the article.

One way to address the issue of software versioning and other issues such as ambiguity in the documentation of programs would be to reach out to authors directly for confirmation. Unlike **Dewald1986** and **McCullough03**, however, who requested data and programs from the original authors, we only tried to obtain data that was available on the journal’s archive, and replicators were explicitly instructed not to contact the authors. Thus, we did not follow-up on missing datasets, private data, protected data, or data available upon request. All attempts at reproducing the analysis were done based exclusively on the materials provided by the authors at the time of publication of the articles. The entire team discussed problems encountered in regular meetings, and sought to find solutions. However, we also limited the time to find a solution for problems encountered to about a week. If unsuccessful in finding a solution, the replication was marked as “not reproducible.” There was no time limit for running programs.

⁸We captured the modified programs created by the replicators in the VCS. We could, therefore, capture objective measures of code changes using, for instance, the number of code lines changed. We have yet to do this.

Figure 1: Example of change to author-provided file

```
> svn diff -r961:HEAD $SVNURL/10.1257/app.5.4.92/replication-xxx/Data/
do_files/main_analysis.do
Index: main_analysis.do
=====
--- main_analysis.do (revision 961)
+++ main_analysis.do (revision 3425)
@@ -6,7 +6,7 @@
     version 11.2

    *place path here:
    -global path "C:\\"
    +global path "\\rschfs1x\usercl\spring\xxx\Data"
    cd "$path\output"
    use "$path\data\thefts_sales.dta", clear
    cap log close
```

2.3 Report on Reproducibility

Once the replication attempt terminated, replicators completed a questionnaire-based report on outcomes (“Exit Questionnaire”) to capture information about the success or failure of the reproduction attempt and other descriptive information.⁹ We asked replicators to describe the clarity and helpfulness of the documentation provided with the supplementary materials of the article. Although this information is also gathered at the assessment stage described in Section 2.1 to assign a subjective measure of ex-ante replication difficulty, a full understanding of an article’s documentation quality is best left to a replicator who has gone through a replication attempt. We also captured the qualitative nature of the code changes (if any).

2.4 Other data related to the articles

In addition to the data collected through the replication process, we also obtained complementary descriptive data for each article and the (academic) characteristics of the authors. We queried the Web of Science (**Web of Science**) database for each of up to five authors per article, and recorded their h-index (**Hirsch2005**) and the number of citations for each author by year, which is the raw data underlying the calculation of the h-index, as well as the search criteria used to find the author. In some cases, a simple search by author name does not yield a unique person (e.g., “Smith, Adam”), and sometimes, the metadata in Web of Science contained errors.¹⁰ We also obtained citation statistics for each article.

⁹The print version of the online questionnaire is provided in Appendix C.

¹⁰For example, we only found one article for “Lawrence E. Katz” in Web of Science as of January 2016, namely the article in [AEJ:AE](#), but did find quite a few more for “Lawrence F. Katz.” While we initially thought this to be the result of some inside joke for senior economists, even the [AEJ:AE](#) website lists the author of the article as “Lawrence F. Katz,” and we have no explanation for how this error could persist in Web of Science. Our search criteria adjusted for this error.

2.5 Recruitment of Replication Lab Members

Over the course of five summers, we recruited undergraduate students (typically but not always **seniors**) for the replication work as part of summer research, to serve as assessors and replicators. The **seniors** members needed to meet some minimal technical qualifications, such as experience working with the relevant programming language and acceptable performance in economics or technically equivalent courses (it turned out that many of our students had never taken an economics class). Team members attended a one-day training course covering the background and purpose of the replication exercise and our approach. They were also given guidance about the subjective aspects of the exercise such as difficulty rating and classifying documentation clarity, and were instructed on some technical matters such as version control with Subversion and the use of remote computing clusters using materials from a Cornell High Performance Computing course designed for social science researchers. The team members were supervised by economics Ph.D. candidates from Cornell University (Kingi, Stanchi, Herbert) and a faculty-level researcher (Vilhuber).

3 Results

3.1 Descriptive Results

Table 2 presents the by-year breakdown of the 303 **AEJ:AE** articles assessed.¹¹ The subjective measure of “reproducibility” (described in Table 1) is presented in Table 3, which shows a reasonably even distribution of articles across the rating scale.

Table 2: Articles by Year

2009	2010	2011	2012	2013	2014	2015	2016	2017	2018	Total
23	32	36	40	10	40	24	36	42	20	303

Table 3: Replication Difficulty Assessment

Difficulty Rating	Number of Articles	Percent
1	64	21.12
2	64	21.12
3	66	21.78
4	37	12.21
5	72	23.76

The vast majority of articles used the Stata programming language for at least some portion of analysis (Table 4). This preponderance of a single language is reflective of broader usage in economics, though the particular dominance of Stata might be specific to the **AEJ:AE**. From a reproducibility perspective, Stata has both advantages and disadvantages. While it is proprietary software, it is relatively cheap and accessible. Many packages to extend its usability are available, many of which are accessible from within the software from both peer-reviewed (Stata Journal) and crowd-sourced (RePEc/SSC) repositories. Unfortunately, in

¹¹There were 341 assessments for these 303 articles. Articles could be assessed multiple times for a number of reasons, including explicit double-coding, “onboarding” of new replicators, and operator error. At this stage, we consolidated duplicates. An assessment of the reliability of the assessors using the double-coding is pending.

contrast to CRAN, the SSC does not currently support versioning of packages, making it sometimes difficult to find the relevant version of a package. Table 4 also indicates that economists tend to provide data in the native format of the programming language used, instead of open formats (CSV and others). Again, the Stata format has proven to be quite robust, as newly released versions of Stata maintain backward compatibility to all previous versions of the data format. Furthermore, the data format is well understood (albeit not open-source), and can be read by many open-source software packages (R, python). The Stata format allows the embedding of richer metadata, which is not feasible for CSV formats. We did not verify that metadata (variable labels, sane variable names, etc.) complied with modern data curation standards.

Table 4: Programming Languages and Data Formats

Software	Programming Language	Data Format
Stata	281	203
Not Reported	14	76
Matlab	11	3
SAS	9	
R	6	
SPSS	2	1
Excel	1	8
Eviews	0	
Fortran	0	0
Mathematica	0	
CSV		12
RDS		2
txt		2

Totals are not equal to the total number of articles because the articles could use more than one programming language or data format.

While all articles had some supplementary materials, not all articles were accompanied by the datasets necessary for replication. Table 5 details whether the initial assessment declared an article to be eligible because the necessary data was present, based on the README and other materials provided by the authors. In general, if the data was not present, it was due to the confidentiality of the data. 80 articles stated that they used confidential or proprietary data and were therefore not considered for replication, along with the 14 articles with missing data for which no explanation was provided. We note that this is based on an ex-ante assessment, not based on an attempt to actual reproduce the analysis (see Table 8 for causes of reproduction failure due to datasets being missing).

Table 5: Was Data Provided?

Reason	2009	2010	2011	2012	2013	2014	2015	2016	2017	2018	Total
Confidential Data	5	10	8	10	1	15	2	11	11	7	80
Data was Provided	16	22	28	27	9	23	21	23	29	11	209
No Data or Reason	2			3		2	1	2	2	2	14
Total	23	32	36	40	10	40	24	36	42	20	303

The post-completion response questionnaire captured information about 180 of the 209 eligible articles. There were 227 post-completion reports filed for these articles, 35 of which were duplicated articles resulting from multiple reproduction attempts. Of these duplicated articles, 21 arrived at different conclusions about the replication success of that article. In these cases, we kept the replication attempts that resulted in the more successful outcome. Specifically, we define an article to be a successful replication if at least one replicator was able to replicate the results. Similarly, if multiple replications of an article arrived

at a “partially replicated” and “not replicated” conclusion (without a successful attempt), then we say it was partially replicated.¹²

Previous authors have pointed toward the need to improve the documentation of submitted data and programs (McCullough2006; ChangLi2015). Replication attempts are made significantly easier when replicators are not required to resolve ambiguity. Table 6 presents a summary of the documentation quality of the materials provided with the articles for which a reproduction attempt is made categorized by the year in which they were published. 132 articles out of 180 (73.3%) provided complete documentation, defined as a ReadMe file with step by step instructions on how to execute every provided program. However, 45 articles (25%) only provided incomplete ReadMe files that either skipped some of the important steps required to run the programs or contained some ambiguous instructions. No documentation was provided in 1 articles.

Table 6: Documentation Clarity

	2009	2010	2011	2012	2013	2014	2015	2016	2017	2018	Total
Complete	10	8	17	21	4	18	15	15	19	5	132
Incomplete	4	5	6	5	3	5	3	5	7	2	45
No Info										1	1
Total	14	13	24	26	7	23	18	20	26	9	180

3.2 Reproduction Results

Table 7 presents the main results of the reproduction exercise. 69 of 180 replication attempts successfully replicated the article’s analysis (38.3%). The success rate of replication conditional on non-confidential data was 42.6% (69 out of 162). A further 68 (42%) were at least partially successful. This means that the replicators were able to execute the computer programs that produced the numerical values reported in the articles, and that any differences in the numerical values were negligible.¹³ The main reason for unsuccessful reproductions is that the data used in the article was either confidential or proprietary, and therefore not available to replicators. Normally, this would imply that no replication attempt is undertaken. In these cases, it would seem that assessors had missed the information that no data was available. Combining with the 80 articles identified as relying on confidential data at the initial assessment stage, a total of 98 out of 303 (32.3%) relied on confidential or proprietary data, and were thus not reproducible by this project. The fact that these articles were not immediately recognized by the assessor is itself an argument for better metadata on journal websites.

¹²We are currently working on resolving two issues with the sample. First, the completion response questionnaire captured information for some articles that were not recorded in the initial assessment questionnaire. This is explained by a subsample of 2013 articles that were incorrectly recorded, which explains the corresponding low 2013 number in Table 2. Second, it appears as if not all of the 209 eligible articles were attempted. These will be assigned to replicators shortly.

¹³Possible causes may lie in software version discrepancies, uninitialized random number generators, different operating systems, or even different machines. We did not identify the causes of the discrepancy.

Table 7: Reproduction Results

Year	Confidential Data	Unsuccessful	Successful	Partial	Total
2009	5	1	4	4	14
2010		3	7	3	13
2011		12	10	2	24
2012	2	3	8	13	26
2013			4	3	7
2014	2	2	7	12	23
2015	2		4	12	18
2016	3	1	9	7	20
2017	2	2	13	9	26
2018	2	1	3	3	9
Total	18	25	69	68	180

Table 7 lists a further 25 articles that were not able to be reproduced for other reasons. Table 8 breaks down the reasons for these unsuccessful reproductions. 4 articles did not provide the (non-confidential) data required to produce their results. Further investigation is needed to identify the reason for this apparent non-conformance to the [AEJ:AE](#) data availability policies, although we point out that some non-confidential data is still subject to terms of use that prevent redistribution (earlier years of IPUMS data and any version of PSID data are just two examples). Errors in the provided computer programs prevented the replication of 1 article, while the data provided in 3 articles was corrupted in some way so that the software available to us was not able to read the datasets. Our replicators did not have access to the software required to run 1 article. For 16 articles, the computer programs successfully ran, but the numerical values were inconsistent with those reported in the articles, and the replicators were unable to find a convincing reason.

Table 8: Reason for Unsuccessful Reproduction

Year	Missing Data	Corrupted Data	Code Error	Software Unavailable	Other	Total
2009		1				1
2010					3	3
2011					12	12
2012	1	1	1			3
2014	1			1		2
2016	1					1
2017	1				1	2
2018		1				1
Total	4	3	1	1	16	25

Turning our attention to successful replications, we tabulate in Table 9 the extent to which modifications to the provided computer programs were required to successfully reproduce the articles. The majority of successful replications required minimal work from the replicators. 45 of the 69 successful replications required, at most, a simple rerouting of directory references. The remaining 24 successful articles required a deeper understanding of the software, and a more in-depth analysis of the code and/or command of the subject matter. These “Complex Changes” to the code required more than simple directory adjustments such as, for example, the debugging of classical code errors or the adjustment of outdated commands to reflect newer versions of software or operating systems.

Table 9: Manipulation of Code Required for Successful Reproductions

Year	Complex Change	Directory Change	No Change	Total
2009	1	1	2	4
2010	1	3	3	7
2011	4	2	4	10
2012	2	3	3	8
2013		1	3	4
2014	2	1	4	7
2015		3	1	4
2016	3	3	3	9
2017	9	2	2	13
2018	2		1	3
Total	24	19	26	69

Good documentation is key to better reproducibility, as emphasized by many authors (McCullough2006; ChangLi2015; Stark2018). In Table 10 we investigate whether better documentation is positively correlated with reproduction success. The results show a positive and statistically significant relationship between reproduction success of an article and the quality of its documentation.

Table 10: OLS: Reproduction Success vs Documentation Quality

	Successful Reproduction
Documentation Clarity = Complete	0.300*** (0.080)
Constant	0.200** (0.070)
<i>N</i>	180

Notes: Dependent variable = 1 if fully replicated.

3.3 Complementary Data

We captured bibliometric measures in 2017 for articles published through 2013, leaving a minimum of 3 years of post-publication years available to measure these metrics. We captured h-index for all authors of a paper, and computed the average per-paper h-index, as well as the lowest and highest when multiple authors were present. We also computed the average annual citations of the paper. Table 11 presents a summary of these measures, categorized by reproduction success. Articles have an average of 2.2 authors and were cited, on average, 4.6 times per year.

Table 11: Publication and Author Metrics

Outcome	Number of Articles	Avg h-index	Lowest h-index	Number of Authors	Avg Annual Citations
Unsuccessful	26	7.1	4.6	2	4.9
Partial	25	7	4.7	2.1	4
Successful	33	7.4	4.3	2.6	4.8

We investigate the relationship between the bibliometric measures and reproducibility measures

and outcomes. We model the count of citations, conditional on h-index measures, reproduction outcome of the paper, the type of data used, and other covariates. In Table 12, we start by controlling for the h-index measures, interacted with an indicator whether the article used confidential data. Results indicate a positive but noisy citation bonus for papers with confidential data. Authors with a high h-index, an indicator of high citation count in the past, also seem to obtain more future citations, but there is no interaction with the use of confidential data.

Table 12: OLS: Citations vs Confidential Data

	Annual Citations		
	(1)	(2)	(3)
avghindex	3.000*** (0.700)		
tophindex		1.000*** (0.400)	
lowhindex			2.000 (1.000)
confidential_data	17.000* (9.000)	13.000* (8.000)	8.000 (9.000)
avghindex:confidential_data	-1.000 (1.000)		
tophindex:confidential_data		-0.400 (0.700)	
lowhindex:confidential_data			-0.200 (2.000)
Constant	5.000 (6.000)	13.000*** (5.000)	20.000*** (5.000)
<i>N</i>	118	118	118

Notes:

*** Significant at the 1 percent level.

** Significant at the 5 percent level.

* Significant at the 10 percent level.

Table 13: OLS: Citations vs Reproduction Success

	Annual Citations					
	(1)	(2)	(3)	(4)	(5)	(6)
avghindex	4.000*** (0.800)			4.000** (2.000)		
tophindex		2.000*** (0.400)			2.000* (0.800)	
lowhindex			2.000* (1.000)			2.000 (2.000)
‘Fully reproduced’	14.000 (10.000)	11.000 (9.000)	9.000 (10.000)			
avghindex:‘Fully reproduced’	-2.000 (1.000)					
tophindex:‘Fully reproduced’		-0.900 (0.700)				
lowhindex:‘Fully reproduced’			-1.000 (2.000)			
‘Full or Partial’				6.000 (15.000)	0.900 (11.000)	-3.000 (12.000)
avghindex:‘Full or Partial’				-1.000 (2.000)		
tophindex:‘Full or Partial’					-0.200 (0.900)	
lowhindex:‘Full or Partial’						0.040 (2.000)
Constant	-0.300 (7.000)	9.000 (6.000)	15.000** (7.000)	-0.100 (14.000)	12.000 (10.000)	22.000** (11.000)
Observations	77	77	77	77	77	77

***Significant at the 1 percent level.

**Significant at the 5 percent level.

Notes:

Conditional on not using confidential data, how does the reproducibility of an article affect its future citation count? Table 14 shows that, controlling for h-indices, the ability to reproduce an article does not appear to play a significant role in the citation count. In columns (1) through (3), we control only for full reproduction, in columns (4) through (6), for full or partial reproduction. The only correlate with a strongly significant effect appears to be the authors' reputation as captured by h-index.

Table 14: OLS: Log Citations vs Reproduction Success

	Annual Citations					
	(1)	(2)	(3)	(4)	(5)	(6)
avghindex	0.200*** (0.030)			0.200*** (0.060)		
tophindex		0.070*** (0.020)			0.070** (0.030)	
lowhindex			0.100** (0.050)			0.090 (0.080)
‘Fully reproduced’	0.700* (0.400)	0.500 (0.300)	0.500 (0.400)			
avghindex:‘Fully reproduced’	−0.070 (0.050)					
tophindex:‘Fully reproduced’		−0.040 (0.030)				
lowhindex:‘Fully reproduced’			−0.070 (0.080)			
‘Full or Partial’				0.400 (0.500)	0.100 (0.400)	−0.100 (0.500)
avghindex:‘Full or Partial’				−0.060 (0.070)		
tophindex:‘Full or Partial’					−0.020 (0.030)	
lowhindex:‘Full or Partial’						0.002 (0.090)
Constant	2.000*** (0.200)	2.000*** (0.200)	2.000*** (0.300)	2.000*** (0.500)	2.000*** (0.400)	3.000*** (0.400)
Observations	77	77	77	77	77	77

*** Significant at the 1 percent level.

** Significant at the 5 percent level.

Notes:

4 Discussion and Conclusion

[To be Completed]

A Acronyms Used

AEA American Economic Association

AEJ:AE American Economic Journal: Applied Economics

AEJ:EP American Economic Journal: Economic Policy

AJPS American Journal of Political Science

DOI Digital Object Identifier

EJ Economic Journal

JASA Journal of the American Statistical Association

JEEA Journal of the European Economic Association

OS operating system

ReStat Review of Economics and Statistics

URL Uniform Record Locator

VCS version control system

B Assessment Questionnaire

10/8/2014

ReplicationDataQuestionnaire - Google Forms

ReplicationDataQuestionnaire

Please fill out the form to the best of your abilities.

* Required

1. Please enter your NetID

.....

2. **DOI ***

What is the DOI (not the URL!) of the article you are reviewing? (This was sent to you by email: simply copy it here)

.....

3. **TypeOfArticle ***

Does the article contain empirical work, simulations, or experimental work?
Mark only one oval.

☐

Yes

☐

No

After the last question in this section, stop filling out this form.

4. **OnlineAppendix ***

Does the article have an online Appendix?
Mark only one oval.

☐

Yes

Skip to question 5.

☐

No

Skip to question 7.

Information on online materials

5. **OnlineAppendixURL**

Enter the URL of the online Appendix

.....

6. **OnlineAppendixDOI**

Enter the DOI of the online Appendix (this is often the case if the journal provides a DOI to the article, and hosts the appendix)

.....

Online Data

https://docs.google.com/forms/d/1c6wfHmXgcad5unPtVWltml_rU8piOcvNQEqgmK_lu9w/edit

1/12

In the next few sections, we consider the following types of data: (i) input data are data as collected by the authors or another agency (examples: "CPS" or "my survey data" (ii) analysis data are the post-processed and clean data underlying specific regressions. Think of the basic workflow [input data] -> [preparation programs] -> [analysis data] -> [regression programs] -> results. If available, we will request information on up to three input datasets, and one analysis dataset.

7. OnlineData *

Does the online appendix link to one or more downloadable dataset?

Mark *only one oval*.

☐ Yes *Skip to question 9.*

☐ No

8. OnlineDataInside *

Does the article itself mention where to obtain the final analysis data? (for instance, because the data are confidential or proprietary, or because there is a public-use download site for the data)

Mark *only one oval*.

☐ Yes

☐ No *Skip to question 43.*

Information on online datasets

Please describe the first INPUT dataset.

9. DataRunClean *

Does the article and/or its appendices allow you to identify the data needed to start from scratch? (Original input datasets)

Mark *only one oval*.

☐ Yes

☐ No *Skip to question 43.*

Input dataset 1

10. OnlineDataDOI

Please enter the DOI of the downloadable dataset (notation: doi://)

.....

11. OnlineDataHandle

Please enter the Handle of the downloadable dataset (notation: hdl://)

.....

12. OnlineDataURL

Please enter the URL of the downloadable dataset. (this may duplicate one or the other of DOI or HDL, but is a more general way to describe it. notation: http://)

.....

13. DataAvailability

Are the data available without restriction (can be downloaded or requested by anybody without restriction)? [Answer DK (Don't know) if it is not clear from the article how users can access non-downloadable data]

Mark only one oval.

- ☐ Yes
- ☐ No
- ☐ DK

14. DataAvailabilityAccess

Do the data require users to apply for access, purchase, or otherwise sign agreements to access the data? (This should be mentioned in the Readme PDF or in the article) [Answer DK (Don't know) if it is not clear from the article how users can access non-downloadable data.]

Mark only one oval.

- ☐ Yes
- ☐ No
- ☐ DK

15. DataAvailabilityExclusive

Are the data accessible only to the authors? [Answer yes if the authors clearly state that the data are only available to them. Answer No if there is clear evidence that others can access to the data, albeit with restrictions. Answer DK if you can't figure it out from the article.]

Mark only one oval.

- ☐ Yes
- ☐ No
- ☐ DK

16. OtherNotes

Any notes for this dataset that was not covered by the questions above.

.....

.....

.....

.....

.....

17. Do you want to describe another dataset? *

Mark *only one oval*.

☐ Yes

☐ No *Skip to question 34.*

Input dataset 2

18. OnlineDataDOI2

Please enter the DOI of the downloadable dataset (notation: doi://)

.....

19. OnlineDataHandle2

Please enter the Handle of the downloadable dataset (notation: hdl://)

.....

20. OnlineDataURL2

Please enter the URL of the downloadable dataset. (this may duplicate one or the other of DOI or HDL, but is a more general way to describe it. notation: http://)

.....

21. DataAvailability2

Are the data available without restriction (can be downloaded or requested by anybody without restriction)? [Answer DK (Don't know) if it is not clear from the article how users can access non-downloadable data]

Mark *only one oval*.

☐ Yes

☐ No

☐ DK

22. DataAvailabilityAccess2

Do the data require users to apply for access, purchase, or otherwise sign agreements to access the data? (This should be mentioned in the Readme PDF or in the article) [Answer DK (Don't know) if it is not clear from the article how users can access non-downloadable data.]

Mark *only one oval*.

☐ Yes

☐ No

☐ DK

23. DataAvailabilityExclusive2

Are the data accessible only to the authors? [Answer yes if the authors clearly state that the data are only available to them. Answer No if there is clear evidence that others can access to the data, albeit with restrictions. Answer DK if you can't figure it out from the article.]

Mark *only one oval*.

☐ Yes

☐ No

☐ DK

24. OtherNotes2

Any notes for this dataset that was not covered by the questions above.

.....

.....

.....

.....

.....

25. Do you want to describe another dataset?

Mark *only one oval*.

☐ Yes

☐ No *Skip to question 34.*

Input dataset 3

26. OnlineDataDOI3

Please enter the DOI of the downloadable dataset (notation: doi://)

.....

27. OnlineDataHandle3

Please enter the Handle of the downloadable dataset (notation: hdl://)

.....

28. OnlineDataURL3

Please enter the URL of the downloadable dataset. (this may duplicate one or the other of DOI or HDL, but is a more general way to describe it. notation: http://)

.....

29. DataAvailability3

Are the data available without restriction (can be downloaded or requested by anybody without restriction)? [Answer DK (Don't know) if it is not clear from the article how users can access non-downloadable data]

Mark only one oval.

☐ Yes

☐ No

☐ DK

30. DataAvailabilityAccess3

Do the data require users to apply for access, purchase, or otherwise sign agreements to access the data? (This should be mentioned in the Readme PDF or in the article) [Answer DK (Don't know) if it is not clear from the article how users can access non-downloadable data.]

Mark only one oval.

☐ Yes

☐ No

☐ DK

31. DataAvailabilityExclusive3

Are the data accessible only to the authors? [Answer yes if the authors clearly state that the data are only available to them. Answer No if there is clear evidence that others can access to the data, albeit with restrictions. Answer DK if you can't figure it out from the article.]

Mark only one oval.

☐ Yes

☐ No

☐ DK

32. OtherNotes3

Any notes for this dataset that was not covered by the questions above.

.....

.....

.....

.....

.....

Analysis datasets

33. DataRunFinal

Does the article and/or its appendices allow you to identify the data needed to run the final models?

Mark only one oval.

☐

Yes

☐No *Skip to question 41.***Analysis dataset****34. OnlineFinalDataDOI**

Please enter the DOI of the downloadable dataset (notation: doi://)

.....

35. OnlineFinalDataHandle

Please enter the Handle of the downloadable dataset (notation: hdl://)

.....

36. OnlineFinalDataURL

Please enter the URL of the downloadable dataset. (this may duplicate one or the other of DOI or HDL, but is a more general way to describe it. notation: http://)

.....

37. FinalDataAvailability

Are the data available without restriction (can be downloaded or requested by anybody without restriction)? [Answer DK (Don't know) if it is not clear from the article how users can access non-downloadable data]

Mark only one oval.

☐

Yes

☐

No

☐

DK

38. FinalDataAvailabilityAccess

Do the data require users to apply for access, purchase, or otherwise sign agreements to access the data? (This should be mentioned in the Readme PDF or in the article) [Answer DK (Don't know) if it is not clear from the article how users can access non-downloadable data.]

Mark only one oval.

☐

Yes

☐

No

☐

DK

39. FinalDataAvailabilityExclusive

Are the data accessible only to the authors? [Answer yes if the authors clearly state that the data are only available to them. Answer No if there is clear evidence that others can access to the data, albeit with restrictions. Answer DK if you can't figure it out from the article.]

Mark only one oval.

☐ Yes

☐ No

☐ DK

40. OtherNotesFinal

Any notes for this dataset that was not covered by the questions above.

.....

.....

.....

.....

.....

Data formats

Considering all of the datasets described above, please check all boxes that apply.

41. DataFormatInputs

What format are the original input datasets in?

Check all that apply.

☐ Stata

☐ CSV

☐ R

☐ Matlab

☐ SPSS

☐ SAS

☐ Other:

42. DataFormatAnalysis

What format are the final analysis datasets in?

Check all that apply.

- ☐ Stata
- ☐ CSV
- ☐ R
- ☐ Matlab
- ☐ SPSS
- ☐ SAS
- ☐ Other:

Information on programs**43. OnlinePrograms**

Does the online appendix have information on the programs used to run the analysis?

Mark only one oval.

- ☐ Yes *Skip to question 45.*
- ☐ No

44. OnlineProgramsInside

Does the article itself mention where to obtain the programs needed to replicate the study (for instance, at a data or code repository)?

Mark only one oval.

- ☐ Yes
- ☐ No

Information on online programs**45. OnlineProgramsDOI**

Please enter the DOI of the downloadable programs (notation: doi://

.....

46. OnlineProgramsHDL

Please enter the Handle of the downloadable programs (notation: hdl://

.....

47. OnlineProgramsURL

Please enter the URL of the downloadable programs (notation: http://

.....

Documentation

48. DocReadmePresent

Does the downloadable data/program archive include a Readme PDF or TXT file, or some other generic instruction file?

Mark only one oval.

☐ Yes

☐ No

49. DocReadmeContent

Does the Readme PDF (or generic instruction file) list all included files, document the purpose and format of each file provided, and provides instruction to a user on how replication can be conducted?

Check all that apply.

- ☐ lists all included files
- ☐ documents the purpose of each file
- ☐ documents the format of each file
- ☐ provides instructions for replication

Program details

50. ProgramFormat

What format are the programs in?

Check all that apply.

- ☐ Stata
- ☐ R
- ☐ Matlab
- ☐ SPSS
- ☐ SAS
- ☐ Other:

51. ProgramSequence

Does the Readme PDF or one of the other included documents (including one of the programs) provide enough detail to run all the programs?

Mark only one oval.

☐ Yes

☐ No

52. ProgramsDocumentation

Are the programs themselves clearly documented? (There are comments throughout the program that briefly describe what is done at each step)

Mark only one oval.

☐ Yes

☐ No

53. ProgramsHeaderAuthor

Do the programs have a header that identifies the author? (Program metadata)

Mark only one oval.

☐ Yes

☐ No

54. ProgramsHeaderInfo

Do the programs have a header that identifies when they were created and/or modified (Program metadata)

Mark only one oval.

☐ Yes

☐ No

55. ProgramsStructureManual

Do the instructions require the user to do manual modifications to data or programs?

Mark only one oval.

☐ Yes

☐ No

56. GeneralNotes

General notes on this article, that wasn't captured by the questions

.....

.....

.....

.....

.....

10/8/2014

ReplicationDataQuestionnaire - Google Forms

57. **How difficult do you think replicating the article will be? ***

Mark only one oval.

	1	2	3	4	5	
easiest	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	hardest

Powered by



C Exit Questionnaire

10/8/2014

Exit_Questionnaire_Draft - Google Forms

Exit_Questionnaire_Draft

Please fill out the form to the best of your abilities.

* Required

1. Please enter your NetID

.....

2. DOI *

What is the DOI (not the URL!) of the article you reviewed?

.....

3. Code_Success *

Did you manage to eventually get all the programs to run successfully?
Mark only one oval.

☐

Yes

☐

No *Skip to question 11.*

Original Program

4. Program_Run_Clean *

Did the programs run "as is" without needing to make ANY changes?
Mark only one oval.

☐

Yes, no changes were necessary. *Skip to question 7.*

☐

No, I needed to make changes in the code.

Changes to Program Code

5. Directory_Change

Were the changes restricted to simply redirecting file/folder paths?
Mark only one oval.

☐

Yes

☐

No, the changes to the code were more involved.

6. Code_Changes

If the changes were more involved, briefly describe what changes you had to make.

.....

.....

.....

.....

.....

Program vs Paper Discrepancies**7. Output_Accuracy**

Do the numbers produced by your program exactly match their corresponding values in the paper?

Mark only one oval.

☐

Yes

☐

No, some of the numbers are different.

8. Discrepancy_Location

If there are values that do not match, please list their location (ie. table number, column, page).

.....

.....

.....

.....

.....

Skip to question 11.

Reason for replication failure.**9. Which of the following apply?**

Check all that apply.

☐

Missing data set

☐

Error in the program code

☐

Other:

10. Briefly describe the reason why you could not replicate.

.....

.....

.....

.....

.....

Software Issues

11. **Software_Extensions**

Did you have to load any software extensions? (Eg. In matlab, the optimization toolkit is required to run the fmincon command. In Stata, outreg2 needs to be installed before running the command.)

Mark only one oval.

- ☐ Yes
- ☐ No
- ☐ DK

12. **Software_Version**

Did the authors use a different version of software (ie. Stata11 instead of Stata13)?

Mark only one oval.

- ☐ Yes
- ☐ No
- ☐ DK

13. **First_Replicator**

Are you the first replicator?

Mark only one oval.

- ☐ Yes *Skip to question 17.*
- ☐ No

Previous Replicator Questions

14. **Common_Issues**

Did you encounter the same issues as the previous replicator?

Mark only one oval.

- ☐ Yes
- ☐ No

15. Overcome_Issues

Were you able to overcome any problems faced by the previous replicator?
Mark only one oval.

- ☐ Yes
- ☐ No
- ☐ N/A. The previous replicator had no issues.

16. Replication_Helpfulness

Describe the usefulness of the previous replicator's notes. Did you add to them?

.....

.....

.....

.....

.....

Original Author**17. How complete was the original author's readme/generic instruction file?**

Mark only one oval.

- ☐ Complete. Provided all information required to run the programs.
- ☐ Incomplete. Was ambiguous or left out crucial steps.
- ☐ No readme file was provided.

18. What actions could the authors have made to make the replication exercise easier? (Eg. correctly point to folder names)

.....

.....

.....

.....

.....

Overall Rating**19. How difficult do you think the replication exercise was? ***

Mark only one oval.

	1	2	3	4	5	
easiest	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	hardest

20. If this differs from the initial assessment, why?

21. **GeneralNotes**

General notes on this article/replication, that wasn't captured by the questions

D Replication Team

The following members of the Replication Lab provided valuable assistance:

Hautahi Kingi Alice Elaine Chou, Haeyong Shin, Yaxian Xie, Nathan Allan Bach, Cindy Vincens, Yuxin Chen, Flavio Stanchi, Sarah Jane Harrison, Yiwen Jiang, Jack Wendler, Jose Fernandez, Joran Isenberg, Sarah Harrison, Koonj Vekaria, Charley Chen, Yang Guo, Yiwen (Evelyn) Jiang, Noah Kwicklis, Madeline Kwicklis, Koonj Vekaria, Robin Wang, Jack Q. Wendler, Qianyan Yao, Joao Vitor Costa, Evan Shapiro, Yudi (Grace) Wang, Christopher Chang, Chuhan Liu, Daniel Kim, Cassandra Madulka, Robert Goldberg, Xinyi Wan, Siming Zou, Yu Gao, Andrew Wink, Matthew Salazar, Naomi Li, Anderson Park, Carina Chien, Nick Swan, Vendela Norman, Hayley A. Timmons, Jack VanSlyke, Gabriel Bond, Wenxin (Andee) Cao, Mcrid Wang, John Park, Xueshi Su, Sam Mbugua, Jiazhen Tan.