


ORIGINAL ARTICLE

The influence of hidden researcher decisions in applied microeconomics

Nick Huntington-Klein, Assistant Professor¹  | Andreu Arenas, Assistant Professor² | Emily Beam, Assistant Professor³ | Marco Bertoni, Associate Professor⁴ | Jeffrey R. Bloem, Research Economist⁵ | Pralhad Burli, Economist⁶ | Naibin Chen, Graduate Student⁷ | Paul Grieco, Associate Professor⁸ | Godwin Ekpe, Graduate Student⁹ | Todd Pugatch, Associate Professor¹⁰ | Martin Saavedra, Associate Professor¹¹ | Yaniv Stopnitzky, Assistant Professor¹²

¹Seattle University, Seattle, Washington, USA

²University of Barcelona & IEB, Barcelona, Spain

³University of Vermont, Burlington, Vermont, USA

⁴Department of Economics and Management "M. Fanno", Padova University, Padova, Italy

⁵USDA Economic Research Service, Kansas City, Missouri, USA

⁶Idaho National Laboratory, Idaho Falls, Idaho, USA

⁷Pennsylvania State University, 303 Kern Building, University Park, Pennsylvania, USA

⁸Pennsylvania State University, 508 Kern Graduate Building, University Park, Pennsylvania, USA

⁹Northern Illinois University, Dekalb, Illinois, USA

¹⁰School of Public Policy, Oregon State University, Corvallis, Oregon, USA

¹¹Department of Economics, Oberlin College, Oberlin, Ohio, USA

¹²University of San Francisco, San Francisco, California, USA

Correspondence

Nick Huntington-Klein, Seattle University, 901 12th Ave., Seattle, WA 98122, USA.
Email: nhuntington-klein@seattleu.edu

Abstract

Researchers make hundreds of decisions about data collection, preparation, and analysis in their research. We use a many-analysts approach to measure the extent and impact of these decisions. Two published causal empirical results are replicated by seven replicators each. We find large differences in data preparation and analysis decisions, many of which would not likely be reported in a publication. No two replicators reported the same sample size. Statistical significance varied across replications, and for one of the studies the effect's sign varied as well. The standard deviation of estimates across replications was 3–4 times the mean reported standard error.

KEYWORDS

metascience, replication, research

JEL CLASSIFICATION

C81; C10; B41

Abbreviations: CPS, Current Population Survey; IPUMS, Integrated Public Use Microdata Series; NBER, National Bureau of Economic Research; RDD, regression discontinuity design.

1 | INTRODUCTION

A primary goal of empirical work in the social sciences is to generate results that are internally valid, generalizable, and replicable. However, estimated results for the same research question can vary considerably from study to study. Variation in measures of economic effects is not in itself surprising or concerning. Results will naturally differ due to the use of different samples, empirical methods, or, in the presence of heterogeneous effects, different settings. For example, Chetty et al. (2016) performed a reevaluation of the Moving to Opportunity experiment and found positive neighborhood effects, where previous studies had not found those effects, by changing the research design to focus on child's age at the time they moved.

Reviews of the literature often attempt to rationalize these results on the basis of focal factors like setting or analytic method. However, variation may also come from the hundreds of decisions made in the process of analysis, from data cleaning to variable definition, even if the empirical approach is held constant. In one notable example, White et al. (2018) show that the choice of whether to use imputation, and the choice of imputation method, has a dramatic effect on measures of productivity dispersion computed from the US Census of manufacturers. This is no small matter, as the degree of dispersion has important implications for the degree of misallocation and overall efficiency of the US economy in apportioning resources across firms.

In another example, Clemens and Hunt (2019) look at differences between studies in the long and contentious literature on the effects of immigration from the Mariel boatlift on wages in Florida. They find that some of the differences in results could be explained by the use in some studies of a subsample that made the analysis sensitive to racial composition changes. The use of that subsample was not explicitly justified in the original studies.

The wide range of decisions behind every analysis can be thought of as “researcher degrees of freedom.” If these researcher degrees of freedom add significant variation to the research process, it will become difficult to meaningfully compare similar studies and properly synthesize a scientific consensus.

Unreported researcher variation is to some extent inevitable. Papers are limited in length, so only some fraction of researchers' choices can be described explicitly. Even if full replication code is available, peer reviewers may not be capable of judging all of these choices, readers may not be aware how these choices are made and so be able to incorporate them into their understanding of the results, and revisiting these choices after publication occurs rarely, if it is even possible.

To be clear, the choices we have in mind are necessarily made at some point in the process even if analysis does not change after viewing the results, so the issue is distinct from “p-hacking” or “The Garden of Forking Paths” (Gelman & Loken, 2014; Silberzahn et al., 2018), as well as being distinct from publication bias (De Long & Lang, 1992), and would not be solved by the use of preregistered analysis, even if the preregistration recorded all the details of analysis. These choices are not necessarily errors, and so even someone reviewing full replication code may not even have reason to reconsider a potentially consequential decision. Still, without a sense of how much noise these well-intentioned but unreported choices can introduce into empirical research, the literature must face a crisis of confidence when evaluating its own results.

In this paper, we attempt to measure the magnitude of variation due to researcher degrees of freedom in the context of applied microeconomics studies that attempt to isolate a causal effect. By doing so, we aim to provide a better understanding of whether this issue is of serious concern in the empirical economics literature. We take a “many-analysts” approach where multiple researchers use the same data set to answer the same research question, without knowledge of the methods used by other analysts, and without publication contingent on their results. This allows us to examine how different the choices made in data manipulation and analysis are between good-faith researchers, and to examine the impact of these choices on the eventual results.

Previous many-analysts studies include Silberzahn et al. (2018), which recruited psychological researchers to examine whether a data set of referee calls in soccer showed evidence of discrimination against darker-skinned players. Analysts differed in the choice of linear or nonlinear regression, the treatment of nonindependent error terms, selection of covariates, and regression method, leading to odds-ratio effects estimates that varied between .89 and 2.93 even after each team's methods underwent a peer-review process. Botvinik-Nezer et al. (2020) recruited teams to examine a data set of fMRI imaging data on subjects playing a gambling task to evaluate nine *ex ante* hypotheses about the impact of gain/loss framing on risk-taking. Replication teams differed in their use of image smoothing, statistical mapping of activated brain regions, and generally in the pipeline from data through analysis. Again, results differed widely across teams: tests for one of the hypotheses were significant for 84.3% of replication teams, and for three others about 6% of

the time. The remaining five hypotheses had more variation and were significant for between 21.4% and 37.1% of teams. Both of these are preceded in publication outside of academia by Cohn (2016), who had five pollsters evaluate the same poll for the 2016 US presidential election, producing an estimate of the percentage planning to vote for each of the candidates, with variation in, for example, demographic and nonresponse weighting choices again leading to widely varying results, not even agreeing which candidate had come out on top.

Taking another angle in which the hypothesis is kept constant but research design and data are both varied, Landy et al. (2020) have multiple independent research teams design psychology experiments intended to test the same five hypotheses about moral judgments, negotiations, and implicit cognition. Participants in a large sample were assigned to participate in the different experiments. For four out of five of the hypotheses, different teams found effects of opposite signs, with the narrowest range of Cohen's d estimates between -0.37 and $+0.26$. They found that variability was not related to the skill of the research team.

Unlike previous many-analysts studies, we allow researcher freedom in the construction and cleaning of the observational data set. The processing and cleaning of administrative, governmental, or otherwise externally generated data is a common feature of applied microeconomic research, and is a likely source of researcher degrees of freedom.

We use two studies published in high-quality journals as a basis, and produce seven replications of each study. We find considerable variation both in results and in the construction of data. No two replicators ended up with the same sample size, and in several cases large differences in sample construction were driven by decisions that would likely go unmentioned in an eventual publication, or at least be overlooked by reviewers and readers. Analysis decisions also showed large differences. Most major analytic differences, like the use of linear probability models vs. logit, would have likely been mentioned in publication for reviewer scrutiny, but others, like the construction of bins when generating a control for education, would likely be overlooked by a reviewer.

The actual effect on results was mixed. In one of the studies, six of the seven replications had very similar point estimates and overlapping confidence intervals. In the other, results varied much more widely, with both significant positive and negative coefficients.

The variation across replications implies that there is a fair amount of uncertainty in published results that is not represented in estimates of sampling variation. Further, much of this variation comes from areas like data cleaning that are not standardized, and for which many of the decisions that researchers make may be difficult to see and evaluate. Some methods to alleviate this issue are discussed in the conclusion.

2 | RESEARCH QUALITY IN THE SOCIAL SCIENCES

This study is a part of a modern metascientific literature in the social sciences with a particular concern for the quality of results. Such metascientific studies in economics sometimes examine published work either through the lens of its unconsidered statistical properties (Ioannidis et al., 2017; Young, 2018). More often, these checks operate through attempts to replicate published papers, either focused on individual papers, or on many at once (Camerer et al., 2016; Chang & Li, 2017; Dewald et al., 1986).

Replication studies allow us to determine variation in results that arises from the use of different samples, when testing the same hypothesis in new data, or from major analysis choices, when testing the same hypothesis in the same data with a new method (Christensen & Miguel, 2018; Clemens, 2017; Hamermesh, 2007). These replication studies can also reveal variation in results that arise from errors in code, when attempting a “pure” replication (Hamermesh, 2007) to reproduce the original tables and figures. The most well-known example of the latter is likely the attempt by Herndon et al. (2014) to replicate Reinhart and Rogoff (2010), which uncovered a coding error in the original study.

“Pure” replication (or “reproduction”) studies in economics, in which a new study is performed purely to check the results of a prior study using the same data and methods, are relatively rare (Berry et al., 2017; Hamermesh, 2017), and the incentives for performing them are not well-aligned (Duvendack et al., 2017; Gertler et al., 2018).

Pure replication studies have at least become easier to perform as more economics journals have added requirements to include analysis code and data, many following the American Economic Review's additional requirements implemented in 2003, although by 2016, more than half of AER papers that used data were exempted from fully sharing it (Christensen & Miguel, 2018). Christensen and Miguel (2018) also have a review of data-sharing and code-sharing policies among major economics journals, finding that data-sharing policies and willingness to publish replication work has increased but are not universal.

Christensen and Miguel (2018) also review several attempts to reproduce studies in economics, with Chang and Li (2017) in particular attempting pure reproductions in macroeconomics. Christensen and Miguel (2018) find that data are often unavailable, even at journals with data-sharing policies and following requests to authors. Even with data in hand, a significant portion of the studies analyzed could not be reproduced. This implies data cleaning or analysis decisions that are either in error, or are not sufficiently described that a reader can evaluate and reproduce all of the choices being made. Beyond the attempts of these authors to reproduce the work, the replication files available on journal websites are generally hundreds of lines of code long, containing similarly many relevant decisions. A reader of the paper is necessarily looking at a short summary of that code, which cannot possibly fully describe all of the choices being made, no matter the author's efforts.

Whether or not a given replication attempt is successful, these kinds of pure replication will also have difficulty in uncovering variation in results that is driven by researcher degrees of freedom. Pure replications will intentionally make the exact same choices as in the original study unless a clear error is spotted, and replications using new data or methods generally attempt to make the same choices except for the specific data set or method being changed, so as to isolate the source of any difference. Incentives to replicate, which favor results that overturn an original study (Dewald et al., 1986; Gertler et al., 2018; Hamermesh, 2007), do not favor looking into these choices. Even if results are found to be sensitive to researcher degrees of freedom, as long as the original choices are not obviously incorrect, it is difficult to make a convincing case to an editor that the results have been overturned.

Researcher degrees of freedom, however, may have significant impact on the results of a study even if the choices made are not *wrong*, but are simply one reasonable option of many. In psychology, Simmons et al. (2011) find that researcher flexibility among accepted options in experimental design and data analysis allow nearly any hypothesis to be supported. Lenz and Sahn (2017) find that, in a major political science journal, 30–40% of significant results were insignificant bivariate relationships that became significant only with the addition of controls. As the authors describe, significant relationships that emerge only in the presence of controls are considered highly suspect in the field unless there is strong theoretical justification for them, but none of the articles offer theoretical justification, and only one published the bivariate relationship that would allow a reader to see that controls were necessary for significance. In economics, McCullough and Vinod (2003) found that the choice of software package significantly affected results from nonlinear optimization, and that none of the papers they reviewed from the *American Economic Review* tested their nonlinear optimizations in multiple software packages.

All three of these studies refer to choices made “behind the scenes,” comparing decisions in the published paper to analyses *not* run or reported, which produce different results.¹ The fact that these papers were published implies that referees and editors found the chosen methods at least reasonable—these studies find that other reasonable and acceptable choices produce different results, showing that researcher degrees of freedom matter. For a reader to evaluate these decisions would require not just access to replication files and data to evaluate, but a willingness to try different versions of the code and analysis to determine whether and how they matter. This implies a threat to the validity of results that is somewhat different from what either pure replication or a careful evaluation of a paper's methods can do.

3 | METHODS

The methods for this study include (i) selecting papers, and analyses within those papers, to replicate, (ii) designing instructions and information to present to replicators, (iii) recruiting replicators, and (iv) evaluating replicator work.

3.1 | Selecting replication tasks

Project organizers developed a list of desirable attributes for studies to replicate. These included:

- Studies should be published in well-regarded economics journals,² with a preference for publication in the last 20 years.
- Studies should contain a single causal estimate of interest that can be replicated.
- Studies should not be so well-known that replicators are likely to recognize them from the instructions.³
- Studies must use publicly-available data, and ideally data that microeconomists would be used to working with. Because organizers anticipate that most replicators will be American, public American data sets are favored.

- If studies rely on highly specific domain knowledge or obscure methods that replicators would not know on their own, it should be simple enough to explain to replicators in instructions, or be secondary to analysis so it can be removed for a simplified replication.
- The two studies selected should be from different subfields of applied microeconomics.

Because the easiest of these criteria to use in a literature search is that the studies use publicly-available data, organizers used Google Scholar to search for the names of publicly available data sets commonly used in applied microeconomics research, including, generically, “Census”, as well as the National Longitudinal Survey of Youth, the National Education Longitudinal Study, the Panel Study of Income Dynamics, the American Community Survey, and the Current Population Survey (CPS). We did not require studies to use one of these data sets.

In search results, studies that were published in top economics journals, and seemed likely to satisfy the other criteria based on the title and abstract, were examined more closely.

We evaluated the abstracts of a large number of studies with a goal of finding 50 for closer examination. From these 50, two were found to be best satisfy all criteria while also being feasible to replicate: Black et al. (2008) and Fairlie et al. (2011).

Given these two studies, we isolate and simplify the analyses to be given to replicators. More detailed instructions are in the next section.

Black et al. (2008) is a study of the effect of compulsory schooling on teenage pregnancy. The authors use variation in compulsory schooling policy in two environments—the United States and Norway—and estimate the effect of those changes on teenage pregnancy rates. They then attempt to distinguish whether the effect operates by improving human capital or through the “incarceration effect.”

Because the goal of this study is to test for variation between replications rather than to test the robustness of the original results, we base replication instructions on a simplified version of the analysis, making the original studies not directly comparable to the replications. We focus on their primary analysis of the United States, which uses US Census data from 1940 to 1980, calculates women’s age at first birth, and excludes apparent births age 14 or below. State- and decade-level variation in compulsory schooling laws identifies the effect of compulsory schooling. Instructions are based on a replication of the top half of Black et al. (2008), Table 2, Column 3, where “birth by age 18” is the outcome variable. We simplify their analysis by looking at only one compulsory schooling margin—whether the state has a compulsory schooling age of 16 or higher, as opposed to 15 or lower.

Fairlie et al. (2011) is a study of the effect of employer-based health insurance on entrepreneurship. The authors use variation in age to identify the effect. Men aged 65 or older qualify for Medicare, and those aged 64 and 11 months or younger generally do not. Medicare reduces the need for employer-based health insurance, and so authors look for a jump in entrepreneurship at age 65 exactly. CPS data are used to identify men who turn 65 within one of the 4-month runs where they are included in the survey.

Instructions are based on one of their analyses, which is shown in Fairlie et al. (2011, table 6). Men who can be observed having just turned 65 are compared to those observed just under 65 in terms of the rates of self-employment, conditional on being employed at all. We simplify the task for replicators by narrowing the sample window from 1996–2006, as in Fairlie et al. (2011), to May 2004–December 2006. This avoids combining data across samples where variable definitions have changed.⁴

3.2 | Replication instructions

We construct sets of instructions for each replication. The goal of these instructions is to ensure that each replicator knows what the data set and research question of interest are, as well as some identifying assumptions, without restraining their choices too much. Replicators were encouraged to perform each analysis as if they were writing their own paper for publication. The full text of each set of instructions is in the online supporting information.

For both sets of instructions, replicators are told to use any statistics package, and that they should use assistants if they would normally use assistants in their work. They are also told that their analysis should be independent, and should not attempt to identify the original study, or to match (or mismatch) with fellow replicators. The goal is to uncover “how you would estimate this effect, if you’d had this question, this idea for identification, and had chosen this particular sample.”

They are also told to focus on a single “headline” result, of the kind that might be reported in an abstract. Replicators were not directed to perform robustness checks and alternate analyses.

Instructions for the Black et al. (2008) replication direct replicators to download Census 1% files from the Integrated Public Use Microdata Series (IPUMS) for 1940–1970, and the 5% files from 1980. IPUMS provides data files in already mostly-usable format, and the different census years are already appended together (Ruggles et al., 2020). Data should then be limited to female subjects aged 20–30.

Replicators are given the background theory that compulsory schooling laws may reduce the incidence of teenage pregnancy for a number of reasons, and told to estimate the effect of compulsory education age in a state on the proportion of women in that state who have a teenage pregnancy, under the identifying assumption that trends in teen pregnancy are unrelated to the decision to change compulsory schooling policy.

Replicators are given the definition of a teenage pregnancy as “having a child by age 18,” and told to determine the compulsory schooling law being applied as the law in place in the mother’s birth state when they are 14 years of age. This removes several researcher degrees of freedom in deciding the appropriate margin for analysis, but ensures that the estimates will be comparable across replications. A table of compulsory schooling laws by state and decade, from Black et al. (2008), is given to replicators in Word format, and for women who turned 14 between policy years, they are told to use the most-recent policy. Replicators are also told to look specifically at the margin of compulsory schooling at age 16, comparing policies requiring students to stay until they are 16+ against policies requiring some age 15 or below.

Instructions for Fairlie et al. (2011) direct replicators to download CPS monthly files from the National Bureau of Economic Research (NBER) for the months of May 2004 through December 2006. The use of NBER individual monthly files, rather than pre-compiled and combined files from, for example, IPUMS, means that replicators will have to import the raw files and combine them into a single data set, a data-cleaning task in which there may be different researcher decisions made.

Data should then be limited to male subjects who can be observed “in the exact month that they turn 65,” meaning subjects observed both at the ages 64 and 65 in one of the 4-month CPS rolling panels they are present in.

Replicators are given the background theory that employer-provided health insurance may be a barrier to entrepreneurship, and given the background information that Medicare eligibility occurs at exactly 65 years of age for most people. It asks for the effect of Medicare eligibility on the rate of self-employment, conditional on being employed at all. They are given the shared identifying assumption that nothing else of importance changes between the ages of 64- and 11-months and 65.

3.3 | Replicator recruitment

Replicator recruitment began in May 2018. There were two main methods for requesting participation from replicators and directing them to the sign-up website.⁵ In both cases, to improve recruitment success, the recruitment message stressed that the replication project was designed so that it would only take a moderate amount of time, and that successful replicators would be offered coauthorship or acknowledgement.

First, we used the U.S. News and World Report ranking of economics departments to develop a list of 138 top economics departments. We sent an email to the chair of each department, asking them to forward on a recruitment message to their faculty, or to only the applied microeconomists. We do not know how many department chairs complied with this request.

Second, we posted a message on Twitter asking for interested researchers to sign up. The link from the tweet to the recruitment website was clicked 638 times.

On the recruitment form, replicators were asked whether they had any published or forthcoming work in applied microeconomics, were familiar with standard causal inference methodology, whether they had performed replication work before, whether they typically used student assistants, whether they would want to complete one replication or two, and what their typical fields of interest were. The pool of replicators was intended to represent people actually producing applied microeconomics research, and so recruitment was limited to those with published or active work in the field.

In total, 51 researchers signed up to complete a replication, 49 of whom were considered qualified for the task, meaning they reported having at least one published or forthcoming work in applied microeconomics and being familiar with standard causal inference methodology. Of the 51, 37 came from Twitter and 14 from email.

Replication tasks were assigned to replicators first on the basis of field of interest. If replicators indicated that their primary fields of research were relevant to one of the replication tasks but not the other, they were assigned to that task. Replicators who listed topics relevant to neither or both tasks were randomly assigned. Replicators who agreed to do both replications were assigned their first task in the same way. Replicators were not given information from organizers about who else had been recruited, or results from any of the other replicators. The initial due date for replication was the end of January 2019, or 7 months after recruitment. This due date was eventually pushed back to March 2019. The first successful replication was completed May 21, 2018, and the final one was received March 31, 2019.

Of the original 49 qualified researchers, 12 finished a replication: 10 finished one replication each, and two finished two replications each, for a total of seven completed replications of each task. Four of the successful replicators had been recruited by email, and eight had been recruited from Twitter. Project organizers, who knew the content of the original studies, did not contribute any replications. In one case organizers provided assistance to a replicator who was having difficulty importing data files.

Upon completion, replicators were asked to complete an exit survey. When those who had signed up to provide a replication but did not complete one gave reasons for not finishing, they reported almost uniformly that they did not have the available time they had expected to work on the project, and their decision was unrelated to the content of the task. Replicators who did complete a replication gave their reasons for participation. Nine reported interest in replication or the importance of replication as a reason. Five mentioned that the request for participation happened to line up with an opening in their schedules. Four cited that they thought project would be fun or would help develop their skills.

Despite reported attrition being due to a lack of time, the attrition rate introduces the possibility of selection bias among the kinds of researchers who might actually finish the replication task. However, other than differential attrition between email and Twitter recruits, observables were unrelated to attrition. Attrition rates were unrelated to whether researchers had previously performed any replication work outside of classroom assignments (average among successful replicators .455, difference in attrition .003, $p = .978$), their reported prior level of confidence from 1 to 10 that they would complete the task (successful average 8.8, linear slope $-.002$, $p = .972$), whether they reported Health or Education as one of their primary areas of research (Health difference .059, $p = .621$, Education difference .098, $p = .413$), and to experience/representation in the literature as proxied by number of published peer-reviewed papers as of May 2020 (successful mean 8.75 and median 7.5, linear slope $-.002$, $p = .666$, linear slope after inverse hyperbolic sine transformation .044, $p = .388$, median 7.5 vs. 8). The most noticeable difference in attrition rates among observables is that researchers who reported Labor as one of their primary fields of interest were more likely to finish, but even this is not statistically significant at standard levels (difference in attrition .198, $p = .098$). Among successful replicators, the mean and median year in which they received their PhD was 2011 and 2014, respectively. As of October 2020, the mean and median Google Scholar “cited by” count among successful replicators were 366 and 128, respectively, with a minimum of 38 and a maximum of 1291, omitting one replicator who did not have a Google Scholar profile.

While there may of course be differential attrition by unobservables, and there is no attempt here to correct for selection into volunteering in the first place, the set of replicators who finished the task looks very similar to the set of replicators who did not.

3.4 | Analysis

Replicators return to the organizers their raw data files, code for data processing and analysis, and a primary result of interest. Organizers then perform a descriptive analysis of the results and code.

Analysis proceeds first by taking the produced analyses and comparing them in absolute terms, analyzing the degree of overlap between analyses as well as in terms of features like included controls and sample size.

Organizers were able to successfully replicate the reported results of all replicators using the provided code, with one exception, where due to version control issues a line of code included in analysis was omitted from the submitted code. Code was later updated to the final version, after that replicator viewed the Results section and notified organizers.

Then, organizers analyzed the submitted code of each replicator line by line. This allowed organizers to code the decisions made by each replicator in the process of cleaning the data and generating variables for inclusion in analysis.

Results consist of a description of the differing decisions made by replicators, and the implications of those decisions for sample construction and analysis.

4 | RESULTS

In total, there were fourteen completed replications: seven for the compulsory schooling and teenage pregnancy study based on Black et al. (2008), and seven for the health insurance and self-employment study based on Fairlie et al. (2011). While we include the results from the original studies for comparison, matching the original study is not the goal, and the instructions were not designed to exactly match the original study, especially in the case of Fairlie et al. (2011), so they are not entirely comparable.

Figures 1 and 2 show the confidence intervals estimated in each study for the preferred estimates that replicators selected, based on reported point estimates and standard errors.⁶ Many replicators performed additional analyses or robustness tests, but we will focus on the estimators they reported as preferred, which in the instructions were said to be the result that would be put in the abstract if these were individual studies being written up.

Results based on the reported point estimates are mixed. Estimates in the compulsory schooling study vary widely across replications. Four are statistically significant at the 95% level and negative, one is statistically significant and positive, and two are insignificant, one of which has a point estimate very near zero. In both cases, the range of replication values centers around the estimates from the original studies, with the distribution mean in both cases inside of the original confidence interval. For compulsory education, though, the individual estimates do not match that well; many of the estimates do not have overlapping confidence intervals with the original. Keep in mind that the instructions given to replicators did not match the exact original analysis, especially for the health insurance analysis.

The results imply that different researchers answering the same question using the same data set may arrive at starkly different conclusions. Three (of seven) would likely conclude that compulsory schooling had a negative and statistically significant effect on teen pregnancy, two would find no significant effect, and one would find a positive and significant effect. The confidence interval for the mean of the replication distribution is very wide, and swamps the uncertainty from the original estimates. This variation in results demonstrates the crucial role played by researcher degrees of freedom in applied microeconomics research.

Estimates in the health insurance replication are more consistent. One replication has a much larger estimate than the others, and its confidence interval does not overlap with any others. Among the other six, while no two estimates are the same, point estimates are within a fairly narrow band, all of them have overlapping confidence intervals, and the confidence interval for the mean of the replications is reasonably precise. Statistical significance does vary across replications, though. Even with largely overlapping confidence intervals across replications, five researchers (of seven) would likely conclude a significant effect of employer-based health insurance on entrepreneurship, while the other two would find no evidence of this effect.

Differences between the point estimates are not the only matter of interest. We are also interested in the extent to which choices made by replicators differed in the ways they put together their data and designed their analyses, and how these choices affected the differences in results.

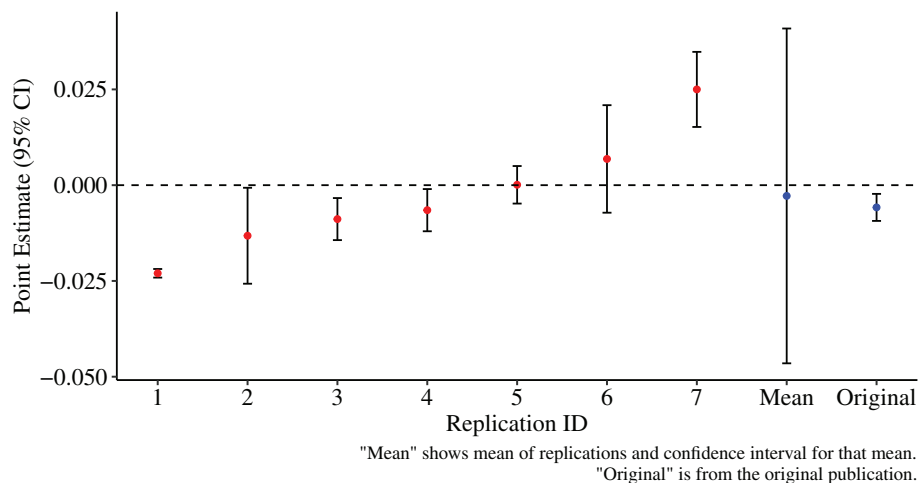


FIGURE 1 Results from compulsory education study

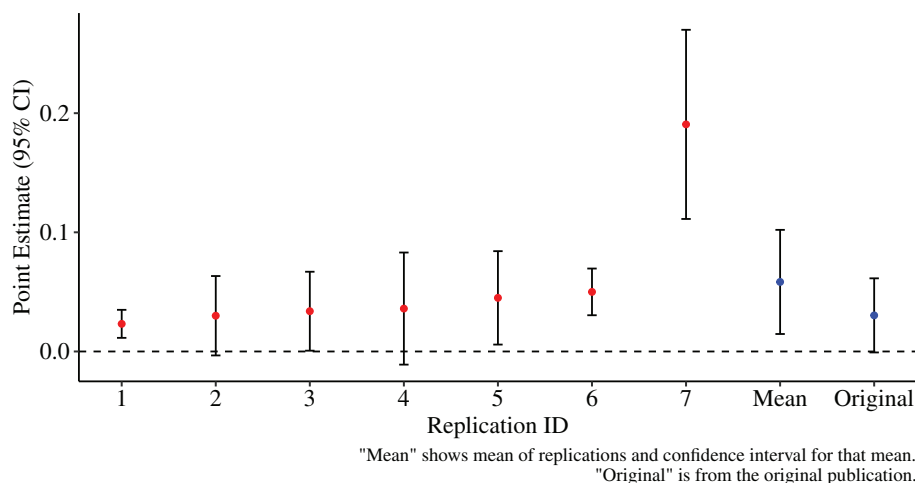


FIGURE 2 Results from health insurance study

4.1 | Compulsory schooling study

In this section, we examine the general study procedures that each replication in the compulsory schooling study took, the ways in which they constructed their data, and the decisions made during analysis.

4.1.1 | Compulsory schooling study procedures

There were seven completed replications of the compulsory schooling study. All seven were completed by the primary replicators, and none reported using assistants. All seven reported that they were familiar with the methods used in the replication, and four of the seven reported that the topic was similar to the work they normally do. One replicator reported after completion that they recognized the original study from the instructions. All seven replications were completed in Stata. The number of replications was not large enough to look for systematic differences between replicators based on their characteristics.

4.1.2 | Compulsory schooling study data construction

Table 1 shows the steps taken in data construction in the compulsory schooling study. In all cases, replicators gathered US Census Data Files from IPUMS and limited the data to adult women subjects born in states with available policy information, as instructed.

Aside from these shared decisions, replicators made different decisions about cleaning the data. The instructions said that women with first births age 14 or below should be dropped. However, one study instead dropped those with first births age 13 or below. Another kept these women but coded them as not being teenage pregnancies. One study did not limit the data to women aged 20–30, as instructed. One study did not match policy dates to individuals in the exact way described in the instructions. Afterwards, one replicator reported having made these decisions because they misread the instructions. Other cases may be due to thinking the differing decisions were more appropriate.

Policy information was provided to replicators. Six replicators left a strange Ohio policy point in the data. In the original study, the table of compulsory schooling laws is 18 for Ohio in every year but 1944, when it is 8. This may be a typo rather than a real policy. Organizers carried this 8 through in the version of the table given to replicators. One replicator (6) changed this to 18 in their main analysis, although several others did point out that it was odd, and said that they might have changed it if they were truly working on their own, but thought that might go against the instructions.

Replicators were instructed to perform the analysis as though they were designing it themselves, and so naturally some data construction decisions not in the instructions are made differently across replicators. One replicator each made the decision to drop subjects in group quarters, to drop the second 1970 census sample, to drop women who never

TABLE 1 Sample creation and shared-variable definition decisions in compulsory schooling study

Study	1	2	3	4	5	6	7
<i>Decisions affecting sample size: (* indicates that this is specified in instructions)</i>							
US census data files from IPUMS*	X	X	X	X	X	X	X
Adult women subjects only*	X	X	X	X	X	X	X
Ages 20–30 only*	X	X	X	X	X		X
Drop women with first-birth age 14 or below*	X	X	X			X	X
Drop women with first-birth age 13 or below				X			
Women with first-birth age 14 or below coded as not being teenage births					X		
Excludes states without policy information (AK, HI)	X	X	X	X	X	X	X
Excludes anyone living in group quarters					X		
Drops all observations from second 1970 Census sample		X					
Drops women without children ever		X					
Keeps only household heads or spouse to household head, with in-house child related to household head						X	
<i>Decisions affecting shared variable definitions</i>							
Changed strange-looking Ohio policy data point						X	
Matches policy years to individuals as in instructions*	X	X		X	X	X	X
Counts age-18 births as “child by age 18”	X	X	X	X		X	X
Sample size used in estimation	1,664,643	831,189	1,696,522	1,701,516	1,669,105	4,271,245	1,640,645
Point estimate	−0.0230	−0.0132	−0.0089	−0.0065	0.0001	0.0068	0.0250
Point estimate under simple shared model (fixed effects for state and birth year, nothing else)	−0.0164	−0.0229	−0.0089	−0.0113	−0.0086	0.0103	0.0177

had a child, and to keep only families where the woman was head or spouse to the head with an in-house child. There was also some minor variation in variable definitions, with one study counting women giving birth at age 18 exactly as not having a child by age 18.

These sample construction decisions led to different sample sizes from every replicator. No two replications had the same sample size, although most are similar. The smallest sample size is 831,139, driven by dropping women without children and one of the census samples. The largest is 4,271,245, driven by including women outside the age range of 20–30. All other samples are fairly similar but not exactly the same, ranging from 1.64 million to 1.70 million.

The small differences in sample size may still be important. Even among the five studies with similar sample sizes, the point estimates vary widely, and even the sign is not consistent. Some of this may be due to differences in analysis rather than differences in sample construction. To account for this, we run the same basic two-way fixed effects model on all seven constructed data sets with no sample weights or other controls.⁷ This reduces differences between estimates, suggesting that some of the differences are due to analysis rather than sample construction. But important differences remain, and the sign is still not consistent. Of particular interest are replications 1 and 7, which do not differ in sample construction in any obvious way, and which likely would have reported identical data construction procedures if these were real studies, but for which sample sizes differ by about 4,000. When using identical models they have similarly-sized estimates of opposite signs.

In the case of this compulsory schooling study, replicators did not perfectly agree on the proper approach to constructing the sample. Some of the differences between approaches, such as dropping women who never had children, would have been reported in a research paper and so could have been evaluated by a reader. Others may not have been.

Studies here that are identical on all data-construction decisions we checked still result in different sample sizes and different point estimates of the treatment effect.

4.1.3 | Compulsory schooling study data analysis

The compulsory schooling study, by design, is based on the concept of a policy that is administered at the state level, such that those in a given birth cohort are exposed, or not exposed, to a certain level of compulsory schooling based on the state they are born in. Except for the fact that treatment does not change monotonically over time within state, this is similar to a difference-in-difference setup. Accordingly, six of the seven replicators used a regression model with two-way fixed effects for state and birth cohort, the standard approach to estimating a difference-in-difference setup with multiple treated groups and variation in treatment timing. The seventh (replication 1) uses two-way fixed effects with state and year of observation.

While no two replicators performed the exact same analysis, all seven replicators made very similar choices in performing the analysis. In addition to all seven using a regression model with state fixed effects and a second set of time fixed effects, all seven clustered standard errors at the state level (replication 5 additionally clustered at the birth year level). Despite a binary dependent variable, all seven used ordinary least squares rather than logit or some other nonlinear model. None of the seven used recent developments in difference-in-difference estimators or standard error adjustments. However, many of these developments (Goodman-Bacon, 2018) were very new at the time replications were performed, and the analysis in question is not exactly the same as difference-in-differences.

The main points of difference between the analyses were whether the second set of fixed effects should be for birth year or year of observation, the Stata commands used to estimate the model, the choice of additional control variables, and the use of sample weights. The two-way fixed effects model was most commonly estimated using the `reghdfe` command (replications 1, 3, 5) and `regress` (2, 4, 6), with 7 using `areg`.

Choice of control variables varied considerably. Table 2 shows the choice of control variables in each replication. As previously mentioned, all studies include state fixed effects and all but one include fixed effects for birth year. Four studies additionally control for the time of observation in some way, either with age or year fixed effects as in 1, 4, 6, and 7. Three studies (2, 4, and 5) include dummies to control for race. One study (5) includes prior time trends by state.

Because all studies used the same design, differences in point estimates can only be driven by differences in data construction or the choice of controls or regression command. To see how much variation is left in point estimates after accounting for data construction differences, we fix the cleaned data set to be that from replication 4, chosen arbitrarily from the seven.

After restricting data to be the same, differences between replications remain (see the final row of Table 2). This indicates that the choice of control variables, even when selecting across different sets that all may seem reasonable, still has a meaningful effect on the published coefficient.

TABLE 2 Control variables included in compulsory schooling study

Study	1	2	3	4	5	6	7
State fixed effects	X	X	X	X	X	X	X
Birth year fixed effects		X	X	X	X		
Race		X		X	X		
Year fixed effects	X			X		X	X
Age fixed effects						X	X
Year-by-age fixed effects							X
State linear time trends					X		
Spouse is household head						X	
Person sample weights used		X				X	X
Point estimate	−0.0230	−0.0132	−0.0089	−0.0065	0.0001	0.0068	0.0250
Point estimate under prepared data from Replication 4	−0.0393	−0.0082	−0.0112	−0.0065	0.0028	−0.0280	−0.0069

Taking both the data construction and data analysis together, both account for a similar share of the initial variation in results, although they overlap in the variation being explained. The sum of squares in the original results is 0.008. Fixing the model used, as in Table 1, reduces this sum of squares by 84%. Instead fixing the sample used, as in Table 2, reduces the sum of squares by 83%. Fixing both would reduce the sum of squares by 100%, by construction.

4.2 | Health insurance study

In this section, we examine the general study procedures that each replication in the health insurance study took, the ways in which they constructed their data, and the decisions made during analysis.

4.2.1 | Health insurance study procedures

There were seven completed replications of the health insurance study. Four were completed by the primary replicators, two (2 and 3) were completed with graduate student assistance, and one (1) had most coding done by a graduate student assistant. Six of the replicators (other than replicator 1) reported that the statistical work was similar to the work they normally do, with two of those reporting that the topic was similar to the work they normally do. All replicators reported not recognizing the original study from the instructions. Six of the replications were completed in Stata, and one (6) was completed in R. The number of replications was not large enough to look for systematic differences between replicators based on their characteristics.

4.2.2 | Health insurance study data construction

Table 3 shows the data construction decisions made by replicators working on the Health Insurance study. In all cases, replicators used monthly NBER CPS files from May 2004 to December 2006, limited to men only. CPS subjects are interviewed for 4 months in a row twice, with a break between the two runs. Since an individual can only turn 65 in the middle of one of those runs, not both, this will produce a small-T rolling panel data set with approximately four observations per individual.

After this point, data construction procedures diverge. The biggest point of divergence is in defining the age range. The instructions specify that subjects should be “observed in the exact month that they turn 65.” However, replications 1–3 and 6 include a wider range of subjects in the data set. The widest range is in Replication 1, which includes subjects aged 54–76. After the fact, replicators reported their reasoning for this decision. Two reported misreading the instructions, and the other two reported that they thought their age range choice was more appropriate. Similarly, the instruction that subjects must be employed was implemented as instead being in the labor force in Replication 4. This replicator reported that this was due to a misreading of the instructions, but that they would have likely made the same choice if writing their own paper.

Replicators were instructed to perform the analysis as though they were designing it themselves, and so naturally some data construction decisions not in the instructions are made differently across replicators. In particular, replicators implemented different kinds of checks on the plausibility of the data. Some dropped individuals with inconsistent demographic data, or who did not appear all four times in the CPS sample, or who were missing income data. Replicators also differed on whether they defined self-employment status using the first worker-class variable in the data, or using both worker-class variables.

The sample sizes differ between the replications, and no two replications have the same sample size. The biggest reason for this is the choice of age ranges, which would have been reported if these replications were written in their own studies. However, even if all age ranges are narrowed to match the instructions, sample sizes are still, in order, 3,543; 5,604; 5,212; 2,493; 1,628; 4,322; and 2,016 (mean 3,545, standard deviation 1,567).

Despite the large differences in sample sizes, the differences in effect sizes are much smaller here than for the compulsory schooling study.⁸ Replication 7 is the only outlier. However, the differences are large enough that some results are statistically significant at the 95% level (1, 3, 5, 6, 7), while others are not (2, 4).

Differences in effects may be due to differences in analysis in addition to differences in sample construction. To account for this, we perform the same basic analysis using the data sets from all seven replications, simply comparing

TABLE 3 Sample creation and shared-variable definition decisions in health insurance study

Study	1	2	3	4	5	6	7
<i>Decisions affecting sample size: (* indicates that this is specified in instructions)</i>							
NBER CPS Monthly Files May 2004–December 2006*	X	X	X	X	X	X	X
Men only*	X	X	X	X	X	X	X
Observed both before and after turning 65*				X	X		
Observed between aged 64 to 65						X	
Observed between ages 63 to 66							X
Observed between ages 60 to 70		X	X				
Observed between ages 54 to 76	X						
Must be employed*	X	X	X ^a		X	X	X
Must be in labor force				X			
Drop observations with panel-inconsistent data					X		X
Drop those not observed four times					X		
Drop those for whom May 2004 is month-in-sample 4 or 8, or for whom December 2006 is MIS 1 or 5	X						
Drop missing income	X					X	
<i>Decisions affecting shared variable definitions</i>							
Treatment compares age 65+ to 64–*		X	X	X	X	X	X
Treatment compares age 65 to other ages	X						
Self-employment given by first worker-class variable	X			X	X	X	X
Self-employment given by both worker-class variables		X	X				
Sample size used in estimation	156,533	90,035	85,400	2,493	1,628	12,288	2,016
Point estimate	0.0232	0.03	0.0338	0.0360	0.0450	0.0501	0.1906
Point estimate under simple shared model	0.1267	0.0751	0.0788	0.0089	0.0522	0.0474	0.0121
Point estimate under simple shared model and shared age range	0.0007	0.0063	0.0074	0.0089	0.0522	0.0091	0.0121

^aThe initial code turned in for this replication had dropped the line in which the sample was limited to the employed, and this X was added after the replicator read this results section. However, because the provided result already used data that included that line, this adjustment did not change anything else.

the proportion of people who are self-employed above 65 vs. below 65 with no controls or sample weights. However, this consistent comparison may actually exaggerate differences, as the studies with large age ranges generally adjust for them with age controls. Accordingly, estimates still vary widely after making the model consistent. So we also perform the shared analysis while narrowing the age ranges to match the instructions. After doing so, while there is still some variation in the effect, results are very similar. This suggests that the differences in results for this replication study are largely due to differences in modeling, and the decision of how wide of an age range is included in the sample.

4.2.3 | Health insurance study analysis

The health insurance study analysis, by design, looks at people just above and just below an age cutoff, which lends itself to a regression discontinuity design (RDD), albeit one with very little variation in the running variable. Some replicators explicitly used regression discontinuity, while others compared the raw average above and below the cutoff, in effect a regression discontinuity with a zero-order polynomial.

Analysis decisions were more heterogeneous for the health insurance study than for compulsory schooling. Table 4 shows the decisions that replicators made.

Four replicators explicitly used regression discontinuity, but each was different, fitting a linear (2 and 4), quadratic (7), or cubic (3) RDD. The other three used a binary treatment indicator (zero-order polynomial RDD). The exception is Replication 1, which as mentioned in the previous section compared 65 to others rather than above/below.

There was also variation in the use of nonlinear models. The dependent variable, self-employment, is binary. Most replicators used linear probability models, but 1 and 3 used probit, while 6 used logit. 2 and 4 used Stata's `rdrobust` function to run regression discontinuity, and the rest used Stata's `regress` function to either perform RDD or use a binary treatment indicator.

There are also many differences between analyses in the controls used, in a way that does not fit well into Table 4, because constructs like race and education are controlled for in different ways in different replications. Controls include:

- Replication 1: Family income (midpoint of bins, treated linearly), education (all included levels), race (white/black/other), marital status, citizenship status, presence of own-children under 18 in household, industry indicators for agriculture, financial services, real estate, health services.
- Replication 2: Education (8th grade / college degree / postgraduate degree / less-than-8th-grade).
- Replication 3: Race (all included levels), education (all included levels), marital status, metropolitan area/non, state.
- Replication 4: None.
- Replication 5: Month of interview, year of interview, race (white/nonwhite), education (below/HS/above), marital status, metropolitan area/non.
- Replication 6: Family income (16 bins), number of people in the household, marital status (all included levels), education (all included levels), citizenship status, census region.
- Replication 7: Month in sample, race (all included levels), education (below/HS/College grad), date.

The list of differences between analyses is long, and are compounded by the fact that some analyses, like cubic RDD, rely on the wide age ranges discussed in the previous section. To evaluate the impact of analysis differences on point estimates, we estimate each model using the data from Replication 5. This process does require changing some of the analyses: specifically, dropping polynomial terms for age in RDD and other contexts. The two replications using `rdrobust` act strangely in this case, and either cannot run (Replication 2), or produce a surprising result that may be due to the command being applied to this particular data set (4). The non-`rdrobust` models, however, produce results with a similar spread to the original estimates, with

TABLE 4 Analysis decisions in health insurance study

Study	1	2	3	4	5	6	7
<i>Decisions affecting analysis</i>							
Regression discontinuity							
Linear		X		X			
Quadratic							X
Cubic			X				
Above/below binary							
Linear probability model		X		X	X		X
Probit/logit	X		X			X	
Heteroskedasticity-robust SEs							
Clustered SEs (individual)					X		
Clustered SEs (state)		X	X				
Stata regress					X		X
Stata probit	X		X				
Stata rdrobust		X		X			
R glm(link = "logit")						X	
Person sample weights			X				X
Point estimate	0.0232	0.0300	0.0338	0.0360	0.0450	0.0501	0.1906
Point estimate under data from Replication 5 (and RDD terms restricted to linear)	0.0488	NA	0.0353	−0.0203	0.0210	0.0307	0.0121

the exception of Replication 7, which drops much closer to 0 (see the final row of Table 4). This is consistent with the previous section in suggesting that much of the variation is due to differences in age ranges combined with analytical choices.

Taking both the data construction and data analysis together, differences in the model and age range account for more variation in results than does data construction, although they overlap in the variation being explained. The sum of squares in the original results is 0.021. Fixing the model and age range used, as in Table 3, reduces this sum of squares by 91%. Instead fixing the sample used, as in Table 4, reduces the sum of squares by 81% (after scaling the sum of squares by 6/5 to account for the missing adjusted value in Table 4).

5 | CONCLUSION

Given the same data and research question of interest, we find considerable variation across researchers in the way that they clean and prepare data and design their analysis. In one of the two studies we examine, this led to considerable variation in results.

In both studies, the estimated sampling variation within studies was small relative to variation between studies. In the compulsory education study, the average reported standard error was 25.1% as large as the standard deviation of reported effects across studies. This figure was 32.5% in the health insurance study. In both cases, standard errors omit a major source of variation in estimates.

It is not surprising that different researchers would carry out an analysis in different ways. Replicators were asked after completing their replication about their reasoning for the analytic and data cleaning choices that were not covered by the instructions and differed among replicators. The most common reasons included familiarity with a given model, differing intuitive or technical ideas about which control variables are appropriate or whether linear probability models are appropriate, and differing preferences for parsimony.

There is nothing inherently wrong about these choices or reasons, although the fact that researchers do not seem to agree on these issues implies additional sources of uncertainty in estimates. These differences only rise to a real cause for concern when they are about things that either would be unlikely to be reported in the resulting study, or would be reported but paid little attention by reviewers and readers. If invisible researcher choices are different and consequential, that means that empirical results in applied microeconomics reflect variation in sample and methods, as expected, but also reflect variation in researcher choice. And while this variation is not the same thing as publication bias or “p-hacking,” as these choices are not necessarily related to an attempt to report a particular result, this does make it easier for an unscrupulous researcher to attempt many different analyses and so get a desired result without detection.

The biggest issue highlighted by these results is the considerable differences between researchers in the way the data was cleaned and prepared. No two researchers had the same sample size in their analysis. Nearly all of the decisions driving data construction would be likely to be omitted from a paper, or skimmed over by a reader. The differences between researchers in data cleaning is the finding here that seems most likely to generalize.⁹

Consideration of what to do about these researcher degrees of freedom must recognize that this is different from several other well-known issues with the reliability of published studies, and so will require different solutions. The results in this paper arose without publication contingent on the result, so the advent of preregistration, while a promising solution to other problems, is not well-suited for handling researcher degrees of freedom. It is also not clear that these differences are because replicators were making *wrong* decisions. So, even in the case where a reader or reviewer goes through the preregistration or data preparation code for the paper, they might not see any problem. For this reason, full availability of code or thorough data appendices would also not address the problem directly, though making a “Data Appendix” a more standard feature of economics papers, even for studies using standard data sources, would likely be beneficial. American Economic Association (2020) has already taken steps in this direction, although they accept code that starts from data that has already been partially processed, allowing some decisions to be hidden.

Other solutions to the problem of researcher degrees of freedom in data construction could support the standardization of data cleaning. IPUMS (Ruggles et al., 2020) and the Current Population Survey Merged Outgoing Rotation Group Files (NBER, 2020) are popular, and remove many researcher degrees of freedom in the processing of raw data, although not from any point later in the data preparation process. Standard preprepared data for other common data sources may help make results more consistent. In application to less-common data sources, though, there is not a well-known set of “best practices” for data cleaning and preparation in economics, at least not to the extent that there is with analysis. An attempt to develop a set of best practices would also help standardize results. This process may also aid in the process of making data more widely available for a broader range of researchers, given the contrasting incentives that individual

researchers have in making their cleaned data available, as opposed to a centralized data-cleaning organization, and given the scale of tasks such as procuring harmonizing records across multiple sources (Hill et al. 2020).

Addressing researcher degrees of freedom that come from analysis seems more achievable, since reviewers and readers are already in the habit of considering the implications of different analytic choices. However, the results in this study imply that more detail in describing the analysis, if only in an appendix, is advisable. More space can be given in economics papers to justifications of decisions like the use of control variables and variable definitions. Another approach is the use of model averaging, in which multiple possible ways of designing the model are averaged together to produce a result (Moral-Benito, 2015).

Finally, before any such changes can be made in response to researcher degrees of freedom in data construction or analysis, readers should be aware of the hidden variation in results that stems from these choices. The effective distribution of an effect is likely wider than the published standard error would indicate, as the evidence presented here suggests.

ACKNOWLEDGMENT

Additional thanks to Andrew Gill, Seth Gitter, Shikhar Mehra, and Virginia Wilcox.

ORCID

Nick Huntington-Klein  <https://orcid.org/0000-0002-7352-3991>

ENDNOTES

- ¹ And in the case of McCullough and Vinod (2003), economics papers rarely report the software or estimation command they used in the text of the paper, despite evidence here that it can be consequential.
- ² The IDEAS/RePEc Aggregate Rankings for Journals (<https://ideas.repec.org/top/top.journals.all.html>) list was used as a guide.
- ³ Indicators that a study may be too well-known include a very large number of citations or extensive media coverage of the study.
- ⁴ While the goal of this paper is not to judge the quality of the original findings being replicated, in preparing instructions for this study we happened to find that the Fairlie et al. (2011) analysis is sensitive to the decision of which years of data to include in the analysis. The effect using their original analysis can even reverse sign depending on the time period used. We did not check if this issue extends to the other analyses in that paper. However, this is an excellent example of results being affected by seemingly innocuous researcher choices.
- ⁵ <https://sites.google.com/view/replication2018>.
- ⁶ In some cases, coefficients from nonlinear models were reported in the original replications; in these cases we calculate marginal effects after the fact for comparability.
- ⁷ We use the Stata command `reghdfe` with “child by age 18” as the dependent variable, allowing the definition to be different for Replication 5, as in Table 1. Being treated is the only independent variable. Birth year and state are absorbed fixed effects.
- ⁸ We calculated marginal effects ourselves in cases where authors reported logit or probit coefficients.
- ⁹ A natural question for this study on the sensitivity of results to hidden decisions is how sensitive these results themselves are to seemingly innocuous decisions, the most glaring of which would be the selection of the two studies. Would these results replicate given a different two basis studies, or different replication teams?

REFERENCES

- American Economic Association. (2020) *Data and code availability policy*. Nashville, TN: American Economic Association.
- Berry, J., Coffman, L.C., Hanley, D., Gihleb, R. & Wilson, A.J. (2017) Assessing the rate of replication in economics. *American Economic Review*, 107(5), 27–31.
- Black, S.E., Devereux, P.J. & Salvanes, K.G. (2008) Staying in the classroom and out of the maternity ward? The effect of compulsory schooling laws on teenage births. *The Economic Journal*, 118(530), 1025–1054.
- Botvinik-Nezer, R., Holzmeister, F., Camerer, C.F., Dreber, A., Huber, J., Johannesson, M. et al. (2020) Variability in the analysis of a single neuroimaging dataset by many teams. *Nature*, 582, 84–88.
- Camerer, C.F., Dreber, A., Forsell, E., Ho, T.-H., Huber, J. et al. (2016) Evaluating replicability of laboratory experiments in economics. *Science*, 351(6280), 1433–1436.
- Chang, A.C. & Li, P. (2017) A preanalysis plan to replicate sixty economics research papers that worked half of the time. *American Economic Review*, 107(5), 60–64.
- Chetty, R., Hendren, N. & Katz, L.F. (2016) The effects of exposure to better neighborhoods on children: new evidence from the moving to opportunity experiment. *American Economic Review*, 106(4), 855–902.
- Christensen, G. & Miguel, E. (2018) Transparency, reproducibility, and the credibility of economics research. *Journal of Economic Literature*, 56(3), 920–980.
- Clemens, M.A. (2017) The meaning of failed replications: a review and proposal. *Journal of Economic Surveys*, 31(1), 326–342.

- Clemens, M.A. & Hunt, J. (2019) The labor market effects of refugee waves: reconciling conflicting results. *ILR Review*, 72(4), 818–857.
- Cohn, N. (2016) *We gave four good pollsters the same raw data. They had four different results*. The New York Times.
- De Long, J.B. & Lang, K. (1992) Are all economic hypotheses false? *Journal of Political Economy*, 100(6), 1257–1272.
- Dewald, W.G., Thursby, J.G. & Anderson, R.G. (1986) Replication in empirical economics: the journal of money, credit and banking project. *The American Economic Review*, 76(4), 587–603.
- Duvendack, M., Palmer-Jones, R. & Reed, W.R. (2017) What is meant by “replication” and why does it encounter resistance in economics? *American Economic Review*, 107(5), 46–51.
- Fairlie, R.W., Kapur, K. & Gates, S. (2011) Is employer-based health insurance a barrier to entrepreneurship? *Journal of Health Economics*, 30(1), 146–162.
- Gelman, A. & Loken, E. (2014) The statistical crisis in science: data-dependent analysis—“a garden of forking paths”—explains why many statistically significant comparisons don’t hold up. *American Scientist*, 102(6), 460–466.
- Gertler, P., Galiani, S. & Romero, M. (2018) How to make replication the norm. *Nature*, 554(7693), 417–419.
- Goodman-Bacon, A. (2018) *Difference-in-differences with variation in treatment timing*. Working Paper 25018, National Bureau of Economic Research.
- Hamermesh, D.S. (2007) Viewpoint: replication in economics. *Canadian Journal of Economics/Revue Canadienne d'Economie*, 40(3), 715–733.
- Hamermesh, D.S. (2017) Replication in labor economics: evidence from data, and what it suggests. *American Economic Review*, 107(5), 37–40.
- Herndon, T., Ash, M. & Pollin, R. (2014) Does high public debt consistently stifle economic growth? A critique of Reinhart and Rogoff. *Cambridge Journal of Economics*, 38(2), 257–279.
- Hill, R., Stein, C. & Williams, H. (2020) Internalizing externalities: designing effective data policies. In: *AEA Papers and Proceedings*, Vol. 110, pp. Nashville, TN: American Economic Association, 49–54.
- Ioannidis, J.P.A., Stanley, T.D. & Doucouliagos, H. (2017) The power of bias in economics research. *The Economic Journal*, 127(605), F236–F265.
- Landy, J.F., Jia, M.L., Ding, I.L., Viganola, D., Tierney, W., Dreber, A. et al. (2020) Crowdsourcing hypothesis tests: making transparent how design choices shape research results. *Psychological Bulletin*, 146(5), 451–479.
- Lenz, G.S. & Sahn, A. (2017) *Achieving statistical significance with covariates and without transparency*. Technical Report, University of California Berkeley.
- McCullough, B.D. & Vinod, H.D. (2003) Verifying the solution from a nonlinear solver: a case study. *American Economic Review*, 93(3), 873–892.
- Moral-Benito, E. (2015) Model averaging in economics: an overview. *Journal of Economic Surveys*, 29(1), 46–75.
- National Bureau of Economic Research. (2020) *CPS merged outgoing rotation groups*. Cambridge, MA: National Bureau of Economic Research.
- Reinhart, C.M. & Rogoff, K.S. (2010) Growth in a time of debt. *American Economic Review*, 100(2), 573–578.
- Ruggles, S., Flood, S., Goeken, R., Grover, J., Meyer, E., Pacas, J. & Sobek, M. (2020) *IPUMS USA: Version 10.0* [dataset]. Technical Report, IPUMS. Minneapolis, MN: IPUMS.
- Silberzahn, R., Uhlmann, E.L., Martin, D.P., Anselmi, P., Aust, F. & 61, m. (2018) Many analysts, one data set: making transparent how variations in analytic choices affect results. *Advances in Methods and Practices in Psychological Science*, 1(3), 337–356.
- Simmons, J.P., Nelson, L.D. & Simonsohn, U. (2011) False-positive psychology: undisclosed flexibility in data collection and analysis allows presenting anything as significant. *Psychological Science*, 22(11), 1359–1366.
- White, T.K., Reiter, J.P. & Petrin, A. (2018) Imputation in US manufacturing data and its implications for productivity dispersion. *Review of Economics and Statistics*, 100(3), 502–509.
- Young, A. (2018) *Consistency without inference: instrumental variables in practical application*. Unpublished.

SUPPORTING INFORMATION

Additional supporting information may be found online in the Supporting Information section at the end of this article.

How to cite this article: Huntington-Klein, N., Arenas, A., Beam, E., Bertoni, M., Bloem, J.R., Burli, P. et al. (2021) The influence of hidden researcher decisions in applied microeconomics. *Economic Inquiry*, 59:944–960. <https://doi.org/10.1111/ecin.12992>