# Transparency in an Era of Data-Driven Policy: The Importance of Reproducible Research

## EXECUTIVE SUMMARY

The advent of cheap digital storage and powerful computer processing has driven an explosion of data collection and empirical research on topics ranging from public health initiatives to stock returns to climate policy. Think tanks such as the Brookings Institution and government entities such as the Environmental Protection Agency encourage and employ data-driven "evidence-based policymaking" to make decisions and recommendations.[1] Empirical research advances our understanding of how the economy works and can improve individual, business, and policy decisions, but only insofar as the research conclusions are accurate and credible. In this Policy Spotlight, we review the evidence of the reproducibility of published social science research and discuss how to improve the reproducibility of future studies. We also use one of our own published studies to demonstrate how researchers can facilitate reproducibility.

## BACKGROUND

As computer storage and processing capacity have grown, so too have the size and sophistication of research projects that harness these tools. Teams of researchers often spend years processing and analyzing datasets that can include upwards of billions of observations. For academic research projects, the results are then described in a paper and submitted to a peer-reviewed journal for publication. During peer review, several experts read the paper and assess the innovation and credibility of the analysis. If the study passes these checks, then the journal decides to publish it.

UNIVERSITY OF ILLINOIS SYSTEM

INSTITUTE OF GOVERNMENT AND PUBLIC AFFAIRS

*Authors*

**David Molitor** is an Associate Professor, Department of Finance, Gies College of Business, University of Illinois Urbana-Champaign.

**Julian Reif** is a Senior Scholar at the Institute of Government and Public Affairs, University of Illinois System; Associate Professor, Department of Finance, Gies College of Business, University of Illinois Urbana-Champaign.

While reviewers are encouraged to scrutinize written content describing a study's results, they seldom examine the computer code and datasets used to produce those results. Even after publication, many journals do not require a study's supporting data and code to be made publicly available. In cases where journals do require supporting data and code, it is often poorly documented and incomplete, and frequently fails to reproduce the published results.[2] As we discuss below, some journals have made significant progress in addressing these issues, but documented reproducibility of study results prior to acceptance remains far from the norm.

The issue of reproducibility is not just academic; it also has ramifications for real-world public policy. For example, following the introduction of no-fault divorce laws in the US in the 1970s, an influential study found that women's standard of living decreased by over 70% after divorce.[3] This study served as a driving force behind several subsequent divorce law reforms.[4] However, a 1996 study re-evaluated the data and concluded that women's standard of living decreased by less than 30%.[5] The large discrepancy between these studies was believed to have been caused by a statistical weighting error in the original study, but the code for the original analysis is no longer available to confirm that suspicion. Had the original code been made publicly available, the mistake could have been caught sooner. Recognizing the importance of reproducing research results for scientific rigor and public trust in science,

> **Transparency plays a key role in increasing the accuracy and credibility of research findings. Requiring researchers to publish their data and code encourages investments in better data management practices at the outset of a project, reducing errors.**

Congress mandated in 2017 that the National Science Foundation undertake an examination of reproducibility and replicability in science.[6] This initiative identified factors that undermine reproducibility and replicability in research and recommended steps for researchers, academic institutions, and journals to improve transparency in scientific research.[7]

Transparency plays a key role in increasing the accuracy and credibility of research findings. Requiring researchers to publish their data and code encourages investments in better data management practices at the outset of a project, reducing errors. Publishing the underlying analysis also allows outside researchers to spot mistakes and makes it easier to build on the original research by democratizing access to data and code. Just as important, making research analyses publicly available for anybody to peruse increases the public's faith in the scientific process.

## WHAT IS SCIENTIFIC REPLICABILITY AND REPRODUCIBILITY?

A replicable study yields the same results when it is repeated by other researchers. There are many reasons why a study may fail to replicate. For example, the results may have been a fluke or may not generalize to other populations or settings. In other cases, the original researchers may have made a mistake in data collection or analysis, invalidating the findings. Understanding the conditions under which a study does or does not replicate sharpens the scientific contribution of the original study.

To illustrate, suppose we are interested in evaluating a new surgical technique. Researchers conduct an initial study in a single hospital with 100 patients and estimate that this technique increases post-surgical life expectancy by one year. If subsequent studies conducted in other hospitals find similar results, widespread adoption of the new technique is likely. By contrast, if other studies find no or negative effects of the new technique, adoption is unlikely. Note that failure to replicate does not necessarily imply the original study was mistaken or poorly run. Perhaps the new technique requires specialized skills only present in the original hospital, or perhaps the procedure is most effective only when performed

on the particular mix of patients commonly found at that original hospital. Understanding why a study fails to replicate helps inform follow-up research: a study that fails to replicate because it was riddled with errors can be safely ignored, but a different study that fails to replicate only in younger populations informs researchers that the results could perhaps still be applied to older populations.

A broad notion of scientific replication involves repeating a study multiple times using data collected from a new population, often in different settings or analyzed using different methods from the original study. As noted in the illustration above, the results may not replicate if they apply only to individuals with particular characteristics. The results may also not replicate if the original findings were a statistical fluke. When sample sizes are small, estimates are noisy and thus an erroneously large estimate is more likely to arise by chance.[8] Researchers who run several experiments may choose to publish only studies that find large effects and to ignore those with null findings, which generates statistical bias often referred to as "publication bias" or the "file drawer problem." In these cases, rerunning the same experiment with a larger sample size will often result in a smaller or even non-existent effect.

A narrower notion of scientific replication is that of computational reproducibility. A study is computationally reproducible if a researcher can use its data and methods to exactly reproduce the original results.

While computational reproducibility ought to be satisfied by most analyses, there are various reasons it may not. For example, researchers may fail to provide code that specifies and automates all steps taken to clean and analyze the study data. Studies that combine multiple datasets with different formats frequently must make a number of important decisions regarding which observations to include and how to define variables. A study of the effect of education on earnings might exclude adults not in the labor force, and a study measuring wealth disparities might choose to ignore future retirement benefits. If these decisions are not well documented, then it is difficult for other researchers to follow the original methodology, which increases the chance that the study results cannot be reproduced.

> **A narrower notion of scientific replication is that of computational reproducibility. A study is computationally reproducible if a researcher can use its data and methods to exactly reproduce the original results.**

Computational reproducibility may also fail if the original analysis itself contains errors. For instance, researchers may load the wrong data, define variables incorrectly, improperly handle missing data values, estimate the wrong model, or copy and paste the wrong results into a table.

### How well do studies replicate and reproduce results?

Unfortunately, recent work suggests many studies—particularly those published in the social sciences—are less replicable than previously recognized. In 2015, Open Science Collaboration conducted replications of 100 empirical studies published in three well-known

psychology journals.[9] Only 36% of replications had significant results, compared to 97% of the original studies. A 2018 study by Klein and colleagues attempted to replicate 28 class psychology studies, but succeeded only half of the time.[10] Similarly, a 2018 study by Camerer and others could only replicate 62% of selected experimental studies published in Nature and Science.[11]

Because there are many reasons—not all of them necessarily problematic—that a study may fail to replicate according to the broad notion of replicability, we focus the remainder of this Policy Spotlight on the narrower notion of computational reproducibility. If a study cannot reproduce its results using its own original data and code, then it becomes difficult to have faith in its results. Unfortunately, while most studies should be computationally reproducible, many are not.

For example, a 2019 award-winning article by Rampini and coauthors in the *Journal of Finance*, a premier finance journal, was recently retracted after an outside researcher discovered that the code provided by the original authors of the study does not reproduce its main findings.[12] The journal released a statement in response to the retraction:[13]

> "The reason why retractions have almost never occurred in the past is unlikely due to the absence of errors in past published research. More likely, errors have not been exposed in the past because replication was not attempted or failures to reproduce or replicate results have not reached the public."
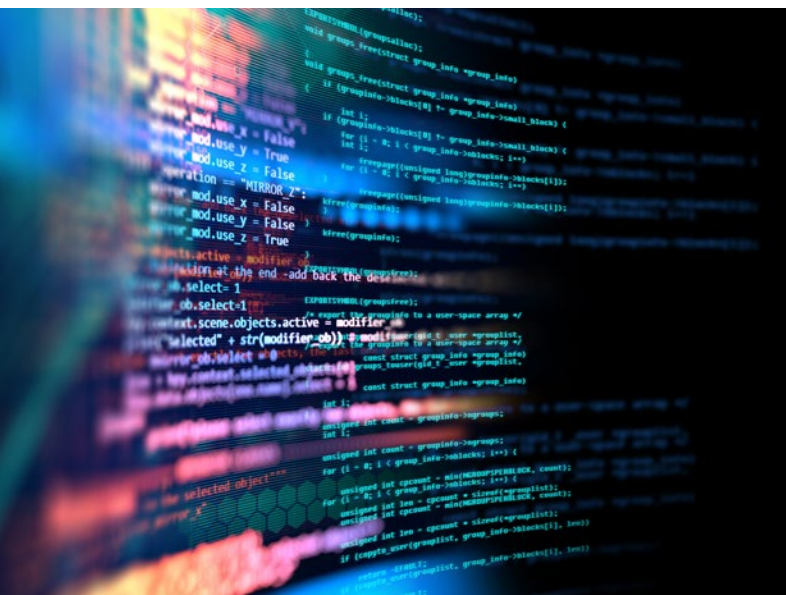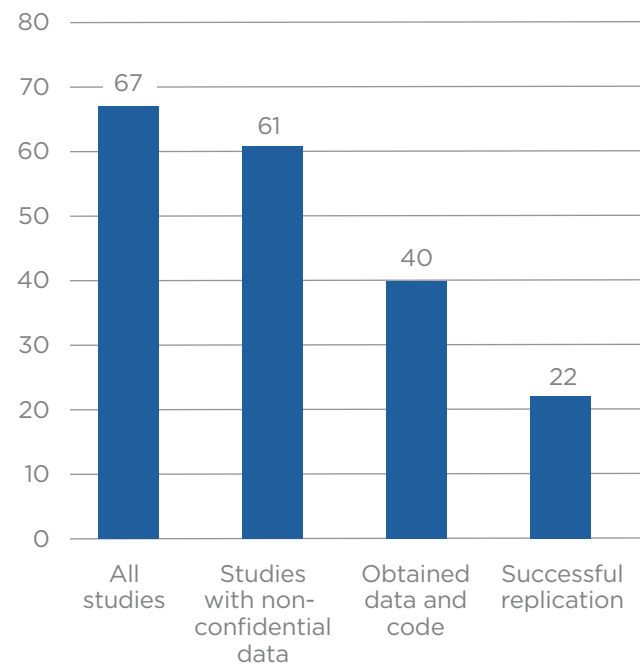


*Figure 1. Number of macroeconomic studies with reproducible results in 13 economics journals*



*Notes: This figure reports results from a study by Chang and Li, who sought to reproduce results from 67 macroeconomic papers published in 13 economics journals.*

The errors in this study were discovered quickly in part because the original authors were required to provide the data and code underlying their results. Unfortunately, in many cases the original data and code are unavailable. For example, researchers Andrew Chang and Philip Li conducted reproducibility checks of 67 macroeconomic studies published in 13 economics journals (see Figure 1).[14] Even after receiving assistance from the original authors, they were unable to reproduce the results for over half of the papers in their sample. The researchers began their reproducibility checks by searching for study data on journal and author websites. If that search was unsuccessful, they then emailed individual study authors. After this process, the researchers were still unable to obtain the data and code for one-third of the studies, even though most were published in journals that required authors to provide these files as a condition of publication.

Even when data and code are unavailable to check a study's analysis, some errors can still be detected. Tim van der Zee and colleagues reanalyzed four published studies from the Cornell

Food and Brand Lab, led by principal investigator Brian Wansink who, prior to Cornell, was on the University of Illinois faculty.[15] Although the researchers were unable to obtain the data from the authors of the published studies, they nonetheless deduced about 150 inconsistencies across the four studies using the statistics reported in those studies. While these studies may have been particularly flawed—Cornell determined that Wansink had committed scientific misconduct and he later resigned—analysis errors such as the ones documented by van der Zee and colleagues are probably common.[16]

### Why do many studies fail computational reproducibility?

Many people would agree that data and code from academic studies should be readily shared and well documented. However, there are many obstacles, and few incentives to encourage these behaviors. Academics are rewarded primarily for publishing papers, and those publications generally do not depend on the transparency of the raw analysis itself. In its statement about the Rampini and coauthors retraction discussed above, the *Journal of Finance* noted that it was not possible for its reviewers to have caught the coding error because, consistent with journal policy, the code was not available during the review process. Indeed, we are unaware of any social science journal that requires a study's data and code to be made available during the peer review process. However, a small number of journals, such as those run by the American Economic Association, have recently begun requiring computational reproducibility checks for papers that have been peer reviewed and "conditionally accepted" for publication.

Several barriers discourage the documentation and sharing of data and code. In many disciplines, researchers are not taught how to write professional code or create reproducible analyses. Documenting and writing a computerized analysis that runs successfully from beginning to end and reproduces all of a study's results can require anywhere from a few hours to several weeks of work, depending on the project and the knowledge base of the researchers. Researchers working with datasets that include personally identifiable information may need to scramble or omit certain variables in their shared dataset. If a mistake in

the analysis is found after the paper's acceptance, then the researcher must inform the journal, which may result in delays or even in the rescinding of the paper's acceptance. Overcoming these barriers may be particularly challenging for early-stage researchers and those who lack research staff, which raises the concern that reproducibility requirements could undermine important research conducted by those with fewer resources.

Sharing data and code can also create future costs for researchers. Once a study's data and code have been released, it becomes easier for other researchers to review and criticize the study. These critiques may, in turn, require the original authors to further clarify or defend their analysis. If a significant mistake is discovered, the authors may need to issue a correction or retract their paper. While issuing corrections and retractions is part of a healthy science—scientists do not want to promulgate incorrect results—doing so often comes at a reputational cost to the original authors, especially since corrections and retractions are currently rare.

These challenges do not mean that researchers receive no benefit from sharing their data and code. Making an analysis publicly available generates goodwill among other researchers, increases faith in the authors' analysis, and potentially increases future citation counts. However, these positive returns have not yet led to a norm of publicly available, reproducible research among social scientists.

> **A small number of journals, such as those run by the American Economic Association, have recently begun requiring computational reproducibility checks for papers that have been peer reviewed and "conditionally accepted" for publication.**

*Potential solutions for improving reproducibility*

The costs of publishing and maintaining a transparent analysis are borne by researchers while the benefits flow primarily to the rest of society. This mismatch encourages researchers to put suboptimal effort into research transparency. Indeed, the numerous reproducibility problems cited earlier in this Policy Spotlight suggest that the transparency of many research studies is quite poor.

One remedy is to update standards in the social sciences profession to better reward transparent research. When reviewing researchers for hire or promotion, the reproducibility of the candidate's research could be included as one component of the evaluation—for example, does the researcher make the data and code underlying her research studies available on her website? Performing reproducibility checks prior to publication, as some journals have begun to do, is another way to encourage transparency. In addition, journals could request that data and code be made available to reviewers or to a data editor during the initial or intermediate phases of submission, rather than only after a decision has already been made to accept the paper. This approach would generate incentives to carefully review data and code earlier in the research process.

Journals could also commit more strongly to publishing corrections and retractions when mistakes are discovered—whether by outside researchers or the original researchers. While corrections and retractions are occasionally published, the process is onerous and seldom successful.[17] A commitment to publishing corrections would encourage extra care during the initial analyses.

Although post-publication reviews of studies can—and should—be carried out on any platform, journals can play an important arbiter role in this process. For example, if transparent, shared analyses become the norm, the possibility arises that some people may critique published papers in bad faith. Researchers are, of course, always free to ignore baseless criticism. But journals can encourage researchers to respond to thoughtful criticism by publishing informed critiques along with responses from the original researchers. A healthy post-publication dialogue can inform future research as well as increase the visibility of the original study.

Because corrections and retractions are currently rare and often reserved for egregious errors or academic misconduct, a commitment to publicizing smaller or inadvertent errors may require a significant change to academic norms. Rather than being viewed as a penalty or humiliating mistake, corrections would need to become viewed as a standard part of the research process. It would also be helpful to introduce a mechanism that alerts other researchers who may have cited the original analysis. Many retracted papers continue to be cited following the retraction, presumably because the citing authors are unaware of the retraction.[18]

A second remedy is to reduce the costs of publishing a transparent analysis. Institutions such as the Inter-university Consortium for Political and Social Research (ICPSR) and the Abdul Latif Jameel Poverty Action Lab (JPAL) provide resources such as data repositories and statistical training to help researchers develop and maintain their analyses. Institutional support for reproducibility efforts can also be designed to prioritize researchers who may face the largest burden from reproducibility requirements, such as early-stage researchers or those who do not have their own research staff.

Another way to reduce the costs of publishing a transparent analysis is to provide technical training. Graduate programs could incorporate materials on best programming and reproducibility practices into their classes. Similarly, providing a template to researchers will reduce the transaction costs for researchers. To that end, we have assembled a code and data repository for one of our recently published studies. The following section explains how this repository was constructed and how it can serve researchers, policymakers, and interested public citizens.

### Reproducible research example: workplace wellness

One of the challenges to sharing reproducible data and code—especially for those without knowledge or experience—is a lack of examples.



To help researchers overcome this obstacle, we created a polished repository that includes public use data and code that reproduces published results from one of our studies. We hope that this repository will be useful to researchers interested in extending our results and to policymakers interested in using our study to inform their decisions. Below, we first describe our study and explain its main findings. We then describe the public data repository we created and how researchers and policymakers can use it.

Some of the most significant drivers of healthcare spending are related to chronic diseases such as obesity and smoking-related health issues.[19] Because people spend a lot of time in the workplace, many businesses employ workplace wellness programs to reduce medical spending, improve employee health, and increase worker productivity.[20] Today, these programs cover more than 50 million U.S. workers, and the workplace wellness industry's annual revenue has grown to more than $8 billion.

Workplace wellness programs have been the subject of numerous prior studies, but most of these studies did not employ a randomized design. To better understand the effects of these programs, we evaluated a randomized controlled trial of a comprehensive workplace wellness program at the University of Illinois. We randomly assigned one set of employees to a group that was eligible to participate in a comprehensive workplace wellness program. These employees received monetary rewards in return for participating in an on-site health screening and completing wellness classes. The second set of employees was assigned to a control group that was not eligible to participate. We collected survey and administrative data on both sets of employees in order to measure the program's effects. These data included health insurance claims, university employment and sick leave data, gym attendance records, and more. Because of the sensitive nature of the information we collected, study data were stored on an offline computer in a locked office.

Our analyses produced three main results:

1. The people most eager to participate in the workplace wellness program already had low healthcare costs.
2. After 30 months, the program had no detectable effects on medical spending; health behaviors; biometric health measurements like weight, blood pressure, and cholesterol; or productivity.
3. Randomization was key to uncovering accurate

estimates: if we had employed an observational design, as most prior studies did, we likely would have drawn incorrect conclusions.

After publishing these initial results, our research team created a public use version of the data we collected. To ensure our study subjects could not be reidentified, we grouped and stored our outcome variables into separate datasets (e.g., claims variables, online survey variables, etc.) and omitted some variables entirely, such as salary.
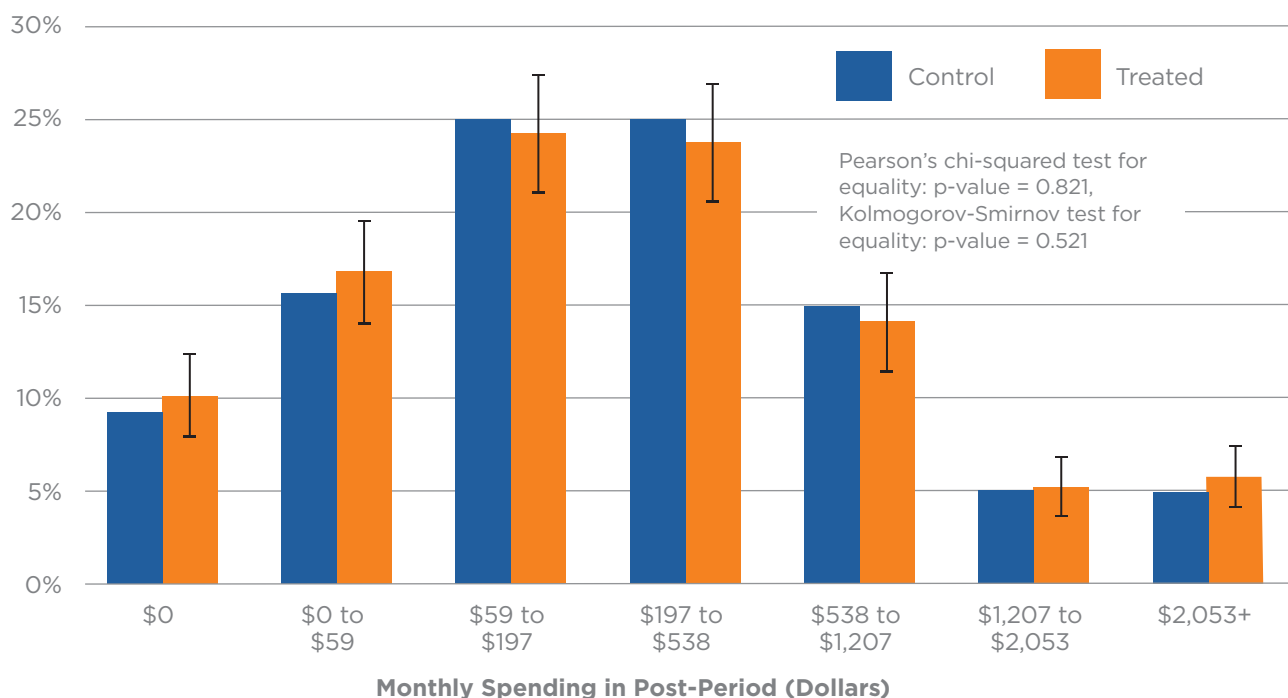
These public use data are now freely and readily available online.[21] Our website also includes accompanying computer code that reproduces the results from our first publication. For example, Figure 2 reproduces Figure 5c from the Jones, Molitor, and Reif study using code provided on our website.[22] This figure illustrates one of the central results of the study: medical spending among people who were eligible to participate in a workplace wellness program (the "treated" group) was not statistically distinguishable from medical spending among people who were ineligible (the "control" group).

## CONCLUSION

The amount of empirical, data-driven research in the social sciences is large and continues to grow. Unlike a theoretical model or logical argument, whose validity depends primarily on written text subjected to peer review, the validity of many empirical papers relies crucially on data and code that in many cases is not reviewed, published, shared, or reproducible. A lack of transparency makes it difficult for outsiders to correct mistakes or to build on prior work in their subsequent research.

As a condition of publication, some peer-reviewed journals have begun requiring accepted studies to make their data and code publicly available when possible and to pass a reproducibility check. Imposing these requirements earlier in the review process and providing better technical training to researchers are two additional ways to encourage better transparency. Finally, this Policy Spotlight includes an example of a study that serves as a pedagogical resource for writing reproducible analyses and provides publicly available data for those interested in building on its research.

*Figure 2. Reproduction of figure from a 2019 study by Jones, Molitor, and Reif using publicly available data and code*



*Notes: This figure reproduces Figure 5c from a 2019 study by Jones, Molitor, and Reif. The figure presents a histogram of average monthly medical spending in the twelve months following a workplace wellness intervention. The figure was created using publicly available data and code from the authors' study website, available at: https://perma.cc/PZK2-2374*

## ENDNOTES

1 "About Us," *The Brookings Institute*, https://perma.cc/2BYF-FBR5; Timothy Puko, "EPA to give preference to scientific studies that disclose data," *Wall Street Journal*, January 5, 2021.

2 Paul Gertler, Sebastian Galiani, and Mauricio Romero, "How to make replication the norm," *Nature* 554 (February 21, 2018): 417-19, https://www.nature.com/articles/d41586-018-02108-9; Andrew C. Chang and Phillip Li, "Is Economics Research Replicable? Sixty Published Papers From Thirteen Journals Say 'Often Not'," *Critical Finance Review* (forthcoming).

3 Lenore J. Weitzman, *The Divorce Revolution* (Collier Macmillan, 1985).

4 Felicia Lee, "Influential Study on Divorce's Impact Is Said to Be Flawed," *New York Times*, May 9, 1996.

5 Richard R. Peterson, "A Re-Evaluation of the Economic Consequences of Divorce," *American Sociological Review*, 61, no. 3 (June, 1996):528-36, https://www.jstor.org/stable/2096363.

6 *American Innovation and Competitiveness Act*, Pub. L. 114-329, Sec. 116, https://www.govinfo.gov/content/pkg/PLAW-114publ329/pdf/PLAW-114publ329.pdf.

7 National Academies of Sciences, Engineering, and Medicine, *Reproducibility and Replicability in Science* (National Academies Press, 2019), https://pubmed.ncbi.nlm.nih.gov/31596559/.

8 John P. A. Ioannidis, "Why most published research findings are false," *PLoS Medicine*, 2, no. 8 (August, 2005): e124, http://dspace.mit.edu/bitstream/handle/1721.1/58674/9-63-fall-2005/contents/readings/ioannidis.pdf.

9 Open Science Collaboration, "Estimating the reproducibility of psychological science," *Science*, 349, no. 6251 (August, 2015): 943, https://www.researchgate.net/publication/281286234_Estimating_the_reproducibility_of_psychological_science

10 Richard A. Klein *et al*., "Many Labs 2: Investigating variation in replicability across samples and settings," *Advances in Methods and Practices in Psychological Science* 1, no. 4 (December 24, 2018): 443-90, https://doi.org/10.1177%2F2515245918810225.

11 Colin F. Camerer *et al*., "Evaluating the replicability of social science experiments in Nature and Science between 2010 and 2015," *Nature Human Behaviour* 2, no. 9 (2018): 637-44, https://www.nature.com/articles/s41562-018-0399-z.

12 "Retracted: Risk Management in Financial Institutions," *The Journal of Finance* (June, 2021), https://doi.org/10.1111/jofi.13064; Paul Guest, "Risk Management in Financial Institutions: A Replication," *The Journal of Finance* (June, 2021), https://doi.org/10.1111/jofi.13063.

13 Stefan Nagel, "Answers to FAQ about the recent retraction of an article in the JF," *The Journal of Finance* (July 10, 2021), https://perma.cc/T6AF-7DEU.

14 Chang and Li, "Is Economics Research Replicable?"

15 Tim van der Zee, Jordan Anaya, and Nicholas J. L. Brown, "Statistical heartburn: an attempt to digest four pizza publications from the Cornell Food and Brand Lab," *BMC Nutrition* 3, no. 1 (July 10, 2017): 1-22, https://doi.org/10.1186/s40795-017-0167-x.

16 Kelly Servick, "Cornell nutrition scientist resigns after retractions and research misconduct findings," *Science*, September 21, 2018, https://www.sciencemag.org/news/2018/09/cornell-nutrition-scientist-resigns-after-retractions-and-research-misconduct-finding.

17 Andrew Gelman, "Ethics and statistics: It's too hard to publish criticisms and obtain data for republication," *Chance* 26, no. 3 (2013): 49-52.

18 Helmar Bornemann-Cimenti, Istvan S. Szilagyi, and Andreas Sandner-Kiesling, "Perpetuation of retracted publications using the example of the Scott S. Reuben case: Incidences, reasons and possible improvements," *Science and Engineering Ethics* 22, no. 4 (2016): 1063-72.

19 Roland Sturm, "The effects of obesity, smoking, and drinking on medical problems and costs," *Health Affairs* 21, no. 2 (2002): 245-53.

20 Damon Jones, David Molitor, and Julian Reif, "What Do Workplace Wellness Programs Do? Evidence from the Illinois Workplace Wellness Study," *Quarterly Journal of Economics* 134, no. 4 (2019): 1747–91; Julian Reif *et al*., "Effects of a Workplace Wellness Program on Employee Health, Health Beliefs, and Medical Use: A Randomized Clinical Trial," *JAMA Internal Medicine* 180, no. 7 (July 1, 2020): 952-60, https://doi.org/10.1001/jamainternmed.2020.1321.

21 See https://perma.cc/PZK2-2374. This website also provides information for researchers interested in working with a confidential, non-public use version of our dataset.

22 Jones *et al*., "What Do Workplace Wellness Programs Do?"

### *Publisher's Notes*

Any opinions expressed herein are those of the authors and not necessarily those of the Institute of Government and Public Affairs, the authors' employers, including the University of Illinois Urbana-Champaign, or the University of Illinois System.

**Mast Photographs**

Chicago cityscape - Elena Sivitskaia, stock.adobe.com
Illinois State Capitol Dome - Frame from vidio at https://www.youtube.com/watch?v=F2wPy7DfXfQ
Capitol Dome at Dusk - Frame from Adobe Stock video file 187821651, by VIA Films

**Photography** from istockphoto.com

Pg. 1 - Handheld connections, #1274394138 by ipopba
Pg. 2 - Data graphic, #1218489833 by monsitj
Pg. 3 - System engineers, #1271619512 by metamorworks
Pg. 4 - Abstract code, #1224500457 by monsitj
Pg. 5 - Journals, #1255619435 by Olga Kadygrob
Pg. 6 - Scientists, #1096502340 by sanjeri
Pg. 7 - Blurred motion in office, #1225575965 by AzmanJaka