

Reproducibility and Open Science in Economics

Lars Vilhuber¹

¹Cornell University

July 16, 2025

1 Introduction

As a graduate student in economics at Université de Montréal, reading the economics literature was easy. While the main university library had all the relevant subscriptions, our department librarian, Fethy Mili, would populate the library of the economics department with multi-hued rows of working papers. Mili was also one of the key creators of what was initially known as Working Papers in Economics (WoPEc) and BibEc (for printed working papers) (Krichel 1997; Cruz and Krichel 2000; Krichel and Zimmermann 2009), populating the latter since 1993 (Bátiz-Lazo and Krichel 2012, p. 450). The overall network, known as Research Papers in Economics (RePEc), was born contemporaneously with the more widely known arXiv (Ginsparg 2011) and the more centralized Social Science Research Network (SSRN) (*Social Science Research Network* 2025). While electronic working papers were mostly free in those days (there was no way to pay for them), Mili’s work consisted of sending out postage-paid envelopes to all of the various economics departments that were publishing the working papers, and then cataloging the incoming printed materials electronically, for public consumption. In the end, information about the existence of the working papers was freely available, but access to (printed) working papers still required a small fee to cover the cost of shipping.

I also experienced the openness of code sharing, with code samples by prominent authors being available to graduate students, though discovery was much more difficult at the time. The Statistical Software Components (SSC), primarily but not exclusively for STATA packages, appeared in 1998 (Cox 2010; Cox and Jenkins 2022), providing a convenient and open way to catalog, distribute, and provide open access to additional Stata functionality.¹

Data sharing was harder, of course, with lots of floppies² being exchanged, but also the use of departmental FTP³ servers (David Card’s collection of data, or the NBER’s), or even the replication archive of the Journal of Applied Econometrics, instantiated at Queens University in 1994 under the long-running guidance of James McKinnon.⁴ —On the other hand, administrative data, such as the French administrative data used in AKM required travel to Paris, sitting in a room without windows at an assigned time, and typing the code into the system that had access.⁵

When data were available, computing was straightforward: You logged on to the university’s big computer, running some variant of Unix, and used whatever software was available. Software licenses were paid by the university, as was the computing hardware itself. Laptops powerful (and light!) enough to do actual work were only then emerging.

In 2025, there are concerned discussions about the cost of publishing academic articles, of accessing those same academic articles, of the ever increasing use of administrative data (Card et al. 2010b, 2010a; Chetty 2012; Einav and Levin 2014) that would appear to be hidden behind insurmountable access restrictions, the use of “proprietary software”, and the increasing use of large computing infrastructure, all of which would seem to be restricting access to the basic elements of conducting research in economics. In this article, I will draw on my experience working on many aspects of increasing access to data and materials of all kinds, in particular my recent experience as the inaugural data editor of the American Economic Association (AEA) (Duflo and Hoynes 2018), to paint a picture of economics in an era of open science. How different are matters in practice now compared to that early view of the field of economics, back when I was a graduate student? In this article, I will discuss the current state of open science in economics as facilitated by and related to reproducibility. I will touch on the tension between accessibility, sharing, and preservation, and some of the approaches that are being implemented, sometimes tentatively, in economics, and sometimes elsewhere.

1. This kind of functionality was inspired by similar functionality available for other software, like CPAN for Perl and CRAN for R, but not for most statistical software used by economists.

2. https://en.wikipedia.org/wiki/Floppy_disk

3. https://en.wikipedia.org/wiki/File_Transfer_Protocol

4. The JAE archive was migrated to the ZBW’s archives in 2022 and can now be found at <https://journaldata.zbw.eu/journals/jae/>, but legacy files are still visible as of 2024 at <http://qed.econ.queensu.ca/jae/legacy.html>.

5. For non-local authors, this meant traveling for extended time periods as a Ph.D. student (Margolis) or spending a sabbatical in Paris (Abowd), neither of which is a cheap endeavor. authors, this meant traveling for extended time periods as a Ph.D. student (Margolis) or spending a sabbatical in Paris (Abowd), neither of which is a cheap endeavor. Both were, and still are, enjoyable, though.

My view will be biased - I am an active participant in this space, primarily via my current appointment as data (and reproducibility) editor of the American Economic Association, but also as a past participant in networks that have and foster access, and a researcher and editor in the space of disclosure limitation.

The guiding theme will be the **accessibility** of the key ingredients for scholarship: manuscripts (or more generally, documents), data, software, and the necessary technology to combine the latter two in order to produce knowledge as published in manuscripts. My focus will be on the latter three, though I will provide some observations about scholarly publishing in the last section. In the conclusion, I will identify a few areas where there is (continued) movement towards greater openness.

2 Concepts

In order to write about “Open Science,” a definition is needed. Open science is a surprisingly difficult term to define precisely, and multiple overlapping definition are commonly referenced. UNESCO (2022) sees four components to open **science**: *open scientific knowledge* (publications, data, code, and teaching materials “openly available, accessible and reusable for everyone”), open science *infrastructures* (which encompasses both physical infrastructure such as instruments and laboratories, as well as virtual components such as open access publication platforms), science *communication* (knowledge translation), and broad *engagement* beyond the boundaries of the academy. It also recognizes the limitations of such access in a caveat:

... human rights, security, personal privacy, ... In such cases, it may still be possible to share the existence of such information or share it among certain users who meet defined access criteria.

The Open Knowledge Foundation (Open Knowledge Foundation 2024, OKF) defines “open” as (my emphasis)

... anyone can freely access, use, modify, and share for any purpose (*subject, at most, to requirements that preserve provenance and openness*).

Less broadly, Vicente-Saez and Martinez-Fuentes (2018) identify a consensus that defines open science as “transparent and accessible knowledge that is shared and developed through collaborative networks.”

In this article, I will focus on the what UNESCO 2022 calls open science “knowledge” and will briefly discuss “infrastructures.” I will highlight how some elements have been quite widespread in economics for some time. I will try to identify limits to fully open accessibility, some of which are intrinsic to the nature of the research conducted in economics, and describe how widespread such limitations may be. In particular, I will highlight how those access restrictions are not, as many think, an impediment to **open** science, in the sense that aforementioned “collaborative networks” can still access these resources.

A key ambiguity will arise in how big such networks need to be in order to be considered “open.” ~~Clearly, $n = 2$ is not considered a network~~Consider the realized size of several relevant networks in economics. The National Bureau of Economic Research (NBER) defines its affiliated scholars as a network: $n = 1804$ as of January 2025, primarily in North America (National Bureau of Economic Research 2025). However, in 2024, a total of 2,966 authors published 1,223 NBER working papers. J-PAL has approximately $n = 1725$ affiliates at 120 universities on all populated continents (Abdul Latif Jameel Poverty Action Lab 2025). Between 2001 and June 2024, there had been $n = 2023$ researchers on projects that used confidential U.S. Census Bureau in the Federal Statistical Research Data Centers (FSRDC) (U.S. Census Bureau 2024). In an average year (2013-2023), $n = 1238$ students graduate from a U.S. university with a Ph.D. in economics (National Science Foundation 2024, Table 1-5). In the approximately 30 years since inception of RePEc, $n = 526$ authors from 52 institutions in **Norway** have published a paper listed on RePEc (presumably in economics) (IDEAS//RePEc 2025). All of these are measured across different spatio-temporal dimensions. Are they large? Context and purpose matter. Some may intersect. How many U.S. graduate students in the past 10 years have published an NBER working paper and are now at a Norwegian institution? It is harder to measure the actual key criterion: How many people can potentially enter the network, become part of it, or, as will be the key use of this concept later, how many people can potentially access data of certain types,

accessible via some sort of network. The actual size, relative to the number of potential entrants, is only a proxy for that, and often, a detailed analysis of entry criteria is necessary. For instance, what does it take to enter the network of research data centers (in the US, in France, in Canada, etc.) in order to be able to access the same data as others. Finally, it may be relevant to measure how diverse these networks are. This matters for the ability of insiders of a certain network to criticize each other. For this, size is not a good measure — even $n = 2$ may be sufficient.

Journals play a key role in this space, and will be an important background to my discussion and experiences, possibly also my biases. Most top economics journals have a data (and possibly code) availability policy.⁶ The AEA’s policy was first implemented in 2005 (Bernanke 2004; American Economic Association 2005). While the focus of early availability policies was on the data, the code often came along for the ride, albeit not always in its most complete form. I am the AEA’s inaugural data editor, appointed in 2018 (Duflo and Hoynes 2018). The AEA implemented a new policy in 2019 (American Economic Association 2019, 2019). Many other economics journals appointed data editors around that time, and multiple journals coordinated on a common policy core called Data and Code Availability Standard (DCAS) (Koren et al. 2022), and revised their policies to align with DCAS, see American Economic Association (2024) for the AEA.⁷ A key part of these newer policies was increased pre-publication monitoring of the content of replication packages (Duflo and Hoynes 2018; Christian et al. 2018).

One way to start to move away from a binary perspective of access is to consider **time** as summary metric that captures what is needed in order to access generic resources, whether data, manuscripts, or computing resources. Time might be needed to write an application to access data, or time might be needed in order to obtain access to large-scale computing resources. Time might be needed in order to obtain grant funding that allows to purchase such resources. I choose time, rather than money, as the metric, since it might appear to be slightly more egalitarian, given that much of science has (theoretical) access to subsidies and grants. In the other dimension, the number of people who have some probability of accessing the resource (the **size of the network**) can be taken as an approximate measure of openness, regardless of how interesting or valuable they might consider the data to be. Figure 1, taken from Vilhuber (2023), serves to illustrate this idea, for access to data, with various institutions that facilitate that access mapped out into the space of time vs. size of network. I will return to this throughout the discussion.

3 Data Access

One subcomponent of open science, and locus of much attention throughout the literature in the social sciences, are “open data.” This in principle easy - why should the data used in research not be open? However, the various caveats that policies and principles include are important to recognize. (Open Knowledge Foundation 2024, OKF) mentions “requirements that preserve provenance and openness,” which does not take into account privacy. UNESCO (2022) does note “human rights, security, personal privacy.” On the other hand, even much data that is available to almost anybody on the Internet may not actually be “open.” Consider the S&P 500, viewed in newspaper and many websites (e.g. S&P Dow Jones Indices LLC 2025), is not “open data” because it does not allow for free re-use. OKF defines “open data” as requiring machine readability, absence of licensing charges, and free re-use, but does not mandate availability via download on the internet, absence of all fees, nor absence of any technical measures, such as a requirement to register and agree to abide by these rules (Open Knowledge Foundation 2024).

I will discuss two sub-areas within this space: Secondary data use, and primary data generation. Much of economic research uses data collected by others, such as survey organizations and national statistical offices, but also private company data and various administrative data sources (“organic data”, Groves 2011a, 2011b). Primary data generation is more frequent in behavioral and development economics.

6. For a review of the history of data and code availability policies in economics, see Vlaeminck (2021).

7. These initiatives are not restricted to economics, of course. Political science (*Data Access & Research Transparency (DA-RT)* 2014; Jacoby 2015; Basile, Blair, and Buckley 2023), sociology (*Sociological Science* 2018; Weeden 2023), and general initiatives like the Transparency and Openness Promotion (TOP) guidelines (Nosek et al. 2015).

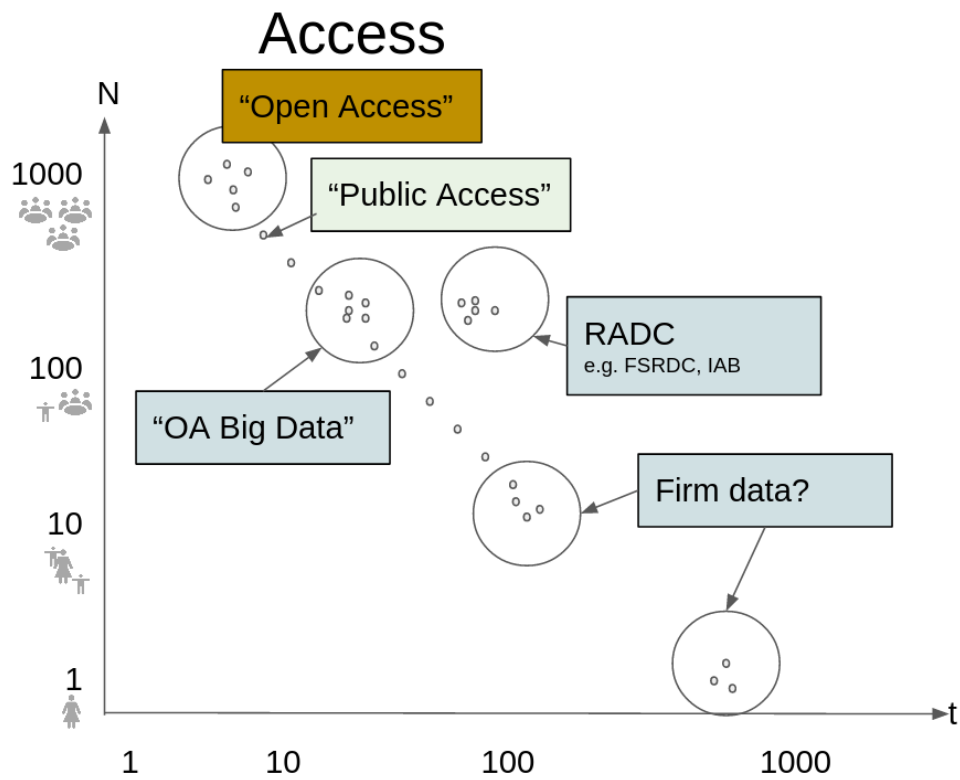


Figure 1: Conceptual trade-off between number of individuals accessing data, and time required to do so. Figure first published in Villhuber (2023).

Secondary data: The data produced by the United States government are in the public domain (i.e., without any restrictions on usage or attribution), which makes it “open” in the above sense (Copyright Act of 1976, [Wikipedia 2025](#)). Many countries have switched their government data to default to open data ([Statistics Canada 2012](#); [UK Government 2014](#)). However, many well known “public-use” data are not always “open”: IPUMS has a redistribution restriction (encapsulated in “Terms of Use”, not a license), and some geographic data by international statistical offices remain under more stringent licensing requirements in other countries (e.g., United Kingdom).

Many well-known survey organizations impose redistribution and usage restrictions that are not consistent with the ~~OKD~~-OKF definition. Many such redistribution restrictions apply to datasets with more detailed personal information collected through surveys, and are meant to ensure continued compliance with ethical rules of behavior, often with informed consent agreements by survey participants and local privacy laws. Notable examples include PSID ([Institute for Social Research 2024](#)), World Value Survey ([Haerpfer et al. 2024](#)), Demographic and Health Surveys (DHS) Program ([DHS Program 2024](#)), German Socio-Economic Panel (Goebel et al. 2019; [Goebel et al. 2024](#)), all of which have broad (cost-free) usage, conditional on registration and compliance with usage restrictions.⁸ Table 1 shows a few examples.

Table 1: Example restrictions

Dataset	Variant	Usage.restrictions.and.justification
DHS		No redistribution and access by unauthorized individuals. To ensure that provider meets host-country agreements, such as compliance with usage by real people with legitimate research purposes.
WVS	Except for joint WVS-EVS files	The data can be used freely for non-commercial purposes such as research, publication, teaching. Data redistribution is prohibited: publication of the original WVS datasets at other online platforms is against the WWSA Constitution.
PSID	Public data	Use the data in the PSID data sets only for scientific research and aggregate statistical reporting; Make no attempts to identify study participants;
PSID	Restricted data	In order to safeguard the confidentiality of respondents at the highest level, some data are provided only under conditions of a restricted use contract, including human subjects review and data security plan.
SOEP	European version	No redistribution and access by unauthorized individuals. Compliance with German Federal Data Protection Act.
SOEP	International version	Same as European version, but 95% sample, and institutional agreement required.

The steps needed to obtain access range from click-through agreements to acknowledge compliance with CC-BY licenses (consistent with open access definitions) to the need to write a paragraph about the purpose on how the data is to be used (maybe consistent), to various commitments to not redistribute the data, all while being cost-free. The first three rows in Table 1 apply some variant of the latter: While anybody can obtain access, the restriction on redistribution is not consistent with definitions of open access. Yet there are no impediments to actually using the data in research.

The last three rows of Table 1, however, go further. Researchers using SOEP data must satisfy certain geographical requirements, such as presence of the researcher in the countries covered by General Data Protection Regulation (GDPR), otherwise they can only access a modified version of the data. Similar restrictions apply to the restricted-use data of many other surveys, for instance, the US-based National Longitudinal Survey of Youth (NLSY). Users of restricted PSID data must comply with even more stringent rules: secure approval from their institution’s ethics board, and use of secure computing environments.

8. I come back to the case of the DHS Program in Conclusion.

Clearly, these restrictions will inhibit broader use of the data per se. Yet these restrictions are not globally very restrictive: There are in Economics alone 1238 new US-based researchers being given the ability to request access to geo-restricted NLSY data every year: newly minted Ph.Ds, as noted earlier. A similar number in the EU likely obtain the right to access the geo-restricted SOEP data as well. Despite the restrictions on the use of SOEP data, there are currently 6,443 papers that in some fashion have used the data.⁹ The PSID bibliography¹⁰ lists over 1,300 dissertations and over 5,300 articles. In the space depicted in Figure 1, access is very much towards the left of the figure.¹¹

Organic data: However, most economists, when asked about data subject to restrictions, will think of “proprietary data,” a term often applied to any data that may be subject to restrictions of use and re-use. Access to most administrative data is typically not “open” in the sense of ~~OKD~~OKF, but are they open enough, given the privacy concerns that are attached to these data? Many have argued that access is not broad enough (Card et al. 2010b; Einav and Levin 2014), while acknowledging the difficulty of addressing privacy and security of the data at scale. Abowd and Schmutte (2018) discuss the challenge of making the choice of between accuracy of (public) statistics and data, and the privacy loss inherent in doing so. Nagaraj, Shears, and Vaan (2020) argue that more openness improves scientific progress, and Nagaraj and Tranchero (2023) study this in the context of the US system for providing access (FSRDC). They note that 4% of US-based empirical authors have had some access to the FSRDC system. I am not aware of similar studies for other countries, such as France (the equivalent system is the Centre d’accès sécurisé aux données (CASD), Gadouche 2019) or Canada (Currie and Fortin 2015). In absolute terms, these networks host several hundred researchers every year. ~~As per (CENSUS) For the 781 projects using Census Bureau data in the FSRDC,~~ 2,084 researchers have had access to ~~projects involving Census Bureau data since (YEAR) confidential data between 1998 and 2024.~~¹² Nagaraj and Tranchero (2023) mention 861 papers ~~in scientific journals.~~ Similar numbers can be obtained for the French (~~xxx-researchers and 6841 researchers from 1109 institutions on 1797 projects with~~ 417 publications)¹³ and Canadian (2201 active researchers ~~from 42 universities~~ as of January 2025, ~~and with~~ 3245 papers published between 2000 and May 2024)¹⁴ networks.¹⁵ The number of publications from access to these networks is smaller in absolute terms than those from PSID and SOEP, though likely higher in impact (Nagaraj and Tranchero 2023).¹⁶

Primary data collection: The discussion of choices made by survey organizations should in principle be applicable when smaller teams of economists, not entire survey organizations, do the primary data collection. Similar to survey organizations, such teams have to balance the privacy of their respondents with the benefits of open science, in particular the broader knowledge to be gained from open access to the data. Many, so it would seem, provide much of the data in replication packages, subject to de-identification (see Kopper, Sautmann, and Turitto 2020; Bjarkefur et al. 2021, for examples), though typically not with stronger disclosure avoidance measures similar to those employed by statistical agencies and larger survey institutions (for a brief discussion of the issues and one possible solution, see Mukherjee et al. 2023). For research teams, ethics boards and institutional review boards (IRBs) have a role to play (Grant and Bouskill 2019), with

9. Source: https://www.diw.de/en/diw_01.c.789503/en/publications_based_on_soep_data_soeplit.html accessed on 2025-02-08.

10. Source: <https://psidonline.isr.umich.edu/publications/Bibliography/search.aspx> accessed on 2025-02-08

11. In fact, authors sometimes forget to abide by the rules for these datasets. As AEA Data Editor, I get notified via “take-down requests” from data providers 12-15 times per year, including in 2024 from the PSID, and have posted information on how to achieve compliance, at least in some cases, at <https://aeadataeditor.github.io/posts/2024-11-01-psid-requests>. Most of the cases affect papers published prior to my tenure, because I do alert authors to data use agreement violations that I am aware of. Ultimately, however, it is the authors’ obligation to remain compliant with such data use agreements.

12. Own calculations based on Census Bureau data, see <https://labordynamicsinstitute.github.io/fsrdc-external-census-projects/>

13. From <https://www.casd.eu/> and <https://www.casd.eu/toutes-les-publications/liste/0/20/> as of May 2025.

14. Provided by Grant Gibson, CRDCN, on February 10, 2025.

15. The Nagaraj and Tranchero (2023) number only includes papers published by economists. ~~All other-Other~~ numbers are counts of researchers and publications in all disciplines, in non-peer-reviewed publications, and include non-economists.

16. For an analysis of code availability over time for SOEP-based publications, see Fink and Marcus (2025).

some arguing very strongly that greater availability to others (though not blind publication of all data) is required in order to maximize the societal benefits that are the *quid pro quo* for the respondents’ consent to their privacy being invaded (Meyer 2018; Grant and Bouskill 2019). Making such data as broadly available, while respecting the privacy of respondents, is precisely what open access to such data promises, *modulo* appropriate access restrictions or data use agreements similar to those outlined in Table 1. In general, however, primary data collections do not have access to robust third-party systems that would allow for access similar to the access required by PSID and similar organizations, situated between no access and fully public access. Thus, while access may be requested in ad-hoc fashion via the original authors, this is known to be fraught with problems (Watson 2022; Gabelica, Bojčić, and Puljak 2022). An ideal scenario would see researchers deposit the data they collected in third-party repositories, which then handle issues such as verifying ethics approval and secure access mechanisms. Some full-service repositories, such as Inter-university Consortium for Political and Social Research (ICPSR) or various national archives, offer such deposits, though they are rarely used by individual researchers in economics. One example are the data collected and used by Ahrsjö, Niknami, and Palme (2024b). The data were collected from public information from the Stockholm District Court. However, in combining individuals’ information into a database, the result was subject to GDPR, and could not be included in the replication package (Ahrsjö, Niknami, and Palme 2024a). However, the authors were able to deposit their pseudonymous data at the Svensk nationell datatjänst (Swedish National Data Service), as a restricted dataset (Ahrsjö, Niknami, and Palme 2023). The SNDS will verify compliance with Swedish law (e.g., GDPR), without requiring future involvement of the original collectors of the data.

How onerous are access restrictions in general? In my work as Data Editor, it matters less than it would seem at first blush. The AEA asks that authors, when submitting, provide some information on how restrictive the data used in the article are, and whether the authors are able to provide the data editor with a not-for-publication copy, for the purpose of verification.¹⁷ Table 2 lists the four levels of restrictions, with some guidance on how to classify specific data. Access via research data centers, described earlier, generally fit within the ‘moderately difficult to obtain’ category. Many private-sector datasets, relying on personal interactions with individuals within the company, fall into the ‘very difficult to obtain’, because others may never be able to obtain them. Also in these lists are access-restricted data that have been discontinued, and are thus no longer available to anybody, such as Zillow’s ZTRAX program, terminated in 2021 (Zillow 2021). Authors can list multiple categories, and we review these, possibly adjusting them before we record them in our internal database.

Table 2: Restriction categories	
Category	Explanation
Very Easy to Obtain	Request takes just a few minutes with no associated costs and the expected response is within a few days
Moderately Easy to Obtain	Request takes less than an hour with minimal cost
Moderately Difficult to Obtain	Request requires a multipage application; request needs university approvals; request involves significant cost; there is some uncertainty as to whether the proposal will be granted
Very Difficult to Obtain	Request must be made in person and/or access is provided only in person; request requires substantial funding; data and/or access mechanism may no longer exist

Figure 2 lists the distribution of the four levels of restrictions plus the no-restriction category, for the 96 AER articles in 2024, interacted with whether the authors offered to share the files with the data editor.¹⁸ Past studies have found that about between 40% and 60% of articles rely on data that are not freely available

17. Data and Code Availability Form

18. Not all offers to share data are taken up, depending on available resources and time. Appendix Table 4 shows the same numbers.

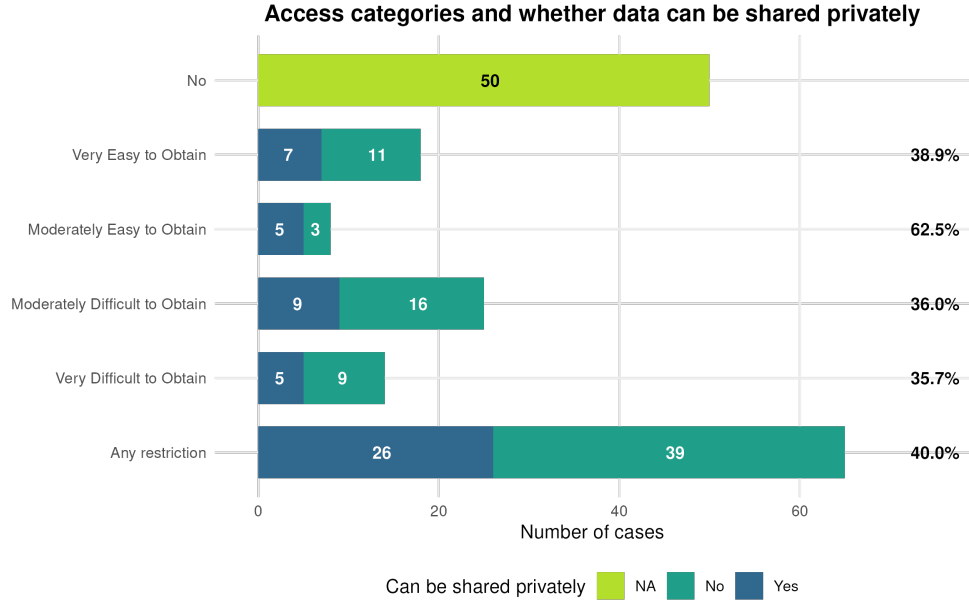


Figure 2: Access categories and whether data can be shared privately

(Herbert et al. 2024; Hamermesh 2013; Chetty 2012). In 2024, slightly more than 50% of all articles had no data restrictions at all (included the data in the replicatin package), with another 19% being ‘very easy to obtain’. The combined percentage of 70% is higher than the earlier long-run average. More importantly, out of the 39 data sources that are difficult to obtain, about a third could be shared with the data editor. It is also worthwhile highlighting that the ‘moderately difficult to obtain’ tends to be data from research data center networks that have fairly robust mechanisms of access, so that in principle and in practice, many researchers do have access to the same data used in these articles, even when the data editor may not have the time to access these data, and when it is not a legal option for the authors to share the data with the data editor. Arguably, the only data sources that are of immediate concern are data where access is uniquely attributable to the authors of the paper, or where access is no longer possible.

4 Access to Software

Turning to software, I will again need to define more precisely what I mean by that. I distinguish two key categories of software: **high-level interpreters** (or more rarely compilers), and **instructions**. The former will in turn comprise software used in two key but distinct features of the scientific production process: **data collection** and **analysis**.

High-level software for analysis are the flagship software that receive the most focus in the social science literature. These are software products such as Stata, MATLAB, R, Julia (see Figure 3). Most of these are interpreted languages at the user level, though user-contributed packages may be compiled (R, Julia). Increasingly present is Python, which may be both compiled and interpreted; less common are purely compiled languages such as C or Fortran. We also include here dedicated geographical information system (GIS) software, mostly used to create maps, though some analytical tasks can also be performed. It arguably may also comprise custom plugins to other software, such as Dynare (Adjemian et al. 2024; Cherrier, Saïdi, and Sergi 2023).

While many economists use data collected by others, some are primary data generators or collectors, for instance through laboratory or field experiments, as well as surveys. **Data collection** software in this

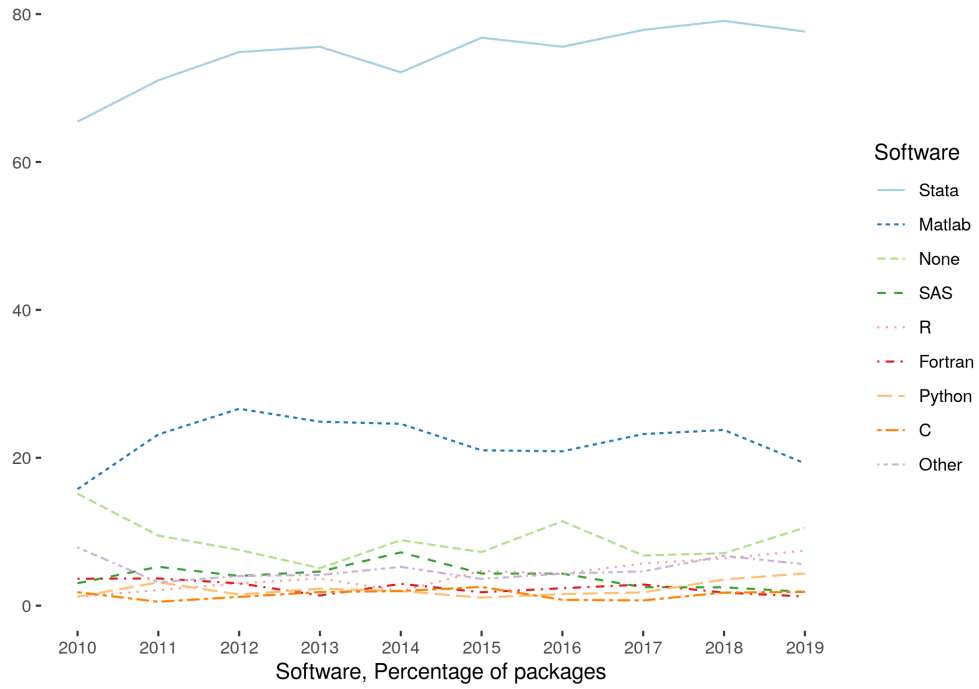


Figure 3: Software usage in AEA journals, 2010-2019, percentage of replication packages. Data do not sum to 100% as a replication package may use several different software. Originally published in Vilhuber, Turitto, and Welch (2020), corrected version see Vilhuber (2020a).

context are the survey software tools (Qualtrics, LimeSurvey, SurveyCTO, KoboToolbox) used to conduct surveys (including as part of lab or field experiments), as well as customized experimental software, such as oTree (Chen, Schonger, and Wickens 2016) and z-tree (Fischbacher, Bendrick, and Schmid 2021).

Software instructions (code): The core of a paper’s analysis, however, resides in how the more general software package is manipulated, in combination with the data, via **instructions**. I do not call these ‘source code’, as that term tends to be used in conjunction with complex compiled software, such as those listed as **high-level software**. Rather, these instructions are mostly interpreted code, or compiled just-in-time, in the language used by the high-level software, or, as is still sometimes the case, in the form of instructions to humans on how to manipulate a non-programmed user interface for software.¹⁹

Open access, therefore, can work through several channels. Software can have a cost. In economics, it is extremely rare to sell access to the **instructions**. However, it is quite normal for **high-level statistical software** and **data collection software** to be accessible only through purchase. As Figure 3 shows, the top two statistical software products used in replication packages are commercial closed-source statistical software products: Stata and Matlab. In order to be able to re-use the instructions (code) from academic papers, a software license is required, and must be purchased.

How much of an impediment to open access is this? Without loss of generality, to illustrate, consider Stata. An academic single-user yearly license as of January 2025 was \$690 for a US-based student. The price is the same for a student in (much poorer) Greece. For a student in Vietnam, the price drops to \$220.²⁰ To put this into perspective, the average monthly incomes in each of these countries, which students are unlikely to earn, are \$6,704, \$1,883, and \$343, respectively, and a reasonably powerful laptop a student is likely to have was around \$1,000.

Table 3: Software licensing internationally

Country	Average.monthly.income	Stata.license	Laptop	Percent.license	Percent.laptop
USA	6,704	690	849	10.3	12.7
Greece	1,883	690	934.8	36.6	49.6
Vietnam	343	220	706.6	64.1	206

Notes: Amounts in USD as of 2025. Source for Stata license prices: Stata website.
Sources for laptop prices: HP.com US and German websites, Shopee.vn for Vietnam.

An informal survey of several economists working with colleagues and students in Latin America, Africa, and Asia uniformly suggested that access to commercial statistical software was often difficult for students, alleviated somewhat when students were fortunate to attend private universities. In fact, in a recent survey (~~CIDR~~)-conducted by CIDR (Sheth et al., n.d.), respondents were asked to name the top two factors that could support early-career African scholars’ publication success. ~~46% of respondents chose ‘Providing access to datasets and data management tools’ as one of the two factors. University staff and students (across multiple disciplines) were also asked which aspect they most needed funding for, and 38% mentioned ‘Data analysis software’.~~ 55% of African scholars noted that providing access to research resources, such as journals, datasets, and software would improve their journal acceptance rate. Within top universities, on the other hand, most researchers, including graduate students, will have access to these software products, and even many undergraduates may be able to leverage university computer labs to access these software. ~~Outside At~~ lower-ranking institutions and outside of academia, however, it may be more constraining, even in Europe or the United States. Precise information is hard to come by, ~~but Figure 4 provides some insights whether this might, in fact, hold. The data plotted relates, by various regions, the relative importance of that region in global Stata and R downloads~~²¹. The top panel summarizes this by “Global North” and “South,” and the

19. In economics, this appears to occur most often for GIS software, and sometimes for data extraction software used by data providers.

20. See source text for data.

21. See Appendix for a detailed explanation of several important caveats.

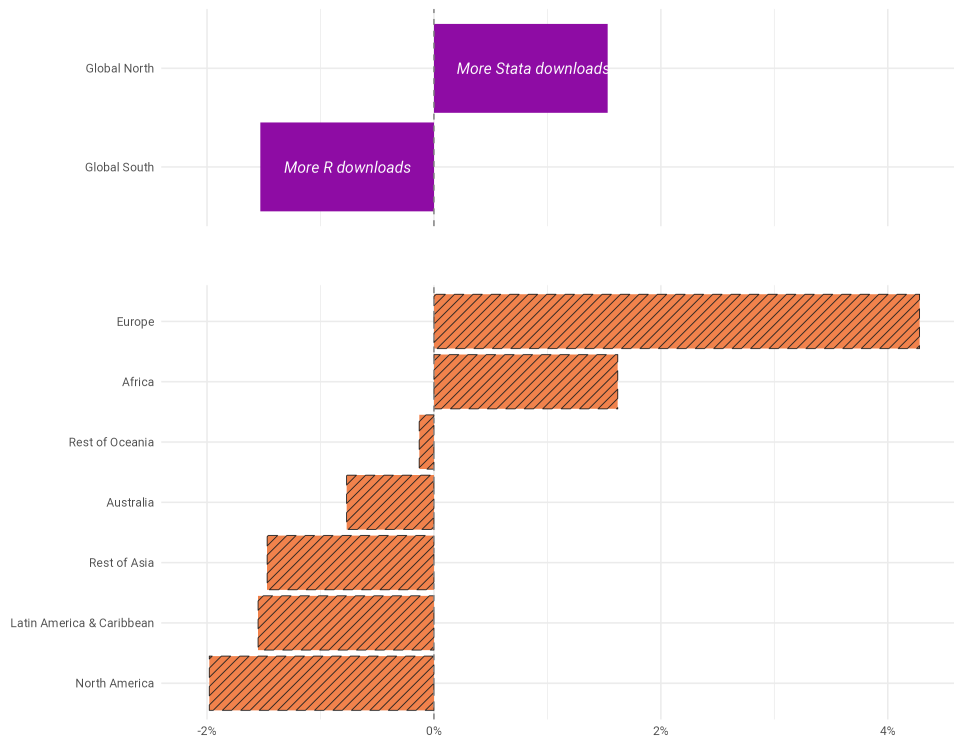


Figure 4: Comparison of relative regional downloads for Stata packages and R software. Source: published data by Kit Baum and Posit PBC, own calculations.

juxtaposition would suggest that the countries of the Global South are heavier users of the free R software, rather than the paid-for Stata software. However, the bottom panel moderates that view somewhat when breaking this out by continents. It turns out that the richest continent (North America) is the heaviest R user, and that Africa is a heavier user of Stata. Thus, there is, at best, very weak evidence that the cost of software matters, in the aggregate.

22

The landscape is even harder to assess for **data collection software**. Open source packages otree and z-tree are often used for lab experiments. otree is cited by between 1,400 (OpenAlex, as of 2025-01-29) and 2,400 articles (Google Scholar, as of 2025-01-29), z-Tree is cited between 9,000 (OpenAlex) and 12,000 times (Google Scholar), but it is not clear how to benchmark that, given poor software citation practices in economics. Lab as well as field experiments will often use online survey platforms to augment experiments, or as primary data collection tool. Qualtrics, one of the big commercial platforms for surveys, is mentioned (but not cited) over 200,000 times (Google Scholar). Open source alternative LimeSurvey is mentioned about 31,000 times.

The second channel is through **software instructions**. Conditional on having the right high-level software and the appropriate resources (see next section), most economics papers have openly accessible code. The advent of data and code availability policies most certainly has helped, but also reflected attitudes

22. To a smaller degree, the problem also arises for **compilers**. While there are open-source compilers for Fortran and C (GNU Fortran, LLVM Flang), many economists will use what they consider to be more efficient commercial compilers, which may require the payment of subscription fees (e.g., Intel, NAG Fortran, PGI (now Nvidia)). More recently, some of the commercial compilers (Intel, Nvidia) have become freely available. Some numerical (non-linear) optimizers (e.g., Gurobi, Knitro) may also require payment of fees, though some limited academic free use may be available.

already present in the researcher community. As noted earlier, the AEA’s earliest policy was announced in 2004 (Bernanke 2004), but the oldest replication package curated by the AEA accompanied an 1999 article (Frankel and Romer 1999a, 1999b) preceding the policy by several years.²³ The newer policies (American Economic Association 2019; Koren et al. 2022) required that more code be provided (starting with raw, rather than cleaned data), and the systematic provision of supplementary materials, such as the survey code to use with data collection software. These policies require that the code be made available in a reasonably liberal license, though no specific open source license is specified. Furthermore, most of the top journal’s repositories of replication packages are at trusted repositories, and not behind journal paywalls. Thus, software instructions in economics, as associated with the scientific output in journals, is generally cost-free, and almost always available under an open access license.

5 Access to Other Resources

In the template README published by myself and several other data editors in economics (Villhuber et al. 2022), we emphasize that a README should provide enough instructions for a reasonable person to re-implement the analysis described in the replication package. Authors need to take into account that cutting edge methods, including technology, may require more information and instruction than for standard methods. Authors likely do not need to describe how to run an analysis in Stata, given its ubiquitous use in economics, and the ease with which instructions can be found more generally, but they may need to provide detailed step-by-step information if the technology used is rare or bespoke. For instance, the emerging use of large language models (LLMs) and artificial intelligence (AI) methods is far from the economic mainstream as of the writing of this article. Recent articles are still identifying ways economists scan actually use these tools, both for personal productivity (Korinek 2023) and as part of the technical toolkit (Athey and Imbens 2019; Dell 2024).

The most recent modern toolkits are not the only resource constraints that might restrict broad access. While it might be argued that the use of proprietary software is restrictive, it is one of many resource constraints that can be binding for some researchers. Researchers in lower-ranked (and lower-funded) research institutions, including in low- and middle-income countries (LMICs), may well not have the funds to purchase proprietary software, but access to computers may be equally constraining. The template README requests information on the type of computer that was used by the original researcher, to provide a benchmark to future re-users. Acquiring access to sufficient memory (random access memory (RAM)), storage, and use over time of those resources can be expensive, even when renting such resources in cloud environments (which very few researchers appear to be doing). Traditionally, that access may be embedded within a single purchased computer, which may have (in 2024) around 32GB of RAM, 1-2 TB of storage, and have 4-12 compute cores available exclusively to the owner. More complex analyses may require access to shared compute clusters (using hundreds or thousands of compute cores), very large storage arrays (measured in the two- to three-digit TB range), and may require up to 1024 GB of RAM. Cutting edge analyses may require specialized chips, such as one or more graphical processing units (GPUs), or even a cluster of GPUs. I have observed analyses that may run data cleaning or data acquisition processes for months at a time.

The vast majority of articles published in economics journals usually require no more computing resources than a modern laptop provides, in all the dimensions enumerated in the previous paragraph. In fact, a formal quantitative measurement of resource usage in economics articles is surprisingly hard to obtain, as most researchers are not very good at reporting the resources they have used to conduct their research. In part, this is because measuring such usage is non-trivial, but to a larger extent, I postulate that this is because most research institutions provide such resources to their researchers in a “convenient way,” and researchers conduct research within those constraints. More importantly, however, it suggests an important constraint on how “open” access can be for some if not all economics research.

23. The oldest replication package in economics may be Koenker (1988) in the JAE replication archive, see <https://aeadataeditor.github.io/posts/2023-02-02-oldest-replication-package-jae>, which actually contained data processing code, but not the data analysis code. That was only published in a later replication package, Koenker and Zeileis (2009).

Some newer research requires vastly different types of resources. Studies using raw satellite data may require more than 10TB of data storage (Khachiyan et al. 2022b; Khachiyan et al. 2022a), may need more than 20,000 compute hours on a cloud provider (Rudik 2020b, 2020a), or the use of one (Dell 2024, 2025) or dozens (Khachiyan et al. 2022a) GPUs.²⁴ Access to the code and data for the papers mentioned is open: The AEA-related replication packages for these articles are licensed under a standard Creative Commons Attribution (CC-BY) license. Some of the data not included in the Dell (2025) replication package is on Huggingface, also under a CC-BY license (Silcock et al. 2024). Open access to satellite data is one of the canonical examples of the benefits of open access (Nagaraj, Shears, and Vaan 2020). In these cases, the computational resources may restrict the benefits of the open access of data and code.

Are such computational constraints a problem? No consistent analysis exists that correlates resource requirements to academic outcomes such as citations, primarily because it is very hard to measure consistently the resource requirements of economic articles. The very small sample in the previous paragraph may serve to illustrate this, but without controls for scientific merit, is purely an indicator. Silcock et al. (2024) had been downloaded 98 times in December 2024, six months after the arXiv paper associated with it was published (Silcock et al. 2024). Rudik (2020a) has had 1432 views, 124 downloads for replication package, as the manuscript (Rudik 2020b) has 15 citations. The replication package Khachiyan et al. (2022a) has 2124 views, 175 downloads, while the manuscript (Khachiyan et al. 2022b) has 5 citations (all as of January 2025). For comparison, the average article in one of the AEA’s journals has 908 views and 106 downloads (Vilhøber 2025, Table 4).

6 The Benefits of Open Science in Economics

~~To illustrate some of the benefits of~~ The discussion about benefits of open science often centers around data availability. The World Bank, in its annual World Development Report, identifies data availability for research, for commerce, for education as a key contributor, and highlights that many LMIC continue to have impediments to the reliable provision of open access data (World Bank 2021, , pg. 62). The (theoretical) optimal level of data availability intersects with privacy, making the optimal level of data availability a non-trivial balance between public and private benefits, and private costs (Duch-Brown, Martens, and Mueller-La). The more recent discussions (and court cases) surrounding the use of data in the training of large language models have only re-emphasized this tension (Panettieri 2025).

A different thread in the discussion brings up normative reasons for transparency, for instance around Mertonian (Merton 1942) norms of openness (see Miguel 2021, for an overview). Openness at all stages of research may act as a moderator for publication bias (Miguel 2021; Brodeur et al. 2016).

Some of the benefits of open science have been measured indirectly, through increased citations, say. Some recent studies find some advantages for studies with linked (openly available) data (Piwowar and Vision 2013; Colavizza et al. 2025). The economics literature has emphasized the benefits of (balanced) open science. Nagaraj, Shears, and Vaan (2020) discuss the canonical example of improved data access through (free) public-use data in the context of satellite imagery. Patents can be usefully investigated, since they are both openly viewable but also access-limiting by their very nature. Economists have looked at how restrictiveness, duration, and type of patents affect scientific progress. In general, restrictions reduce social welfare (Williams 2013; Murray et al. 2016), even enabling anti-competitive behavior that directly identify welfare loss (Xie and Gerakos 2020). The much broader applicability of open science practices is also much more recent, and as of yet, hard to measure. Nevertheless, it appears to be widely accepted in economics, as evidenced by very strong positive attitudes documented in Ferguson et al. (2023), as selectively depicted in Figure 5. The top left part of Figure 5 shows behavior (black) and opinions (color scale) in regards to data sharing among economists, with more than 50% of economists having shared data, but over 90% being very much or moderately (the two green colors) in favor of sharing data. The right panel illustrates the evolution over time, depicting the proportion

24. As of January 2025, the type of GPU used by Dell (2025), costs between USD 4500 and USD 7700, or between 2 and 7 times as much as a standard laptop.

25. Some of these studies need to be taken with a grain of salt, since they typically measure whether data is referenced, not whether it is actually openly available.

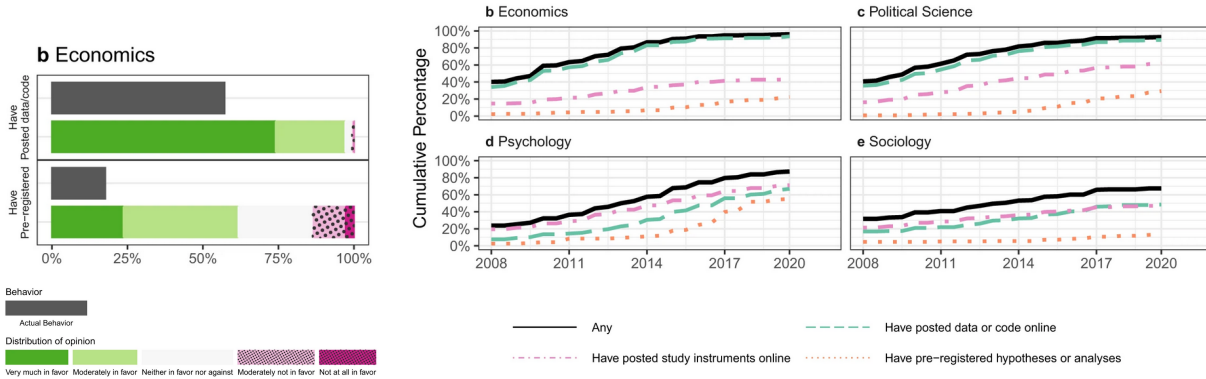


Figure 5: [Extracts from Ferguson et al. \(2023, Figures 1 and 7\), reused under CC-BY 4.0.](#)

of social scientists in the four disciplines (economics, political science, psychology, and sociology) who had adopted an open science practice as of a particular date, showing again the rapid increase in active data sharing amongst economists from around 60% in 2011 to the 2020 number of over 90%.

I want to add to this discussion of benefits three concrete case studies that advance science in economics, while relying heavily on the openness of prior research. The three articles in question leverage the open science, I describe two recent articles in economics that leverage the availability of code and data, with licenses that allow for re-use, to improve econometric methods for future researchers. The first paper relies on the empirical recomputation of prior papers to assess a theoretically ambiguous potential bias in inference. The second paper also focuses on inference, and uses actual data from previous papers to simulate the relevance of the impact. Both then provide new software (R and Stata packages), under open source licenses, to “fix” the problem for future studies. The third study selects studies again where data are available, and leverages the centralized availability of such replication packages to select the studies in their sample.

Roth (2022) uses 12 previously published papers to assess whether the usual tests for pre-existing differences in trends when using difference-in-differences methods are properly controlling for power, and the empirical impact on subsequent inference. The theoretical bias is ambiguous, so an empirical evaluation is necessary. The study both leverages the open availability of materials in economic journals, but also illustrates the limitations imposed by imperfect adherence to openness. To wit, Roth writes that he searched for

“the phrase “event study” in papers published in the *American Economic Review*, *American Economic Journal: Applied Economics*, and *American Economic Journal: Economic Policy* between 2014 and June 2018 ... The search returned 70 total papers that include a figure that the authors describe as an event-study plot.”

but continues to then be limited by lack of data in the majority of cases:

“I exclude 43 papers for which data to replicate the main event-study plot were unavailable. (Roth 2022, pg. 307)”²⁶

While it remains unclear whether the excluded papers are non-compliant with the AEA’s policy at the time, or whether they have legitimate reasons not to provide the data (Roth does not provide the raw result of his search), the paper is able to make an important methodological point (723 / 415 citations as of January 2025, per Google Scholar/ OpenAlex) because it is able to fully recompute the results in previous papers, apply new tests and methodologies, and come to meaningful recommendations and tools — Roth provides an (open source) R package to implement his methodology.

26. He also excludes another 15 papers for reasons not related to data availability.

Chaisemartin and Ramirez-Cuellar (2024) use open access information on RCTs (AEA registry) to find 15 RCTs of a particular type (clustered paired or small strata), of which 4 have publicly available data (and reproducible artifacts). They then provide results both on simulations using these data, and in particular, re-estimate the regressions used in those studies and apply their proposed solution, showing that the number of significant effects is reduced by one-third. In other work (Chaisemartin and D’Haultfoeulle 2020, 2024), the results from various other papers are also recomputed to empirically demonstrate the relevance of the proposed methods, and software packages (e.g. Chaisemartin et al. 2025) are developed and made openly available.²⁷ Chaisemartin and D’Haultfoeulle (2020) has been cited between 2,600 and 4,600 times (OA, GS).

~~In both~~ Finally, Goldsmith-Pinkham, Hull, and Kolesár (2024) investigate contamination bias (“each treatment’s effect are contaminated by nonconvex averages of the effects of other treatments”). To do so, they search the centralized repository of AEA replication packages²⁸ to identify packages that, crucially, contain data. They then reproduce one of each original paper’s specifications, conduct several tests, and conclude that “economically and statistically meaningful contamination bias [is present] in two of the three observational studies while showing no evidence for bias in any of the experimental studies.” (pg. 4046). Goldsmith-Pinkham, Hull, and Kolesár (2024) has been cited between 65 and 118 (OA, GS) times.²⁹

~~In each~~ of these examples, the ability to access prior data, code, and information is critical to improving future scientific progress, ~~but is~~. ~~In some cases, it remains~~ limited by both historical and unavoidable limitations on openness, ~~as well as scaling limitation based on absence of “push button” reproducibility.~~

~~“These studies were identified by a systematic search of papers in the AEA Data and Code Repository”~~

7 Open Infrastructure: Publications

The challenges of openly accessible written scholarship are manifold, with the current focus on “Plan S”, master publication agreements, and in the US, similar efforts under the moniker of the “Nelson memo” (Brainard and Kaiser 2022; Brainard 2024). I note that the economics profession has a very long history of making much of the written knowledge available at very low cost via working papers (Vilhøber 2020b), with the first working papers at the reputable NBER working paper series going back to 1973 (Welch 1973).

~~At the time of this writing~~ Over the past eight years, I have ~~or have had~~ three editorial appointments. I am the Data Editor for the journals of the American Economic Association (AEA) ~~the~~ (Duflo and Hoynes 2018), ~~a column editor for the open access~~ Harvard Data Science Review (HDSR) (Vilhøber et al. 2023), and until January 2025, the joint executive editor for the open access and multi-disciplinary Journal of Privacy and Confidentiality (JPC), ~~and I am a column editor for the open access for which I continue to manage the publication infrastructure~~ (J. Abowd et al. 2025). I will use each of these to highlight a particular pattern in broadening access to publications, without any claim to generality.

The AEA is a not-for-profit organization, as are many other learned societies. It self-publishes eight journals, plus the proceedings of the annual conference, without relying on a commercial publishing house. Depending on the measurement, three of these publications are in the top ten journals in economics (Mogstad et al. 2022). Its publication costs account for about half of its overall operating expenses, and are only partially offset by directly attributable subscription and membership fees (Cherry Bekaert, LLP 2024). In fact, 6 of the top 10 journals in economics (Mogstad et al. 2022) are published by societies (JEL, JEP, Econometrica, AER, Restud, JOLE), some of which have as sole or primary purpose the publication of the journal.³⁰ A further three journals are primarily associated with economics departments (QJE, JPE,

27. Note that as of January 2025, many of the packages do not have an explicit open source license — or any license — applied, a common feature of economists working in the open source world. My presumption is that they simply assume that everybody knows that the code is openly available.

28. “These studies were identified by a systematic search of papers in the AEA Data and Code Repository.” [pg. 4043]

29. Google Scholar citations are directly reported from a view of the article’s page on Google Scholar, and may include citations to multiple versions. OpenAlex citations are the sum of citations to all recorded works on OA with the same title and by the same authors.

30. JOLE is a bit of an outlier, in that one becomes a member of the Society of Labor Economists by subscribing to the journal, rather than the other way around.

RESTAT), which arguably may not be driven by pure profit. The sole outlier in the top ten is the Journal of Financial Economics (JFE), which is owned by Elsevier, a big commercial publisher. It should be noted that the European Economics Association (EEA) severed its relationship with Elsevier in 2003 for its official journal, creating a journal that is fully owned by the association, adding to the list of society-owned journals in economics (Tirole et al. 2003). Access to these journals is generally still on a subscription basis (JEP is the exception, being free to read), but given the primarily not-for-profit organization of its owners, personal subscriptions (often via society membership) are quite low, compared to journals in many other sciences. For instance, as of 2024, a personal subscription to the Review of Economic Studies is \$156 or €141 per annum; a yearly membership to the AEA, providing access to the seven subscription journals and the proceedings is \$25 for students and researchers in low-income countries, and \$150 at the highest personal income tier. As outlined earlier for the AEA, these subscription fees cover only a small fraction of the production costs. Nevertheless, even these (arguably low) costs do not satisfy “Plan S” or “Nelson memo” requirements, which require no access cost to the end consumer, and in the case of “Plan S”, also require a liberal license allowing for re-use.³¹

Interestingly in the context of the previous sections, all of the society-owned journals in the previous paragraph have appointed data (reproducibility) editors.³²

Since 2018, I have been the executive editor of the JPC, an open access multi-disciplinary journal, having taken over the journal from Stephen Fienberg (Vilhuber 2018).³³ As of 2024, the journal does not charge submission fees, and is free to read (what is called “diamond open access.” Articles default to a Creative Commons Attribution-NonCommercial-NoDerivatives (CC-BY-NC-ND) license, though authors are allowed to choose a more liberal license, for instance to comply with “Plan S” (which does not allow for the “no-derivatives” part). As executive editor, I have been responsible for all aspects of running the journal, not just finding referees for the articles that I am responsible for. The journal is made available through open-source software called Open Journal System, hosted by its creators at Simon Fraser University’s Public Knowledge Project, preserved via industry-standard mechanisms (CLOCKSS, a non-profit) in case the journal ever needs to shut down, indexed in a variety of academic indexes, including via assignment of DOI. Copy-editing is done through a mixture of professional copy-editors and volunteer work by editors and board members. All editors, including myself, are unpaid, and referees are, like in much of the publishing industry, unpaid volunteers. Yet I do pay bills, for each of the above components of a properly managed, indexed, and preserved academic journal — and professional copy-editors and university staff do not work for free. I am thus quite aware of the absolute minimum cost of running a (small) journal. Over the years, funding has come from a variety of chaired professorships at Carnegie Mellon (Fienberg), Cornell (Abowd), and Harvard (Dwork). In order to make such funding more robust, a non-profit society was created to better and more robustly structure the funding situation (J. M. Abowd et al. 2024). Time will tell if this will stabilize the funding situation, while maintaining the foundational commitment to open access. Others, in particular Sociological Science, have shown that it is feasible to sustainably publish high-quality research

8 Conclusion

I have described in this article how open data and code are in the academic literature in economics, and how access to software and hardware resources can be limiting factors. Almost all code is openly accessible in top journals. The vast majority of data is accessible with little to no effort, and a large proportion of the remaining restricted data can be accessed by networks that include thousands of researchers. I provide a few concrete examples where the openness of the data and code available allows others to directly build on prior

31. “The author(s) [...] grant(s) [...] a free, irrevocable, worldwide, right of access to, and a license to copy, use, distribute, transmit and display the work publicly and to make and distribute derivative works, [...] subject to proper attribution of authorship” (Max-Planck-Gesellschaft 2023).

32. Two additional societies, not previously mentioned, also employ data editors: the Canadian Economics Association (CEA) and the Western Economics Association International (WEAI).

33. Fienberg, together with Cynthia Dwork and Alan Karr, founded the journal in 2009 (Abowd, Nissim, and Skinner 2009). Fienberg passed away in 2016 (Slavković and Vilhuber 2018). In 2023, I recruited Rachel Cummings to jointly manage the journal.

results. Broader assessments of the benefits of open access are more difficult to measure, in part because the right controls are hard to construct, in part because the community is still only starting to learn how to technically leverage the openness in a large-scale fashion.

Nevertheless, access to networks of data access and financing remains one of the key worries. I am a regular participant in discussions within three research networks that provide access to restricted data (FSRDC, Canadian Research Data Center Network (CRDCN), and CASD). Core discussions center around equity and access, and how to balance those criteria while preserving the privacy of the respondents for which these networks act as curators.

One under-appreciated aspect of open access is that it enables persistence. For any data that is subject to a gatekeeping mechanism, however objective, impartial, and lightweight it may be at the present time, such mechanisms can and do disappear. Every one of the networks that I mention above has a mandate to work within budget constraints, and those budgets are determined fundamentally by external forces, typically government-based funding agencies. A particular striking example, as of the writing of this article, is access to data from the USAID-funded DHS program. Data access to DHS data was classified as ‘moderately easy to obtain’ (as per Table 2) until February 2025, as it took only a day or two to obtain access subject to a lightweight data use agreement. My team at the AEA regularly went through the process to obtain data used by other researchers, a process that was easy to navigate even for an undergraduate researcher on my team. In February 2025, the second Trump administration shut down USAID and “paused” the DHS program, with very little notice, and no recourse. While the DHS program system was still accessible to those with prior access permission in late February 2025, no further access requests were accepted. That turns it into a ‘very difficult to obtain’ dataset.³⁴ These events are a note of caution that any kind of redistribution restriction may very well negatively impact future availability to the research community. Whereas journals can subscribe to mechanisms that allow past issues to remain accessible even when the journal is shuttered,³⁵ and open access software can be preserved by communities with an interest in continued use,³⁶ no such mechanism exists for most restricted-access data. While open code may ensure we can recompute, and advances in computational infrastructure will bring the current cutting edge into the space of everyday-accessibility, data that are not open can and will disappear. That is concerning.

References

- Abdul Latif Jameel Poverty Action Lab. 2025. *Affiliated Professors*. Accessed January 17, 2025. <https://web.archive.org/web/20250117031829/https://www.povertyactionlab.org/affiliated-professors/all>.
- Abowd, John, Cynthia Dwork, Alan Karr, Jerome Reiter, and Aleksandra Slavković. 2025. “Changes to the Journal of Privacy and Confidentiality.” Publisher: Society for Privacy and Confidentiality Research, *Journal of Privacy and Confidentiality* 15 (1). Accessed July 14, 2025. <https://doi.org/10.29012/jpc.991>.
- Abowd, John M., Cynthia Dwork, Alan F. Karr, Kobbi Nissim, Jerome Reiter, Aleksandra Slavković, and Lars Vilhuber. 2024. “Launching the Society for Privacy and Confidentiality Research to Own the Journal of Privacy and Confidentiality.” Number: 3, *Journal of Privacy and Confidentiality* 14 (3). Accessed January 1, 2025. <https://doi.org/10.29012/jpc.939>.
- Abowd, John M., Kobbi Nissim, and Chris J. Skinner. 2009. “First Issue Editorial.” Number: 1, *Journal of Privacy and Confidentiality* 1 (1). Accessed January 1, 2025. <https://doi.org/10.29012/jpc.v1i1.562>.
- Abowd, John M., and Ian M. Schmutte. 2018. *An Economic Analysis of Privacy Protection and Statistical Accuracy as Social Choices*. Document 48. Labor Dynamics Institute, Cornell University. <https://hdl.handle.net/1813/58669>.

34. As of March 2025, efforts are underway to preserve the DHS program data, for instance via IPUMS, and to resuscitate both the access to historical data, as well continue funding new surveys.

35. See CLOCKSS.

36. See f.i. Data Rescue Project.

- Abowd, John M., Ian M. Schmutte, William N. Sexton, and Lars Vilhuber. 2019. “Suboptimal Provision of Statistical Accuracy when it is a Public Good.” *submitted*.
- Acquisti, A., C. Taylor, and L. Wagman. 2016. “The Economics of Privacy.” *Journal of Economic Literature* 54 (2): 442–492. <https://doi.org/10.1257/jel.54.2.442>.
- Adjemian, Stephane, Michel Juillard, Frederic Karame, Willi Mutschler, Johannes Pfeifer, Marco Ratto, Normann Rion, and Sebastien Villemot. 2024. *Dynare: Reference Manual, Version 6*. Working Paper 80. CEPREMAP. Accessed January 26, 2025. <https://www.dynare.org/wp-repo/dynarewp080.pdf>.
- Ahrsjö, Ulrika, Susan Niknami, and Mårten Palme. 2023. *Data för: Identity in Court Decision-Making*. Artwork Size: 39.6 MiB, 64 variables, 15456 cases Pages: 39.6 MiB, 64 variables, 15456 cases. Accessed February 9, 2025. <https://doi.org/10.58141/27J7-8606>.
- . 2024a. *Code for Identity in Court Decision-Making*. Accessed February 9, 2025. <https://doi.org/10.3886/E193364V1>.
- . 2024b. “Identity in Court Decision-Making.” *American Economic Journal: Economic Policy* 16 (4): 142–164. Accessed February 9, 2025. <https://doi.org/10.1257/pol.20220802>.
- American Economic Association. 2005. *Data and Code Availability Policy (2005-2019)*. Accessed February 5, 2025. <https://www.aeaweb.org/journals/data/archive/2005>.
- . 2019. “Updated AEA Data and Code Availability Policy.” Tex.nonote: (accessed: 2019-12-08 via Archive.org) tex.notusedurl: <https://www.aeaweb.org/news/member-announcements-july-16-2019> tex.timestamp: 2019-09-21T23:07:30Z, *AEA Member Announcements: Updated AEA Data and Code Availability Policy (July 16, 2019)*, accessed September 21, 2019. <https://web.archive.org/web/20191208160745/https://www.aeaweb.org/news/member-announcements-july-16-2019>.
- . 2024. *Data and Code Availability Policy*. Accessed February 5, 2025. <https://www.aeaweb.org/journals/data/data-code-policy>.
- Athey, Susan, and Guido W. Imbens. 2019. “Machine Learning Methods That Economists Should Know About.” *Annual Review of Economics* 11 (1): 685–725. Accessed December 30, 2024. <https://doi.org/10.1146/annurev-economics-080217-053433>.
- Basile, Linda, Alasdair Blair, and Fiona Buckley. 2023. “Research transparency and openness.” *European Political Science* 22 (2): 177–181. Accessed March 11, 2025. <https://doi.org/10.1057/s41304-023-00424-x>.
- Bátiz-Lazo, Bernardo, and Thomas Krichel. 2012. “A brief business history of an on-line distribution system for academic research called NEP, 1998-2010.” *Journal of Management History* 18 (4): 445–468. Accessed July 16, 2018. <https://doi.org/10.1108/17511341211258765150788>.
- Bernanke, Ben S. 2004. “Editorial Statement.” Publisher: American Economic Association, *The American Economic Review* 94 (1): 404–404. Accessed September 1, 2020. <https://www.jstor.org/stable/3592790>.
- Bjarkefur, Kristoffer, Luiza Cardoso de Andrade, Benjamin Daniels, and Maria Ruth Jones. 2021. *Development Research in Practice: The DIME Analytics Data Handbook*. Handbook. Accepted: 2021-05-19T14:47:41Z ISBN: 9781464816949. Washington, DC: World Bank. Accessed March 30, 2022. <https://doi.org/10.1596/978-1-4648-1694-9>.
- Brainard, Jeffrey. 2024. *U.S. science funding agencies roll out policies on free access to journal articles*. Accessed December 30, 2024. <https://www.science.org/content/article/u-s-science-funding-agencies-roll-out-policies-free-access-journal-articles>.
- Brainard, Jeffrey, and Jocelyn Kaiser. 2022. *White House requires immediate public access to all U.S.-funded research papers by 2025*. Accessed December 30, 2024. <https://www.science.org/content/article/white-house-requires-immediate-public-access-all-u-s--funded-research-papers-2025>.

- Brodeur, Abel, Mathias Lé, Marc Sangnier, and Yanos Zylberberg. 2016. “Star Wars: The Empirics Strike Back.” *American Economic Journal: Applied Economics* 8 (1): 1–32. Accessed June 14, 2020. <https://doi.org/10.1257/app.20150044>.
- Card, David E., Raj Chetty, Martin S. Feldstein, and Emmanuel Saez. 2010a. “Expanding Access to Administrative Data for Research in the United States.” *SSRN Electronic Journal*, accessed August 15, 2024. <https://doi.org/10.2139/ssrn.1888586>.
- . 2010b. *Expanding Access to Administrative Data for Research in the United States (long version)*. Mimeo. UC Berkeley. Accessed August 15, 2024. <https://eml.berkeley.edu/~saez/card-chetty-feldstein-saezNSF10dataaccess.pdf>.
- Chaisemartin, Clément de, and Xavier D’Haultfoeuille. 2020. “Two-Way Fixed Effects Estimators with Heterogeneous Treatment Effects.” *American Economic Review* 110 (9): 2964–2996. Accessed February 8, 2025. <https://doi.org/10.1257/aer.20181169>.
- . 2024. “Difference-in-Differences Estimators of Intertemporal Treatment Effects.” *The Review of Economics and Statistics*, 1–45. Accessed February 9, 2025. https://doi.org/10.1162/rest_a_01414.
- Chaisemartin, Clément de, Xavier D’Haultfoeuille, Diego Ciccía, Felix Knau, Méline Malézieux, and Doulo Sow. 2025. *chaisemartinPackages/did_multiplet_dyn*. Technical report. Original-date: 2023-12-20T15:28:30Z. Github. Accessed February 8, 2025. https://github.com/chaisemartinPackages/did_multiplet_dyn.
- Chaisemartin, Clément de, and Jaime Ramirez-Cuellar. 2024. “At What Level Should One Cluster Standard Errors in Paired and Small-Strata Experiments?” *American Economic Journal: Applied Economics* 16 (1): 193–212. Accessed February 8, 2025. <https://doi.org/10.1257/app.20210252>.
- Chen, Daniel L., Martin Schonger, and Chris Wickens. 2016. “oTree—An open-source platform for laboratory, online, and field experiments.” *Journal of Behavioral and Experimental Finance* 9:88–97. Accessed January 26, 2025. <https://doi.org/10.1016/j.jbef.2015.12.001>.
- Cherrier, Béatrice, Aurélien Saïdi, and Francesco Sergi. 2023. ““Write Your Model Almost as You Would on Paper and Dynare Will Take Care of the Rest!” A History of the Dynare Software.” *OEconomia*, nos. 13-3, 801–848. Accessed January 26, 2025. <https://doi.org/10.4000/oeconomia.16123>.
- Cherry Bekaert, LLP. 2024. *THE AMERICAN ECONOMIC ASSOCIATION: Financial Statements and Supplementary Information*. Report of the Independent Auditor. American Economic Association. Accessed December 30, 2024. <https://www.aeaweb.org/content/file?id=20801>.
- Chetty, Raj. 2012. *Time Trends in the Use of Administrative Data for Empirical Research*. Accessed July 19, 2018. http://www.rajchetty.com/chettyfiles/admin_data_trends.pdf.
- Christensen, Garret, Allan Dafoe, Edward Miguel, Don A. Moore, and Andrew K. Rose. 2019. “A study of the impact of data sharing on article citations using journal policies as a natural experiment.” Publisher: Public Library of Science, *PLOS ONE* 14 (12): e0225883. Accessed October 7, 2021. <https://doi.org/10.1371/journal.pone.0225883>.
- Christian, Thu-Mai, Sophia Lafferty-Hess, William Jacoby, and Thomas Carsey. 2018. “Operationalizing the Replication Standard: A Case Study of the Data Curation and Verification Workflow for Scholarly Journals.” *International Journal of Digital Curation* 13 (1). Accessed June 4, 2018. <https://doi.org/10.2218/ijdc.v13i1.555>.
- Colavizza, Giovanni, Iain Hrynaszkiewicz, Isla Staden, Kirstie Whitaker, and Barbara McGillivray. 2020. “The citation advantage of linking publications to research data.” Tex.ids= colavizza2020 publisher: Public Library of Science, *PLOS ONE* 15 (4): e0230416. Accessed October 7, 2021. <https://doi.org/10.1371/journal.pone.0230416>.

- Cox, Nicholas J. 2010. "A Conversation with Kit Baum." *The Stata Journal* 10 (1): 3–8. Accessed June 14, 2020. <https://doi.org/10.1177/1536867X1001000102>.
- Cox, Nicholas J., and Stephen P. Jenkins. 2022. "The Stata Journal Editors' Prize 2022: Christopher F. Baum." *The Stata Journal: Promoting communications on statistics and Stata* 22 (4): 727–733. Accessed November 8, 2024. <https://doi.org/10.1177/1536867X221140932>.
- Cruz, Jose Manuel Barrueco, and Thomas Krichel. 2000. "Cataloging Economics Preprints." *Journal of Internet Cataloging* 3 (2-3): 227–241. Accessed July 16, 2018. https://doi.org/10.1300/J141v03n02_08.
- Currie, R.F., and S. Fortin. 2015. *Social Statistics Matter: A History of the Canadian RDC Network*. Tex.ids: Currie2015. Canadian Research Data Centre Network = Réseau canadien des Centres de données de recherche. <http://rdc-cdr.ca/sites/default/files/social-statistics-matter-crdcn-history.pdf>.
- Data Access & Research Transparency (DA-RT)*. 2014. Accessed March 11, 2025. <https://www.dartstatemnt.org>.
- Dell, Melissa. 2024. *Deep Learning for Economists*. Working Paper. Accessed December 7, 2024. <https://doi.org/10.3386/w32768>.
- . 2025. *Data and Code for: Deep Learning for Economists*. Accessed February 7, 2025. <https://doi.org/10.3886/E210922V1>.
- DHS Program. 2024. *Demographic and Health Surveys*. United States Agency for International Development (USAID). <https://dhsprogram.com/data/Access-Instructions.cfm>.
- Duch-Brown, Néstor, Bertin Martens, and Frank Mueller-Langer. 2017. *The Economics of Ownership, Access and Trade in Digital Data*. SSRN Scholarly Paper. Rochester, NY. Accessed May 16, 2025. <https://doi.org/10.2139/ssrn.2914144>.
- Duflo, Esther, and Hilary Hoynes. 2018. "Report of the Search Committee to Appoint a Data Editor for the AEA." *AEA Papers and Proceedings* 108:745. Accessed July 22, 2018. <https://doi.org/10.1257/pandp.108.745>.
- Einav, Liran, and Jonathan Levin. 2014. "Economics in the age of big data." *Science* 346 (6210): 1243089. Accessed September 12, 2021. <https://doi.org/10.1126/science.1243089>.
- Ferguson, Joel, Rebecca Littman, Garret Christensen, Elizabeth Levy Paluck, Nicholas Swanson, Zenan Wang, Edward Miguel, David Birke, and John-Henry Pezzuto. 2023. "Survey of open science practices and attitudes in the social sciences." Number: 1 Publisher: Nature Publishing Group, *Nature Communications* 14 (1): 5401. Accessed September 6, 2023. <https://doi.org/10.1038/s41467-023-41111-1>.
- Fink, Lukas, and Jan Marcus. 2025. "Replication code availability over time and across fields: Evidence from the German Socio-Economic Panel." Eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1111/ecin.13267>, *Economic Inquiry* 63 (2): 357–386. Accessed February 12, 2025. <https://doi.org/10.1111/ecin.13267>.
- Fischbacher, Urs, Katharine Bendrick, and Stefan Schmid. 2021. *z-Tree 5.1: Tutorial and Reference Manual*. Mimeo. University of Zurich. Accessed January 26, 2025. <https://www.ztree.uzh.ch/static/doc/manual.pdf>.
- Frankel, Jeffrey A., and David H. Romer. 1999a. "Does Trade Cause Growth?" *American Economic Review* 89 (3): 379–399. Accessed February 8, 2025. <https://doi.org/10.1257/aer.89.3.379>.
- . 1999b. *Replication data for: Does Trade Cause Growth?* Accessed February 8, 2025. <https://doi.org/10.3886/E113211V1>.
- Gabelica, Mirko, Ružica Bojčić, and Livia Puljak. 2022. "Many researchers were not compliant with their published data sharing statement: a mixed-methods study." *Journal of Clinical Epidemiology* 150:33–41. Accessed April 4, 2024. <https://doi.org/10.1016/j.jclinepi.2022.05.019>.

- Gadouche, Kamel. 2019. “Le Centre d’accès sécurisé aux données (CASD), un service pour la data science et la recherche scientifique.” *Courrier des statistiques (INSEE)*, no. 3, accessed June 3, 2020. <https://www.insee.fr/fr/information/4254227?sommaire=4254170>.
- Ginsparg, Paul. 2011. “It was twenty years ago today ...” ArXiv: 1108.2700, *arXiv:1108.2700 [astro-ph, physics:cond-mat, physics:gr-qc, physics:hep-ph, physics:hep-th, physics:physics, physics:quant-ph]*, accessed April 6, 2021. <http://arxiv.org/abs/1108.2700>.
- Goebel, Jan, Markus M. Grabka, Stefan Liebig, Martin Kroh, David Richter, Carsten Schröder, and Jürgen Schupp. 2019. “The German Socio-Economic Panel (SOEP).” *Jahrbücher für Nationalökonomie und Statistik* 239 (2): 345–360. Accessed January 31, 2025. <https://doi.org/10.1515/jbnst-2018-0022>.
- Goebel, Jan, Markus M. Grabka, Stefan Liebig, Carsten Schröder, Sabine Zinn, Charlotte Bartels, Jascha Dräger, et al. 2024. *Socio-Economic Panel, data from 1984-2022, (SOEP-Core, v39, International Edition) Sozio-oekonomisches Panel, Daten der Jahre 1984-2022 (SOEP-Core, v39, International Edition)*. Accessed January 31, 2025. <https://doi.org/10.5684/SOEP.CORE.V39I>.
- Goldsmith-Pinkham, Paul, Peter Hull, and Michal Kolesár. 2024. “Contamination Bias in Linear Regressions.” *American Economic Review* 114 (12): 4015–4051. Accessed February 11, 2025. <https://doi.org/10.1257/aer.20221116>.
- Grant, Sean, and Kathryn E. Bouskill. 2019. “Opinion: Why institutional review boards should have a role in the open science movement.” *Proceedings of the National Academy of Sciences* 116 (43): 21336–21338. Accessed February 24, 2020. <https://doi.org/10.1073/pnas.1916420116>.
- Groves, Robert M. 2011a. “Designed Data” and “Organic Data”. Accessed February 9, 2025. <https://web.archive.org/web/20250202073014/https://www.census.gov/newsroom/blogs/director/2011/05/designed-data-and-organic-data.html>.
- . 2011b. “Three Eras of Survey Research.” *Public Opinion Quarterly* 75 (5): 861–871. Accessed February 9, 2025. <https://doi.org/10.1093/poq/nfr057>.
- Haerpfer, Christian, Ronald Inglehart, Alejandro Moreno, Christian Welzel, Kseniya Kizilova, Jaime Diez-Medrano, Marta Lagos, Pippa Norris, Eduard Ponarin, and Bi Puranen. 2024. *World Values Survey Wave 7 (2017-2022) Cross-National Data-Set*. Accessed January 31, 2025. <https://doi.org/10.14281/18241.24>.
- Hamermesh, Daniel S. 2013. “Six Decades of Top Economics Publishing: Who and How?” *Journal of Economic Literature* 51 (1): 162–172. Accessed March 11, 2025. <https://doi.org/10.1257/jel.51.1.162>.
- Herbert, Sylvérie, Hautahi Kingi, Flavio Stanchi, and Lars Vilhuber. 2024. “Reproduce to validate: A comprehensive study on the reproducibility of economics research.” *Canadian Journal of Economics/Revue canadienne d’économique*, caje.12728. Accessed August 5, 2024. <https://doi.org/10.1111/caje.12728>.
- IDEAS//RePEc. 2025. *Within Country and State Economics Rankings: Norway*. Accessed January 17, 2025. <https://web.archive.org/web/20250117034629/https://ideas.repec.org/top/top.norway.html>.
- Institute for Social Research. 2024. *Panel Study of Income Dynamics*. Institute for Social Research, University of Michigan.
- Jacoby, William. 2015. *AMERICAN JOURNAL OF POLITICAL SCIENCE GUIDELINES FOR PREPARING REPLICATION FILES*. Technical report. American Journal of Political Science. Accessed June 4, 2018. <https://ajpsblogging.files.wordpress.com/2015/03/ajps-guide-for-replic-materials-1-0.pdf>.
- Khachiyan, Arman, Anthony Thomas, Huye Zhou, Gordon Hanson, Alex Cloninger, Tajana Rosing, and Amit Khandelwal. 2022a. *Data and Code for: Using Neural Networks to Predict Micro-Spatial Economic Growth*. Accessed January 15, 2025. <https://doi.org/10.3886/E158002V1>.

- Khachiyan, Arman, Anthony Thomas, Huye Zhou, Gordon Hanson, Alex Cloninger, Tajana Rosing, and Amit K. Khandelwal. 2022b. “Using Neural Networks to Predict Microspatial Economic Growth.” *American Economic Review: Insights* 4 (4): 491–506. Accessed January 15, 2025. <https://doi.org/10.1257/aeri.20210422>.
- Koenker, Roger. 1988. *Asymptotic theory and econometric practice (replication data)*. Version Number: 1. Accessed February 8, 2025. <https://doi.org/10.15456/JAE.2022313.1129100068>.
- Koenker, Roger, and Achim Zeileis. 2009. *On reproducible econometric research (replication data)*. Version Number: 1. Accessed March 11, 2025. <https://doi.org/10.15456/JAE.2022319.1306071580>.
- Kopper, Sarah, Anja Sautmann, and James Turitto. 2020. *J-PAL Guide to de-identifying data*. Technical report. J-PAL Global. Accessed August 8, 2020. <https://www.povertyactionlab.org/sites/default/files/research-resources/J-PAL-guide-to-deidentifying-data.pdf>.
- Koren, Miklós, Marie Connolly, Joan Lull, and Lars Vilhuber. 2022. “Data and Code Availability Standard.” Publisher: Zenodo Version Number: 1.0, accessed February 8, 2025. <https://doi.org/10.5281/ZENODO.7436134>.
- Korinek, Anton. 2023. “Generative AI for Economic Research: Use Cases and Implications for Economists.” *Journal of Economic Literature* 61 (4): 1281–1317. Accessed December 14, 2024. <https://doi.org/10.1257/jel.20231736>.
- Krichel, Thomas. 1997. “WoPEc: Electronic Working Papers in Economics Services.” *Ariadne*, no. 8, accessed July 19, 2018. <http://www.ariadne.ac.uk/issue8/wopec>.
- Krichel, Thomas, and Christian Zimmermann. 2009. “The Economics of Open Bibliographic Data Provision.” *Economic Analysis and Policy* 39 (1): 143–152. Accessed July 16, 2018. [https://doi.org/10.1016/S0313-5926\(09\)50049-5](https://doi.org/10.1016/S0313-5926(09)50049-5).
- Max-Planck-Gesellschaft. 2023. *Berlin Declaration*. Accessed December 31, 2024. <https://openaccess.mpg.de/Berlin-Declaration>.
- Merton, Robert K. 1942. “A note on science and democracy.” *Journal of Legal and Political Sociology* 1:115–126.
- Meyer, Michelle N. 2018. “Practical Tips for Ethical Data Sharing.” *Advances in Methods and Practices in Psychological Science* 1 (1): 131–144. Accessed March 12, 2019. <https://doi.org/10.1177/2515245917747656>.
- Miguel, Edward. 2021. “Evidence on Research Transparency in Economics.” Publisher: American Economic Association, *Journal of Economic Perspectives* 35 (3): 193–214. Accessed July 13, 2025. <https://doi.org/10.1257/jep.35.3.193>.
- Mogstad, Magne, Joseph Romano, Azeem Shaikh, and Daniel Wilhelm. 2022. “Statistical Uncertainty in the Ranking of Journals and Universities.” *AEA Papers and Proceedings* 112:630–634. Accessed December 31, 2024. <https://doi.org/10.1257/pandp.20221064>.
- Mukherjee, Soumya, Aratrika Mustafi, Aleksandra Slavković, and Lars Vilhuber. 2023. *Assessing Utility of Differential Privacy for RCTs*. Technical report arXiv:2309.14581. ArXiv:2309.14581 [cs, econ, stat]. arXiv. Accessed April 8, 2024. <https://doi.org/10.48550/arXiv.2309.14581>.
- Murray, Fiona, Philippe Aghion, Mathias Dewatripont, Julian Kolev, and Scott Stern. 2016. “Of Mice and Academics: Examining the Effect of Openness on Innovation.” *American Economic Journal: Economic Policy* 8 (1): 212–252. Accessed May 16, 2025. <https://doi.org/10.1257/pol.20140062>.
- Nagaraj, Abhishek, Esther Shears, and Mathijs de Vaan. 2020. “Improving data access democratizes and diversifies science.” *Proceedings of the National Academy of Sciences* 117 (38): 23490–23498. Accessed March 4, 2022. <https://doi.org/10.1073/pnas.2001682117>.

- Nagaraj, Abhishek, and Matteo Tranchero. 2023. *How Does Data Access Shape Science? The Impact of Federal Statistical Research Data Centers on Economics Research*. Working Paper. Accessed February 9, 2025. <https://doi.org/10.3386/w31372>.
- National Bureau of Economic Research. 2025. *Affiliated Scholars*. Accessed January 17, 2025. <https://web.archive.org/web/20250110073711/https://www.nber.org/affiliated-scholars?page=1&perPage=50>.
- National Science Foundation. 2024. *Doctorate Recipients from U.S. Universities: 2023*. Report NSF 25-300. Accessed January 17, 2025. <https://ncses.nsf.gov/pubs/nsf25300>.
- Nosek, B. A., G. Alter, G. C. Banks, D. Borsboom, S. D. Bowman, S. J. Breckler, S. Buck, et al. 2015. "Promoting an open research culture." *Science* 348 (6242): 1422–1425. Accessed October 9, 2018. <https://doi.org/10.1126/science.aab2374>.
- Open Knowledge Foundation. 2024. *Defining Open in Open Data, Open Content and Open Knowledge: Open Definition 2.1*. Accessed December 30, 2024. <https://opendefinition.org/od/2.1/en/>.
- Panettieri, Joe. 2025. *Generative AI Lawsuits Timeline: Legal Cases vs. OpenAI, Microsoft, Anthropic, Nvidia, Perplexity, Intel and More*. Accessed May 16, 2025. <https://sustainabletechpartner.com/topics/ai/generative-ai-lawsuit-timeline/>.
- Piowar, Heather A., and Todd J. Vision. 2013. "Data reuse and the open data citation advantage." *PeerJ* 1:e175. Accessed January 28, 2019. <https://doi.org/10.7717/peerj.175>.
- Roth, Jonathan. 2022. "Pretest with Caution: Event-Study Estimates after Testing for Parallel Trends." *American Economic Review: Insights* 4 (3): 305–322. Accessed November 16, 2023. <https://doi.org/10.1257/aeri.20210236>.
- Rudik, Ivan. 2020a. *Data and Code for: Optimal climate policy when damages are unknown*. Version Number: 1. Accessed January 15, 2025. <https://doi.org/10.3886/E111185V1>.
- . 2020b. "Optimal Climate Policy When Damages are Unknown." *American Economic Journal: Economic Policy* 12 (2): 340–373. Accessed January 15, 2025. <https://doi.org/10.1257/pol.20160541>.
- S&P Dow Jones Indices LLC. 2025. *S&P 500 [SP500]*. Federal Reserve Bank of St. Louis. Accessed January 31, 2025. <https://fred.stlouisfed.org/series/SP500>.
- Sheth, Ketki, Amy Shipow, Aimable Nsabimana, Maya Ranganath, and Amos Njuguna. n.d. *Understanding Barriers to and Expanding Opportunities for Publishing Quality Research from African Scholars (tentative title)*. Collaboration for Inclusive Development Research. Center for Effective Global Action, UC Berkeley. Accessed July 15, 2025. <https://cega.berkeley.edu/article/cega-niera-launch-cidr-survey/>.
- Silcock, Emily, Abhishek Arora, Luca D'Amico-Wong, and Melissa Dell. 2024. *Newsire [Database]*. Accessed January 15, 2025. <https://doi.org/10.57967/HF/2423>.
- Slavković, Aleksandra, and Lars Vilhuber. 2018. "Remembering Stephen Fienberg." Number: 1, *Journal of Privacy and Confidentiality* 8 (1). Accessed January 1, 2025. <https://doi.org/10.29012/jpc.685>.
- Social Science Research Network*. 2025. Page Version ID: 1273115470. Accessed March 2, 2025. https://en.wikipedia.org/w/index.php?title=Social_Science_Research_Network&oldid=1273115470.
- Sociological Science. 2018. *Manuscript Preparation*. Accessed December 22, 2018. <https://www.sociologicalscience.com/for-authors/manuscript-preparation/>.
- Statistics Canada. 2012. *Statistics Canada Open Licence*. Last Modified: 2021-10-29. Accessed February 20, 2022. <https://www.statcan.gc.ca/en/reference/licence>.
- Tirole, Jean, Torsten Persson, Peter Neary, and Richard Blundell. 2003. "Editorial." *Journal of the European Economic Association* 1 (1): iii–iv. Accessed January 1, 2025. <https://doi.org/10.1162/154247603322256747>.

- U.S. Census Bureau. 2024. *uscensusbureau/fsrdc-external-census-projects at 2625c2169f2d1b22cac262b90f7af87b7e969d6b*. Github. Accessed January 17, 2025. <https://github.com/uscensusbureau/fsrdc-external-census-projects/tree/2625c2169f2d1b22cac262b90f7af87b7e969d6b>.
- UK Government. 2014. *Open Government Licence for public sector information V3*. Accessed February 20, 2022. <https://www.nationalarchives.gov.uk/doc/open-government-licence/version/3/>.
- UNESCO. 2022. *Understanding open science*. Online Open Access programme and meeting document. UNESCO. Accessed December 31, 2024. <https://doi.org/10.54677/UTCD9302>.
- Vicente-Saez, Ruben, and Clara Martinez-Fuentes. 2018. “Open Science now: A systematic literature review for an integrated definition.” *Journal of Business Research* 88:428–436. Accessed December 31, 2024. <https://doi.org/10.1016/j.jbusres.2017.12.043>.
- Vilhuber, Lars. 2018. “Relaunching the Journal of Privacy and Confidentiality.” Number: 1, *Journal of Privacy and Confidentiality* 8 (1). Accessed January 1, 2025. <https://doi.org/10.29012/jpc.706>.
- . 2020a. *Migrating historical AEA supplements*. Release v20200515. Cornell University. <https://github.com/AEADDataEditor/aea-supplement-migration/releases/tag/v20200515>.
- . 2020b. “Reproducibility and Replicability in Economics.” *Harvard Data Science Review* 2 (4). <https://doi.org/10.1162/99608f92.4f6b9e67>.
- . 2023. “Reproducibility and transparency versus privacy and confidentiality: Reflections from a data editor.” Published online, *Journal of Econometrics* 235 (2): 2285–2294. Accessed June 6, 2023. <https://doi.org/10.1016/j.jeconom.2023.05.001>.
- . 2025. “Report of the AEA Data Editor.” *AEA Papers and Proceedings*.
- Vilhuber, Lars, Marie Connolly, Miklós Koren, Joan Llull, and Peter Morrow. 2022. *A template README for social science replication packages*. Technical report v1.1.0. Publisher: Zenodo. Zenodo. Accessed May 17, 2023. <https://doi.org/10.5281/zenodo.7293838>.
- Vilhuber, Lars, Ian Schmutte, Aleksandr Michuda, and Marie Connolly. 2023. “Reinforcing Reproducibility and Replicability: An Introduction.” *Harvard Data Science Review* 5 (3). Accessed August 2, 2023. <https://doi.org/10.1162/99608f92.9ba2bd43>.
- Vilhuber, Lars, James Turitto, and Keesler Welch. 2020. “Report by the AEA Data Editor.” *AEA Papers and Proceedings* 110:764–75. <https://doi.org/10.1257/pandp.110.764>.
- Vlaeminck, Sven. 2021. “Dawning of a new age? Economics journals’ data policies on the test bench.” *LIBER Quarterly: The Journal of the Association of European Research Libraries* 31 (1): 1–29. Accessed August 31, 2021. <https://doi.org/10.53377/lq.10940>.
- Watson, Clare. 2022. “Many researchers say they’ll share data — but don’t.” Bandiera_abtest: a Cg_type: News Publisher: Nature Publishing Group Subject.term: Funding, Ethics, *Nature* 606 (7916): 853–853. Accessed April 4, 2024. <https://doi.org/10.1038/d41586-022-01692-1>.
- Weeden, Kim A. 2023. “Crisis? What Crisis? Sociology’s Slow Progress Toward Scientific Transparency.” *Harvard Data Science Review* 5 (4). Accessed November 27, 2023. <https://doi.org/10.1162/99608f92.151c41e3>.
- Welch, Finis. 1973. *Education, Information, and Efficiency*. Working Paper w0001. National Bureau of Economic Research. Accessed December 30, 2024. <https://www.nber.org/papers/w0001>.
- Wikipedia. 2025. *Copyright status of works by the federal government of the United States*. Page Version ID: 1270780906. Accessed February 9, 2025. https://en.wikipedia.org/w/index.php?title=Copyright_status_of_works_by_the_federal_government_of_the_United_States&oldid=1270780906.

- Williams, Heidi L. 2013. “Intellectual Property Rights and Innovation: Evidence from the Human Genome.” *Journal of Political Economy* 121 (1): 1–27. Accessed May 16, 2025. <https://doi.org/10.1086/669706>.
- World Bank. 2021. *World Development Report 2021: Data for Better Lives*. Technical report. World Bank. Accessed May 16, 2025. <https://hdl.handle.net/10986/35218>.
- Xie, Jin, and Joseph Gerakos. 2020. “The Anticompetitive Effects of Common Ownership: The Case of Paragraph IV Generic Entry.” *AEA Papers and Proceedings* 110:569–572. Accessed May 16, 2025. <https://doi.org/10.1257/pandp.20201029>.
- Zillow. 2021. *Zillow’s Assessor and Real Estate Database*. Accessed March 4, 2022. <https://perma.cc/DA2S-XMK2>.

Appendices

Table 4: Access categories and whether data can be shared privately

Access restrictions	No	Yes	Total	Percent
No		n/a	50	
Very Easy to Obtain	11	7	18	(38.89%)
Moderately Easy to Obtain	3	5	8	(62.50%)
Moderately Difficult to Obtain	16	9	25	(36.00%)
Very Difficult to Obtain	9	5	14	(35.71%)
Any restriction	39	26	65	(40.00%)

* Percentages are calculated as the number of ‘Yes’ divided by the total number of responses. An article can have multiple categories of data; the sum of responses is therefore higher than the number of articles.

Inferring Stata and R usage

Stata usage is inferred from downloads of Stata packages from the SSC web server, which is the sole official location to obtain these packages. Private mirrors may exist, and not all Stata packages are installed from SSC - both the Stata Journal and Github are likely to be significant sources. Data are obtained from log files from the SSC web server, provided by Kit Baum. R usage is inferred from downloads of R binaries (for Windows and MacOS). While this is likely to be less frequent than Stata packages, relative patterns are of interest. There is no easy way to obtain the full list of downloads for all packages, other than cycling through several thousand such packages. The data stem from one of dozens of Comprehensive R Archive Network (CRAN) mirrors, though this one, managed by Posit PBC, is the first one listed in the list of mirrors that are offered to users. Data in both cases is for February 2025.

As a first step, these downloads were mapped to countries, and then aggregated by regions (Table 5 for Stata, Table 6 for R). Mapped onto a world map, this provides a pretty picture, though not very informative (Figure 6). The data suggest that downloads from China are not fully captured by the Posit-managed mirror, possibly because of peculiarities of the Chinese internet infrastructure. While this is possible for other countries as well, it is not fully detectable. I have excluded China from all calculations, but the remaining computations should be regarded as indicative as best.

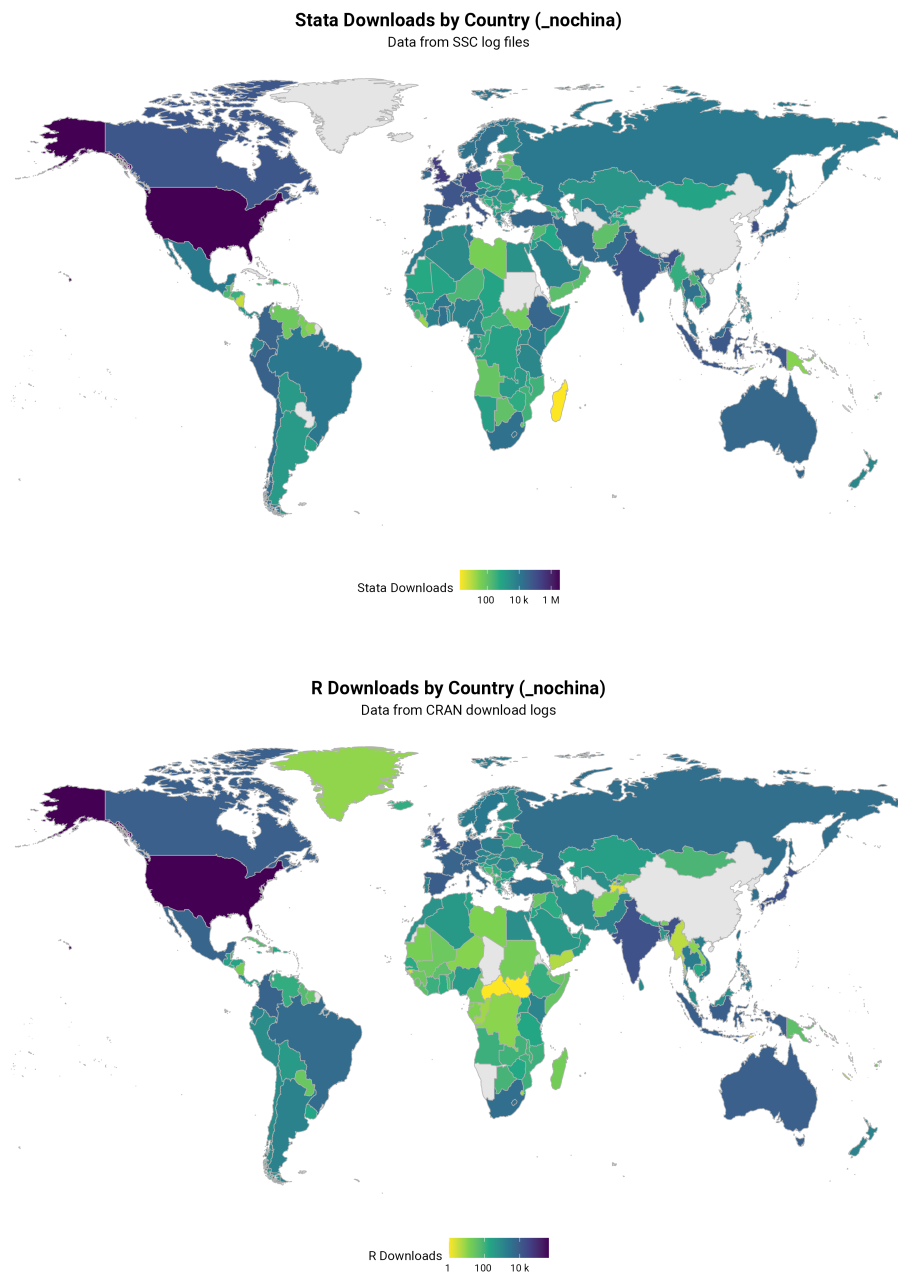


Figure 6: [Downloads of Stata packages and R software, by region, without China. Source: published data by Kit Baum and Posit PBC, own calculations.](#)

	custom_region	countries	Regional Downloads	Fraction
1	China	2	5438811	50.79
2	North America	3	3509194	32.77
3	Europe	40	912383	8.52
4	Rest of Asia	45	521557	4.87
5	Africa	51	168663	1.57
6	Latin America & Caribbean	28	127078	1.19
7	Australia	1	26703	0.25
8	Rest of Oceania	5	4763	0.04

Table 5: Regional Statistics for Stata Downloads

	custom_region	countries	Regional Downloads	Fraction
1	North America	3	438562	66.62
2	Europe	50	83367	12.66
3	Rest of Asia	48	72724	11.05
4	Latin America & Caribbean	42	25316	3.85
5	China	1	18641	2.83
6	Africa	51	10124	1.54
7	Australia	1	8166	1.24
8	Rest of Oceania	11	1406	0.21
9	Other	1	9	0.00

Table 6: Regional Statistics for R Downloads

I then aggregate countries and regions into "Global North" and "Global South" (Table 7 for Stata, Table 8 for R). Differencing the two tables gives the relative importance of Stata versus R, namely that there are

	North/South	Countries	Downloads	Fraction
1	Global North	44	4448280	84.40
2	Global South	129	822061	15.60

Table 7: Regional Statistics for Stata Downloads (excl. China)

relatively more downloads of Stata packages in the Global North than there are of R binaries, suggesting (very tentatively) that Stata usage may be higher in the Global North (Table 9) Reverting back to smaller regions, however, paints the more nuanced picture referenced in the main text (Table 10). Tables 9 and 10 are summarized in Figure 4 in the main text.

	North/South	Countries	Downloads	Fraction
1	Global North	54	530095	82.87
2	Global South	152	109570	17.13
3	Other	1	9	0.00

Table 8: Regional Statistics for R Downloads (excl. China)

	North/South	Stata	R	Diff
1	Global North	84.40	82.87	1.53
2	Global South	15.60	17.13	-1.53

Table 9: North/South Share of Downloads: Stata vs R

	custom_region	Stata	R	Diff
1	Europe	17.31	13.03	4.28
2	North America	66.58	68.56	-1.98
3	Africa	3.20	1.58	1.62
4	Latin America & Caribbean	2.41	3.96	-1.55
5	Rest of Asia	9.90	11.37	-1.47
6	Australia	0.51	1.28	-0.77
7	Rest of Oceania	0.09	0.22	-0.13

Table 10: Regional Share of Downloads: Stata vs R