
AN ANALYSIS OF THE EFFECTS OF SHARING RESEARCH DATA, CODE, AND PREPRINTS ON CITATIONS

Giovanni Colavizza[†]
University of Bologna

Lauren Cadwallader
PLOS

Marcel LaFlamme
PLOS

Grégory Dozot
HEIG-VD

Stéphane Lecorney
HEIG-VD

Daniel Rappo
HEIG-VD

Iain Hrynaszkiewicz^{††}
PLOS

April 26, 2024

ABSTRACT

Calls to make scientific research more open have gained traction with a range of societal stakeholders. Open Science practices include but are not limited to the early sharing of results via preprints and openly sharing outputs such as data and code to make research more reproducible and extensible. Existing evidence shows that adopting Open Science practices has effects in several domains. In this study, we investigate whether adopting one or more Open Science practices leads to significantly higher citations for an associated publication, which is one form of academic impact. We use a novel dataset known as Open Science Indicators, produced by PLOS and DataSeer, which includes all PLOS publications from 2018 to 2023 as well as a comparison group sampled from the PMC Open Access Subset. In total, we analyze circa 122'000 publications. We calculate publication and author-level citation indicators and use a broad set of control variables to isolate the effect of Open Science Indicators on received citations. We show that Open Science practices are adopted to different degrees across scientific disciplines. We find that the early release of a publication as a preprint correlates with a significant positive citation advantage of about 20.2% (± 7) on average. We also find that sharing data in an online repository correlates with a smaller yet still positive citation advantage of 4.3% (± 8) on average. However, we do not find a significant citation advantage for sharing code. Further research is needed on additional or alternative measures of impact beyond citations. Our results are likely to be of interest to researchers, as well as publishers, research funders, and policymakers.

1 Introduction

Arising from a diverse set of cultural and technological projects at the turn of the twenty-first century [1, 2, 3], contemporary calls to make scientific research more open point toward a no less diverse range of outcomes. One influential definition characterizes Open Science as “transparent and accessible knowledge that is shared and developed through collaborative networks” [4], encompassing knowledge objects or outputs as well as processes [5]. Another developed by UNESCO defines Open Science as “an inclusive construct that combines various movements and practices aiming to make multilingual scientific knowledge openly available, accessible and reusable for everyone, to increase scientific collaborations and sharing of information for the benefits of science and society, and to open the processes of scientific knowledge creation, evaluation and communication to societal actors beyond the traditional scientific community” [6].

While acknowledging this diversity of ambitions, in what follows we focus on practices resulting in what UNESCO terms “open scientific knowledge” [6]: that is, the making of scientific publications and the materials that underpin them

[†]giovanni.colavizza@unibo.it

^{††}ihrynaszkiewicz@plos.org

available to all, free of charge. These Open Science practices include but are not limited to Open Access publication; the early sharing of results, for example via the use of preprints; openly sharing outputs such as data, code, and protocols to make research more reproducible and extensible; and fostering rigor and transparency in study design, for example via study registration. While the uptake of these practices by researchers varies by field, career stage, and region, their prevalence is growing overall [7, 8]. Drivers of this growth include publisher and funder policies, training and infrastructure support, and cultural change [9, 10]. The proliferation of policies for Open Science has led to a greater need to monitor the effects of these policies on Open Science [11], although comprehensive solutions for measuring Open Science are lacking. Still, researchers, technology providers, research funders, institutions and publishers have begun to monitor the prevalence of Open Science practices (<https://open-science-monitoring.org/monitors/>). These efforts provide new evidence and data sources from which to understand if and how Open Science practices are being adopted, and to explore the extent to which these practices confer effects, impacts or benefits as a result of their adoption.

In this article, we contribute to an emerging strand of research assessing the impact of Open Science practices. We focus on a set of measurable Open Science practices that include data sharing, code sharing, and preprint posting. More specifically, we ask whether adopting any combination of these practices leads to a significantly higher citation impact for an associated publication when compared to similar publications for which authors have not adopted Open Science practices. We answer this question by leveraging a novel dataset known as Open Science Indicators, which is produced by the nonprofit Open Access publisher PLOS in partnership with DataSeer (<https://dataseer.ai>) [12], and by adapting a previously released workflow to mine citation data from the PMC Open Access Subset [13]. An important aspect of our contribution is the assessment of Open Science practices in combination, rather than individually as is usually the case in previous work.

2 State of the Art

There is evidence that adopting Open Science practices has effects or impacts in several domains: academic, societal, and economic [14]. In terms of academic or research impacts, Open Science practices are associated with increased visibility and reuse, as measured for example by the diversity of citations and media attention received by Open Access articles [15, 16]. Open Science has been instrumental in accelerating progress on certain scientific problems [17], in making results more transparent [18], and in addressing what has been termed the replication crisis in certain fields of research [19]. Societal benefits identified in a systematic scoping review include enabling broader participation in research, by supporting citizen science and educational initiatives. However, evidence of societal benefits to policy, health, or trust in research is to date more limited [20]. Economic benefits, identified from economic modeling studies and case studies, include cost and labor savings from Open Access and open (or FAIR) data, as well as increased innovation [21, 22]. However, there is a lack of causal evidence for and prospective studies of these benefits. Open Science practices have also been linked to negative impacts including imposing additional costs [23], reinforcing existing inequalities [24], and homogenizing diverse research traditions [25].

2.1 Data and code sharing

Researchers who adopt Open Science practices may see increased use and impact of their work, which can support career progression. Several studies examine the importance of data (and code) sharing for scientific advancement but diverge to some extent in their findings. Evidence shows that the novel combination of datasets leads to higher impact and visibility [26]. Several studies in specific research disciplines have found correlations between sharing research and increased citations of articles that share data [27, 28, 29]. Implementation of journal policies requiring data sharing has also been correlated with increased citations [30]. Researchers can share data in several different ways but sharing data privately, upon request, and via supporting information files with publications are the most common approaches – despite being considered suboptimal [31, 32]. Sharing research data in a public repository is considered best practice for data sharing but may require additional effort compared to other approaches [33]. However, in previous work, we found that, relative to sharing data upon request or as supporting information files, data sharing in repositories was correlated with a 25.36% citation advantage on average [13].

While we can hypothesize that, similar to data sharing, code sharing would promote the reuse of published research that shares code, there is less evidence about whether code sharing is correlated with any effect on citations. Studies of a single journal [34] or a small number of journals in a single field [35] have found mixed effects. A larger-scale study showed a correlation between links to methods including (but not limited to) code and increased citation, especially when links were still active [36]. Another found that monthly citations of articles increased after their associated code repositories were made public [37].

2.2 Preprints

There is evidence for advantages in terms of visibility for peer-reviewed publications that were previously posted as preprints, as measured by increased citations and altmetrics [38, 39, 40, 41]. This effect was examined in detail during the COVID-19 pandemic [42], when media coverage of health-related preprints also saw a significant uptick [43]. Other forms of impact associated with preprint posting include receiving additional feedback, which research has shown to be constructive if variable in frequency [44, 45]. Studies examining the adoption of preprints by career stage have suggested that they have particular advantages for early-career researchers in terms of career development [46, 47]. At the same time, concerns over how preprints may introduce unvetted findings into the scientific record have pointed to the need for nuanced approaches to evidence synthesis [48, 49] and the communication of retractions [50].

3 Methods and data

To make this study entirely reproducible, we focus only on Open Access publications and release all of the accompanying code. We strictly follow and expand upon a published methodology [13]. This methodology entails selecting a set of publications of interest, calculating publication and author-level citation counts using a larger Open Access collection, and modeling the effect of interest as independent variables. We use PLOS' Open Science Indicators version 5 as a starting point [12, 51]. The OSI publication count totals $N = 124'274$. We also use the PMC Open Access Subset, with all publications up to October 2023 included [52]. The PMC Open Access Subset is used to calculate citation counts for publications and authors. Citation counts calculated using the PMC Open Access Subset have been shown to track global citation counts, and thus to be appropriate when the relative rather than absolute counts are of interest [13]. Publications missing a known identifier (DOI, PubMed reference number, PMCID, or a publisher-specific ID), a publication date, and at least one reference are discarded. These often are editorials, letters, or similar article types. The final PMC Open Access Subset publication count totals $M = 5'020'948$. After an initial analysis, a limited amount of OSI publications are also discarded for being absent in the PMC Open Access Subset or identified as editorials or reviews (i.e., not research articles). Of the $124'376$ publications in OSI, $121'999$ (98.1%) are processed, matched, and used for the modeling analysis that follows.

We use a linear model for quantifying the relative effect of Open Science Indicators on citation counts, as follows:

Dependent variable. Citation counts for each publication are calculated using the full PMC Open Access Subset dataset (M publications above). Citations are based on identifiers, hence only references that include a valid ID are considered. Under these limits, we calculate total citation counts and use this as our main dependent variable. We also calculate citations given within a certain time window from publication (1, 2, and 3 years, also considering the month of publication). This is done in order to conduct a robustness check using citation counts over the same citation accrual time as the dependent variable (e.g., the three-year window for a publication published in June 2015 runs to June 2018 excluded).

Independent variables. We use a set of control variables for modeling. Firstly, publication-level variables are commonly considered in similar studies [53, 54, 55]. We include the year of publication, to account for citation inflation over time; the month of publication (missing values are set to a default value of 6, that is June), to account for the advantage of publications published early in the year that have more time to accrue citations; the number of authors and the total number of references (including those without a known identifier), both usually correlated with citation impact. We also use the Australian and New Zealand Standard Research Classification (ANZSRC) Fields of Research classification system at the publication level, to account for disciplinary variation in the adoption of Open Science practices. We use the broadest level provided, that of the Division, to avoid data sparsity. In the dataset, 22 divisions are found. We group the least frequent five categories into a single category, since they all belong to the Arts and Humanities. We therefore end up with 18 distinct categories that are encoded as dummy variables to account for the fact that a publication can belong to multiple categories. See Table 1 for a list of the categories used from the division-level ANZSRC Fields of Research.

The reputation of authors before publication has also been linked to the citation success of a paper [56]. To control for this, we have to identify individual authors, a challenging task in itself [57, 58, 59, 60, 61]. We focus on a publication-level aggregated indicator of author popularity: the mean H-index of a publication's authors at the time of publication, calculated from the PMC Open Access Subset. In so doing, we minimize the impact of errors arising from disambiguating author names [62, 63], which would have been higher if we had used measures based on individual observations such as the maximum H-index. We therefore use a simple disambiguation technique when compared to the current state of the art, and consider two author mentions to refer to the same individual if both full name and surname are found to be identical within all of the PMC Open Access Subset. We acknowledge the limitations of this method in

Table 1: ANZSRC Fields of Research Divisions to model categories. Note that the total publication count is higher than the number of publications in OSI, since a publication can belong to more than one division.

ANZSRC FoR Division	Publication counts in OSI	Category (regression model)
32 Biomedical and Clinical Sciences	59'377	division_1
31 Biological Sciences	35'081	division_2
42 Health Sciences	29'778	division_3
30 Agricultural, Veterinary and Food Sciences	8549	division_4
46 Information and Computing Sciences	7704	division_5
52 Psychology	6910	division_6
44 Human Society	5645	division_7
40 Engineering	5294	division_8
41 Environmental Sciences	5208	division_9
34 Chemical Sciences	3468	division_10
35 Commerce, Management, Tourism and Services	2811	division_11
37 Earth Sciences	2716	division_12
38 Economics	1960	division_13
51 Physical Sciences	1313	division_14
39 Education	1153	division_15
47 Language, Communication and Culture	988	division_16
49 Mathematical Sciences	901	division_17
43 History, Heritage and Archaeology	838	division_18
48 Law and Legal Studies	792	division_18
33 Built Environment and Design	662	division_18
36 Creative Arts and Writing	459	division_18
50 Philosophy and Religious Studies	399	division_18

Table 2: Descriptive statistics for the dependent variable and a set of publication and author level controls.

	n_cit_tot	n_cit_2	n_authors	n_references_tot	p_year	p_month	h_index_mean
Min.	0	0	1	0	2018	1	0
1st Qu.	0	0	4	34	2019	3	2
Median	2	1.0	6	46	2020	6	3.9
Mean	5.1	2.3	7.1	51.1	2020	5.4	5
3rd Qu.	6	3	9	63	2022	7	6.6
Max.	3683	788	2621	986	2023	12	57
NA's							425

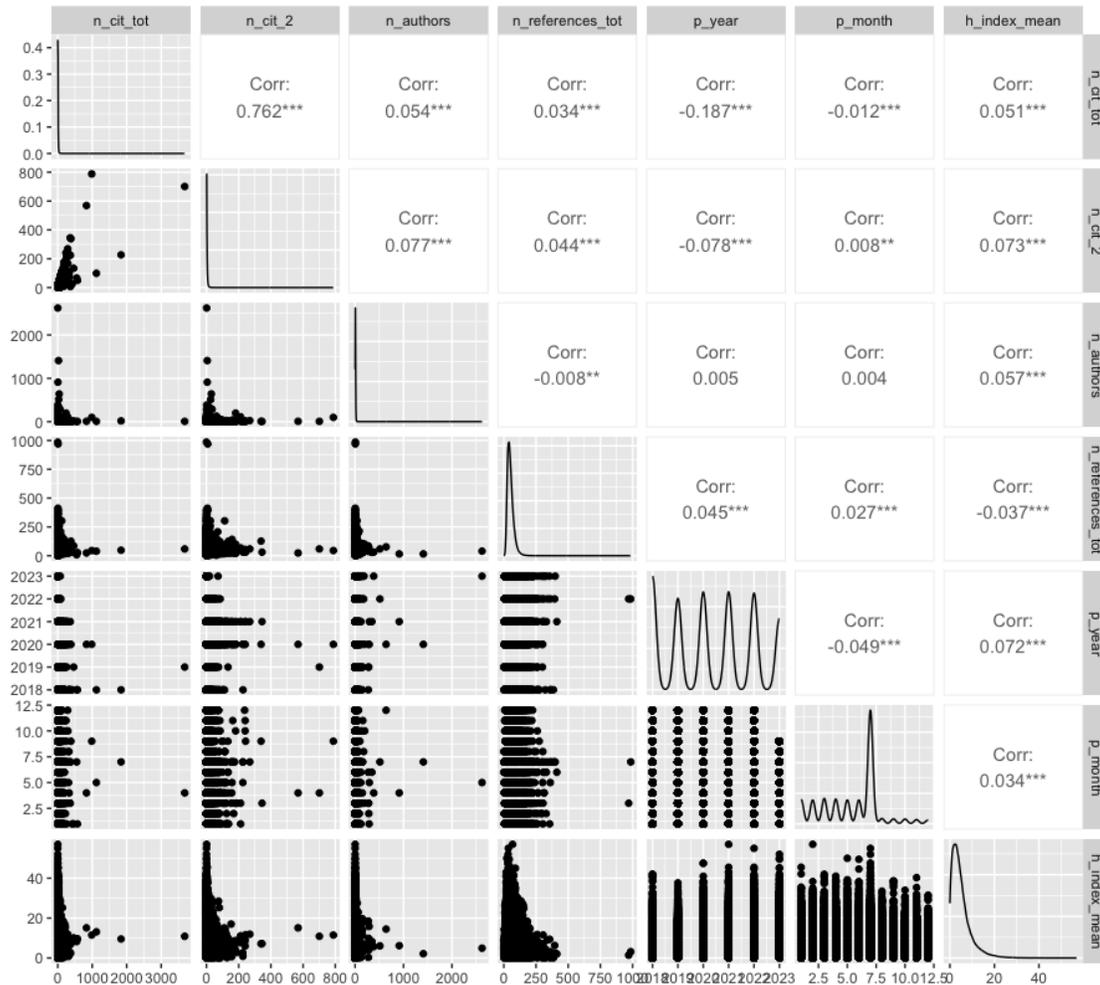
possibly merging different authors with the same name and surname. We identify 8'481'129 seemingly distinct authors in this way.

We finally consider the following journal-level variables: if a publication is published by PLOS (any journal), and if a publication is published in PLOS ONE specifically. Given the preponderance of PLOS publications (101'366, or nearly 85% of publications overall), and specifically PLOS ONE publications (83'843, or nearly 70% of publications overall) in the dataset, we do not use any other journal-level variable.

A set of descriptive statistics for the numerical variables in use is reported in Tables 2 and 3, while their correlations are illustrated in Figure 1. The models we test, besides OLS linear regression and robust linear regression, include ANOVA, Tobit, and GLM with negative binomial, zero-inflated negative binomial, lognormal, and Pareto 2 family distributions. These largely support the findings using linear regression and robust linear regression, which are easier to interpret. Therefore, results from other models are omitted here and can be reproduced using the accompanying codebase. Robust linear regression results differ little from simple linear regression, as is expected given the log transformations we systematically apply on skewed numerical variables, but they are provided for comparison.

Table 3: Descriptive statistics for the OSI controls. C: Code; D: Data; Repo: Repository; P: Preprint.

D_Shared	D_Location	Repo_Data	C_Generated	C_Shared	C_Location	P_Match
N: 39'124	N/A: 39'657	F: 94'715	N: 76'452	N: 107'390	N/A: 107'443	F: 97'235
Y: 82'875	Online: 39'700	T: 27'284	Y: 45'547	Y: 14'609	Online: 11'221	T: 24'764
	Suppl. Info: 42'642				Suppl. Info: 3'335	

Figure 1: Correlation plot among most variables. We see that no two variables are too highly correlated, except as expected for two alternatives for dependent variables (n_cit_2 and n_cit_tot).

4 Results

We start by providing a brief descriptive overview of the Open Science Indicators in the target corpus and then proceed to the modeling section.

4.1 Overview of the Open Science Indicators dataset

As mentioned previously, the OSI dataset we use for analysis comprises 121'999 articles. The majority are articles published in PLOS journals, with the largest proportion originating from PLOS ONE. The remaining articles have been taken from 1232 different journals published by a range of publishers. Rates of adoption for each Open Science practice can be calculated from the dataset to give an overall impression of the degree to which Open Science is practiced.

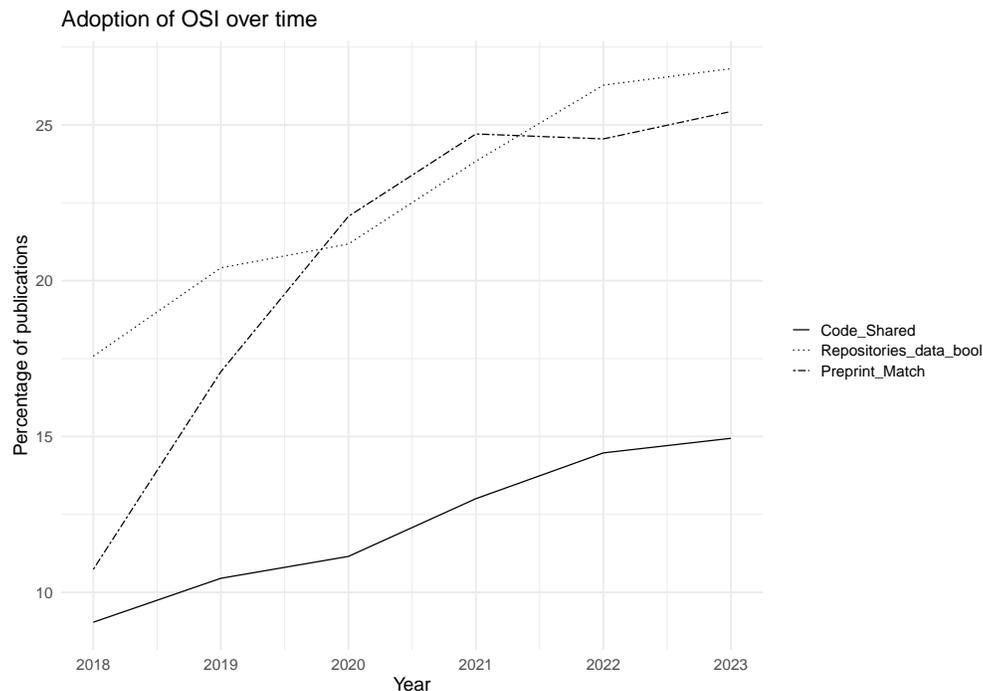
Table 4: Descriptive statistics for the Open Science practices as measured in the OSI dataset.

	Number	% of publications
Publications	121'999	100%
Sharing data (anywhere)	82'875	68%
Sharing data (in a repository)	27'284	22%
Sharing data (online)	33'786	28%
Sharing code	14'609	12%
Has a preprint	24'764	20%

Table 4 outlines the overall rates of adoption for the main Open Science practices in the dataset and shows that data (in a repository) and code sharing have a relatively low adoption rate across the dataset.

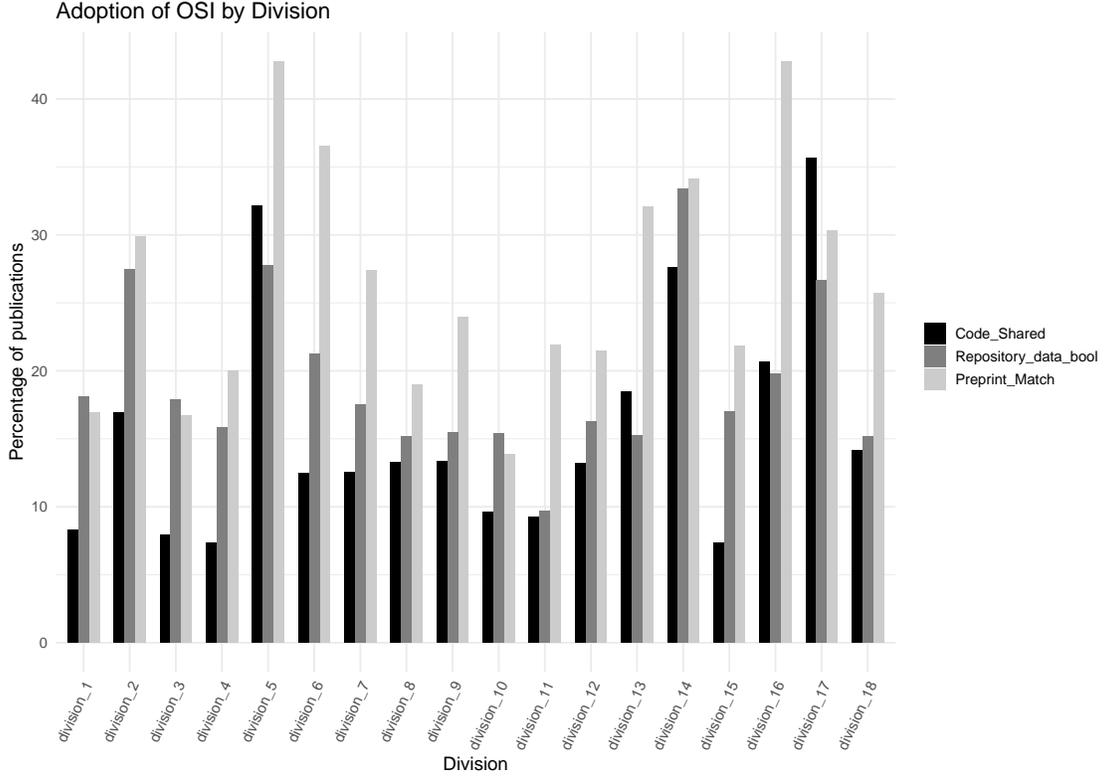
In OSI, the average rates of adoption for Open Science practices observed in the dataset have been increasing over time with changes between 5% and 15% from 2018 to 2023. Data sharing in any form has seen a 5% increase from 2018 to 2023, with data sharing in repositories and online increasing by 9% and 10% respectively. Code sharing (out of all publications) has increased by 6% over the same time period and preprint posting by 15%. Whilst data and code sharing show positive trends over time, the trend for preprint posting shows a large increase between 2018 and 2019 and again from 2019 to 2020, followed by a plateauing since 2021. These trends are also seen when the PLOS cohort and the PMC Open Access Subset cohort are considered separately, although the PMC Open Access Subset cohort shows an increase in preprints in 2023 compared to 2022 which is not present in the PLOS data. We show the general adoption trends in Figure 2.

Figure 2: Adoption of OSI over time. Each OSI remains adopted by a fraction of publications, but adoption grows over time.



The prevalence of different Open Science practices varies by field of research (following the divisions presented in 1). For example, both Division 3 (Health Sciences) and Division 8 (Engineering) have the lowest data sharing rate at 60%, whilst Division 16 (Language, Communication and Culture) has the highest data sharing rate at 82%. Similar degrees in variation are seen for the other indicators with data sharing in a repository ranging from 14% to 43%, code sharing from 7% to 36%, and preprint posting from 10% to 33%. Such wide variation in OSI adoption across divisions suggest that research fields face different challenges in adopting Open Science practices, and some practices may not be equally useful or relevant across fields. In Figure 3, we show trends for the main OSIs across all Divisions, as described in Table 1.

Figure 3: Adoption of OSI by Division, as described in Table 1. Each OSI remains adopted by a fraction of publications, but there is a wide variation across Divisions.



Please refer to the PLOS' Open Science Indicators version 5 documentation for further details [12, 51].

4.2 Modeling

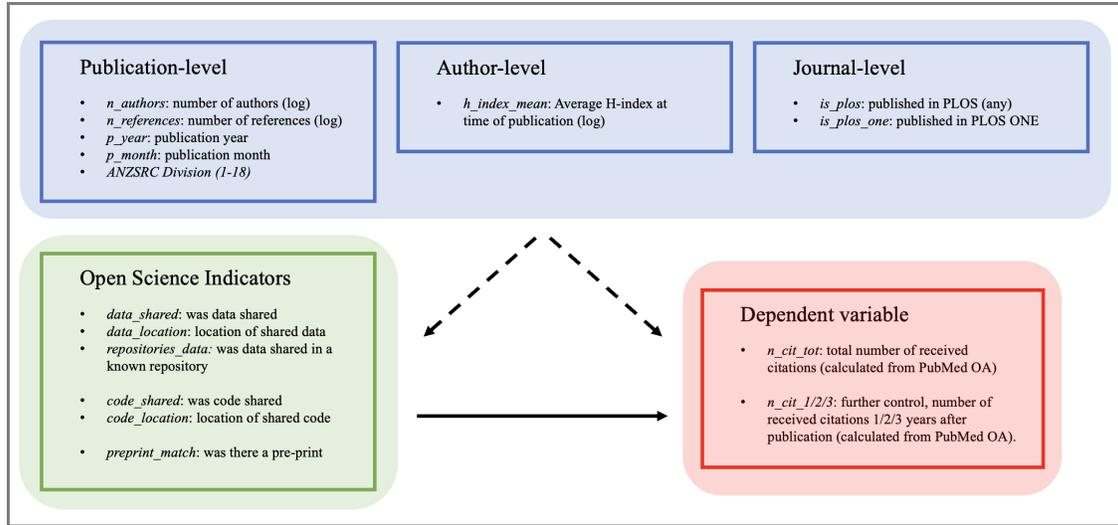
The base model results we discuss are provided in Table 5. It contains the basic author, publication, and journal-level variables we discussed above. It does not contain publication-level division classification. The most complete model we discuss is instead provided in Table 6. Here, we use all the previous variables from the base model and add the publication-level division classification as dummy variables (division 1 to 18, see Table 1. Several more models were tested, primarily as robustness checks, and are discussed in the Appendix.

The base model is described in Equation 1, and the full model is described in Equation 2. Variable transformations are shown, numerical variables are given in *Italics*, and categorical variables are in regular text. Variables are grouped along lines. An illustration of the assumed causal dependency graph among variable groups is given in Figure 4. In the same figure, the variables for which we used a log scaling to limit the effects of outliers are flagged as such. These include the dependent variable (n_cit_tot), which is always used on a log scale.

$$\begin{aligned}
 \log(n_cit_tot + 1) = & \log(n_authors + 1) + \log(n_references + 1) + p_year + p_month + \\
 & \log(h_index_mean + 1) + \\
 & is_plos + is_plos_one + \\
 & data_shared + data_location + repositories_data + \\
 & code_shared + code_location + \\
 & preprint_match
 \end{aligned} \tag{1}$$

$$\begin{aligned}
 \log(n_{cit_tot} + 1) = & \log(n_{authors} + 1) + \log(n_{references} + 1) + p_{year} + p_{month} + \\
 & \log(h_{index_mean} + 1) + \\
 & is_plos + is_plos_one + \\
 & data_shared + data_location + repositories_data + \\
 & code_shared + code_location + \\
 & preprint_match + \\
 & \sum_{i=1}^{18} I(\text{division} = i)
 \end{aligned} \tag{2}$$

Figure 4: An illustration of the assumed causal dependency graph among dependent and independent variables. We distinguish among the dependent variable and its variations (red), independent control variables (blue), and OSI control variables (green). We are interested in the total effect of OSI variables on the dependent variable (n_{cit_tot}), shown by the thick black line, and in controlling for the effect of other independent variables, shown by the dotted black lines.



Starting with the base model in Table 5, we provide results for an OLS model and a robust linear model as a comparison. The results are aligned and show a relatively high explained variance with the base model having $R^2 = .408$. The model shows expected trends, as previously discussed. For example, the higher the year the lower the total citation count on average (about -30% per year increase), or the higher the average H-index of the authors, the higher the citation counts of the paper (this can be interpreted as an elasticity in a log-log model, therefore a 1% increase in the average H-index leads to a $.141\%$ increase in the number of citations, on average). More of interest to us are the OSIs. These show that there is a significant and positive effect of preprints (20.4%) and of sharing data via an online repository (3.9%). These percentage changes for log-linear relationships are calculated as follows: $(\exp(.186) - 1) \times 100 \approx 20.4\%$. These effects are cumulative, so a publication with both a preprint and data shared in a repository would be associated with an average citation increase of 24.3% . On the other hand, the OSI for code sharing did not yield a statistically significant positive citation effect. Our next question is whether these results hold when we account for the large disciplinary variations in the adoption of Open Science practices, which we assess next.

Table 5: Results for the base model.

	Dependent variable:	
	n_cit_tot_log	
	OLS	robust linear
	(1)	(2)
n_authors_log	0.265***	0.254***

	(0.005)	(0.004)
n_references_tot_log	0.192*** (0.005)	0.198*** (0.005)
p_year	-0.357*** (0.001)	-0.370*** (0.001)
p_month	-0.037*** (0.001)	-0.039*** (0.001)
h_index_mean_log	0.141*** (0.003)	0.141*** (0.003)
C(is_plos)True	0.095*** (0.009)	0.107*** (0.008)
C(is_plos_one)True	-0.347*** (0.007)	-0.351*** (0.007)
C(Data_Shared)True	-0.011 (0.034)	-0.012 (0.033)
C(Data_Location)Online	-0.0002 (0.034)	-0.001 (0.034)
C(Data_Location)Supplementary Information	0.024 (0.034)	0.027 (0.033)
C(Repositories_data_bool)True	0.038*** (0.008)	0.038*** (0.008)
C(Code_Shared)True	0.057 (0.107)	0.064 (0.104)
C(Code_Location)Online	-0.105 (0.107)	-0.127 (0.105)
C(Code_Location)Supplementary Information	-0.130 (0.108)	-0.132 (0.105)
C(Preprint_Match)True	0.186*** (0.006)	0.159*** (0.006)
Constant	720.843*** (2.705)	746.799*** (2.638)
Observations	121,999	121,999
R ²	0.408	
Adjusted R ²	0.408	
Residual Std. Error (df = 121983)	0.775	0.723
F Statistic	5,615.550*** (df = 15; 121983)	

Note:

*p<0.1; **p<0.05; ***p<0.01

The full model in Table 6 adds the ANZSRC divisions as 18 dummy variables. The model shows an even higher explained variance with $R^2 = 0.426$. The full model shows trends that largely confirm the results from the base model. We consolidate our estimate for the citation impact of OSI indicators as follows. We find that the early release of a publication as a preprint correlates with a significant positive citation advantage of about 20.2% (± 0.7) on average. We

also find that sharing data in an online repository is associated with a smaller yet still positive citation advantage of 4.3% ($\pm .8$) on average. These effects are cumulative, so a publication with both a preprint and data shared in a repository would be associated with an average citation increase of 24.5%. We do not find a significant effect for sharing code, and we detect significant variations across disciplines in average citation impact. All the remaining coefficients are confirmed in sign and, with minor variation, in magnitude.

Table 6: Results for the full model.

	<i>Dependent variable:</i>	
	n_cit_tot_log	
	<i>OLS</i>	<i>robust linear</i>
	(1)	(2)
n_authors_log	0.207*** (0.005)	0.194*** (0.005)
n_references_tot_log	0.246*** (0.005)	0.252*** (0.005)
p_year	-0.357*** (0.001)	-0.368*** (0.001)
p_month	-0.037*** (0.001)	-0.038*** (0.001)
h_index_mean_log	0.119*** (0.003)	0.120*** (0.003)
C(is_plos)True	0.058*** (0.009)	0.070*** (0.008)
C(is_plos_one)True	-0.304*** (0.007)	-0.306*** (0.007)
C(Data_Shared)True	0.002 (0.033)	-0.005 (0.033)
C(Data_Location)Online	0.010 (0.034)	0.015 (0.033)
C(Data_Location)Supplementary Information	0.020 (0.033)	0.027 (0.033)
C(Repositories_data_bool)True	0.042*** (0.008)	0.041*** (0.008)
C(Code_Shared)True	0.099 (0.105)	0.107 (0.103)
C(Code_Location)Online	-0.110 (0.106)	-0.134 (0.103)
C(Code_Location)Supplementary Information	-0.147 (0.106)	-0.154 (0.104)
C(Preprint_Match)True	0.184*** (0.006)	0.155*** (0.006)

C(division_1)True	0.126*** (0.006)	0.128*** (0.006)
C(division_2)True	0.023*** (0.006)	0.031*** (0.006)
C(division_3)True	0.099*** (0.006)	0.089*** (0.006)
C(division_4)True	0.018* (0.009)	0.030*** (0.009)
C(division_5)True	-0.086*** (0.010)	-0.084*** (0.010)
C(division_6)True	-0.075*** (0.010)	-0.087*** (0.010)
C(division_7)True	0.031*** (0.011)	0.020* (0.011)
C(division_8)True	-0.186*** (0.011)	-0.172*** (0.011)
C(division_9)True	-0.301*** (0.012)	-0.294*** (0.011)
C(division_10)True	0.035** (0.014)	0.045*** (0.013)
C(division_11)True	-0.059*** (0.015)	-0.046*** (0.015)
C(division_12)True	-0.235*** (0.015)	-0.238*** (0.015)
C(division_13)True	-0.137*** (0.018)	-0.133*** (0.018)
C(division_14)True	-0.147*** (0.021)	-0.138*** (0.021)
C(division_15)True	-0.037 (0.023)	-0.017 (0.023)
C(division_16)True	-0.057** (0.025)	-0.064*** (0.024)
C(division_17)True	-0.173*** (0.026)	-0.164*** (0.025)
C(division_18)True	-0.111*** (0.015)	-0.112*** (0.014)
Constant	721.097*** (2.698)	744.009*** (2.635)
Observations	121,999	121,999
R ²	0.426	

Adjusted R ²	0.426	
Residual Std. Error (df = 121965)	0.764	0.714
F Statistic	2,739.270*** (df = 33; 121965)	
<hr/>		
<i>Note:</i>	*p<0.1; **p<0.05; ***p<0.01	

5 Discussion

This study offers a comprehensive analysis of the citation impact of Open Science practices, drawing on a dataset of about 122'000 research articles and using both descriptive and regression analysis. Our findings reveal a consistent citation advantage for articles whose authors adopted Open Science practices, including data sharing in online repositories and preprint posting. This correlation suggests that Open Science practices may significantly enhance the visibility and academic impact of research findings. However, the Open Science practice of sharing code does not seem to lead to a citation advantage in our sample.

5.1 Limitations

Several limitations of our study should be acknowledged. First, while our dataset is extensive, it is heavily weighted toward publications by the Open Access publisher PLOS, and as such it may not fully capture the diversity of research across all fields, potentially limiting the generalizability of our findings. Furthermore, PLOS champions Open Science practices, and the stance that a publisher takes in this regard may have an influence on the observed effects. In particular, PLOS requires all authors, with limited exceptions, to share the research data supporting their articles as a condition of publication, with the use of data repositories as the preferred approach. This is reflected in the higher overall rates of data sharing, and higher rates of data repository use in PLOS articles compared to comparators in the OSI dataset. As data sharing is the norm in PLOS articles and as the use of repositories is not uncommon, a citation advantage for the use of data repositories may be smaller in PLOS articles compared to non-PLOS articles. Posting preprints, however, is an optional practice for researchers publishing with PLOS and most other journals. Code sharing, similarly, is optional in most of the journals in our sample, with rare exceptions such as PLOS Computational Biology, where this practice is mandatory [64].

Additionally, the observational nature of our study precludes definitive conclusions about causality. The observed citation advantage might be influenced by other factors not accounted for in our analysis, such as the intrinsic quality of the research or access to research funding.

5.2 Extension of previous research

The model-explained variance in our results is globally high with respect to similar studies. For instance, there is previous work showing a positive correlation between citation and altmetric impact of publications, and the posting of preprints [38, 39, 40, 41]. The extent of the citation advantage, previously found to be as much as fivefold, is known to vary according to the timing of preprint posting, the discipline, and the preprint server used, among other factors. The smaller magnitude of the effect we find relative to previous studies may relate to the broader range of preprint servers that our sample considers.

Using similar methods to ours, previous work also found a correlation between articles that include statements linking to data in a repository and a citation advantage of up to 25% [13]. We confirm this finding in our study, finding a positive correlation between sharing data in a repository and citation impact. Yet the effect we find is considerably smaller in magnitude. This might be caused by the smaller and more uniform dataset that we use here, which includes all PLOS publications and a smaller comparator set, while this previous work used all PLOS and BMC articles and a dataset of over half a million publications. Other studies have also found a positive citation impact of the use of discipline-specific repositories [27, 28, 29].

While previous work [34, 36, 37] has found as much as a threefold citation advantage for code sharing, we did not confirm this finding in our sample. Following [65], it is possible that outside of fields like computer science authors are more likely to cite or link to shared code directly rather than citing the research paper with which it was associated.

5.3 Implications for future research

Our data and code are shared openly to enable independent replication of our results and extension of our findings as larger or different, comparable sources of data on the adoption of Open Science practices become available. This

includes future versions of the PLOS OSI dataset, as well as outputs from other Open Science monitoring initiatives, such as the French Open Science Monitor ¹ or OpenAIRE ².

As Open Science practices and policies continue to develop, future research could explore longitudinal changes in citation patterns. Further studies could also investigate the relationship between additional Open Science practices and citation impact, extending our understanding of how different aspects of openness contribute to research visibility. Moreover, it would be valuable to examine the impact of Open Science practices on other domains of research dissemination and engagement, such as open commons (e.g., Wikipedia), public policy influence, collaboration networks, and public engagement. We might hypothesize, for example, that non-citation measures of impact – such as forks and downloads – may be more relevant for the sharing of code and software. Contemporary calls for the reform of research assessment (such as <https://coara.eu>) emphasize valuing more diverse research outputs and contributions, as well as more diverse measures of impact. These developments underscore the importance of future research exploring the association of Open Science practices with effects other than citations.

6 Conclusion

In summary, our study contributes to the growing body of literature on the effects or impacts of Open Science by quantifying the citation impact of data sharing, code sharing, and preprint posting. Our results could be readily extended with additional data on Open Science practices detected in a larger sample of non-PLOS Open Access publications. We advocate for further empirical research to build on these findings, particularly work that focuses on causal mechanisms, discipline-specific effects, and broader impacts beyond citation metrics.

Data and Code Availability

OSI dataset: <https://doi.org/10.6084/m9.figshare.21687686.v5>

Code (GitHub): <https://github.com/MediaComem/das-public>

Data (Zenodo): <https://zenodo.org/doi/10.5281/zenodo.10134811>

Funding

PLOS provided funding for data acquisition, modeling, and analysis, and had a role in the study design, analysis, and preparation of the manuscript. PLOS also provided support in the form of salaries for authors LC, ML, and IH.

Competing interests

Three of the authors (LC, ML, and IH) were at the time of publication employed by PLOS, the publisher of PLOS ONE.

Authors' contributions

- GC: Conceptualization, Data curation, Formal analysis, Investigation, Methodology, Software, Supervision, Writing – original draft, Writing – review & editing.
- LC: Formal Analysis, Visualization, Writing – original draft.
- ML: Writing – original draft, Writing – review & editing.
- GD: Data curation, Software, Writing – review & editing.
- SL: Supervision, Writing – review & editing.
- DR: Supervision, Writing – review & editing.
- IH: Conceptualization, Funding acquisition, Methodology, Project administration, Resources, Supervision, Writing – original draft, Writing – review & editing.

¹<https://data.enseignementsup-recherche.gouv.fr/explore/dataset/open-access-monitor-france>.

²<https://monitor.openaire.eu>.

Acknowledgements

We thank Tim Vines, Scott Kerr, Souad McIntosh, and the team at DataSeer for their collaboration in enabling the OSI dataset to be used for this analysis. We also thank Ross Gray at PLOS for reviewing the data and code from our study.

Appendix

We show in this Appendix results for a few more models in order to further confirm our results. Firstly, a base model adding code generated as a variable shows a small yet significant negative effect related to it (Table 7). This effect goes away when controlling for disciplines, therefore we consider it spurious. When considering OSI interactions (Table 8), we find a further negative effect provided by code generated and code shared. This surprising result may be an artifact of the dataset, that we are unsure how to explain. Next, we show how different preprint servers are associated with varying degrees of citation impact (Table 9). Lastly, we check a full model using as dependent variables the citation counts up to 1 year after publication (Table 10). We still find the same results as using the full citation counts, albeit with a smaller magnitude as expected.

Table 7: Results for the base model with code generated OSI.

	<i>Dependent variable:</i>	
	n_cit_tot_log	
	<i>OLS</i>	<i>robust linear</i>
	(1)	(2)
n_authors_log	0.266*** (0.005)	0.255*** (0.004)
n_references_tot_log	0.195*** (0.005)	0.200*** (0.005)
p_year	-0.357*** (0.001)	-0.370*** (0.001)
p_month	-0.037*** (0.001)	-0.039*** (0.001)
h_index_mean_log	0.141*** (0.003)	0.142*** (0.003)
C(is_plos)True	0.095*** (0.009)	0.106*** (0.008)
C(is_plos_one)True	-0.348*** (0.007)	-0.352*** (0.007)
C(Data_Shared)True	-0.005 (0.034)	-0.007 (0.033)
C(Data_Location)Online	-0.003 (0.034)	-0.003 (0.034)
C(Data_Location)Supplementary Information	0.020 (0.034)	0.024 (0.033)
C(Repositories_data_bool)True	0.041*** (0.008)	0.040*** (0.008)
C(Code_Generated)True	-0.022***	-0.017***

	(0.005)	(0.005)
C(Code_Shared)True	0.070 (0.107)	0.075 (0.104)
C(Code_Location)Online	-0.111 (0.107)	-0.132 (0.105)
C(Code_Location)Supplementary Information	-0.137 (0.108)	-0.138 (0.105)
C(Preprint_Match)True	0.188*** (0.006)	0.160*** (0.006)
Constant	720.947*** (2.705)	746.881*** (2.638)
Observations	121,999	121,999
R ²	0.409	
Adjusted R ²	0.408	
Residual Std. Error (df = 121982)	0.775	0.723
F Statistic	5,266.481*** (df = 16; 121982)	
<i>Note:</i>	*p<0.1; **p<0.05; ***p<0.01	

Table 8: Results for the base model with interactions among OSI.

	<i>Dependent variable:</i>	
	n_cit_tot_log	
	<i>OLS</i>	<i>robust linear</i>
	(1)	(2)
n_authors_log	0.266*** (0.005)	0.254*** (0.004)
n_references_tot_log	0.195*** (0.005)	0.200*** (0.005)
p_year	-0.357*** (0.001)	-0.370*** (0.001)
p_month	-0.037*** (0.001)	-0.039*** (0.001)
h_index_mean_log	0.141*** (0.003)	0.142*** (0.003)
C(is_plos)True	0.096*** (0.009)	0.107*** (0.008)
C(is_plos_one)True	-0.350*** (0.007)	-0.352*** (0.007)
C(Data_Shared)True	-0.007 (0.034)	-0.009 (0.033)
C(Data_Location)Online	-0.003	-0.003

	(0.034)	(0.034)
C(Data_Location)Supplementary Information	0.022 (0.034)	0.025 (0.033)
C(Repositories_data_bool)True	0.044*** (0.009)	0.041*** (0.009)
C(Preprint_Match)True	0.190*** (0.007)	0.160*** (0.007)
C(Code_Generated)True	-0.015*** (0.005)	-0.010* (0.005)
C(Code_Shared)True	0.144 (0.108)	0.147 (0.105)
C(Code_Location)Online	-0.120 (0.107)	-0.142 (0.105)
C(Code_Location)Supplementary Information	-0.166 (0.108)	-0.166 (0.105)
C(Repositories_data_bool)True:C(Preprint_Match)True	-0.007 (0.012)	0.002 (0.012)
C(Code_Generated)True:C(Code_Shared)True	-0.079*** (0.017)	-0.077*** (0.017)
Constant	720.952*** (2.705)	746.876*** (2.638)
Observations	121,999	121,999
R ²	0.409	
Adjusted R ²	0.409	
Residual Std. Error (df = 121980)	0.775	0.723
F Statistic	4,683.261*** (df = 18; 121980)	

Note: *p<0.1; **p<0.05; ***p<0.01

Table 9: Results for the base model with preprint servers (considering only those mentioned in 500 or more publications part of the dataset).

	<i>Dependent variable:</i>	
	n_cit_tot_log	
	<i>OLS</i>	<i>robust linear</i>
	(1)	(2)
n_authors_log	0.259*** (0.005)	0.250*** (0.004)
n_references_tot_log	0.201*** (0.005)	0.204*** (0.005)
p_year	-0.360*** (0.001)	-0.372*** (0.001)

p_month	-0.037*** (0.001)	-0.039*** (0.001)
h_index_mean_log	0.143*** (0.003)	0.143*** (0.003)
C(is_plos)True	0.094*** (0.009)	0.106*** (0.009)
C(is_plos_one)True	-0.344*** (0.007)	-0.346*** (0.007)
C(Data_Shared)True	-0.010 (0.034)	-0.013 (0.033)
C(Data_Location)Online	-0.001 (0.034)	0.002 (0.034)
C(Data_Location)Supplementary Information	0.027 (0.034)	0.030 (0.033)
C(Repositories_data_bool)True	0.043*** (0.009)	0.041*** (0.008)
C(Code_Shared)True	0.060 (0.107)	0.059 (0.105)
C(Code_Location)Online	-0.107 (0.108)	-0.119 (0.105)
C(Code_Location)Supplementary Information	-0.131 (0.108)	-0.125 (0.106)
C(Preprint_Match)True	0.689** (0.345)	0.264 (0.338)
C(Preprint_Server)bioRxiv	0.189*** (0.027)	0.191*** (0.027)
C(Preprint_Server)Journal of Medical Internet Research	0.517*** (0.038)	0.480*** (0.037)
C(Preprint_Server)medRxiv	0.470*** (0.030)	0.381*** (0.029)
C(Preprint_Server)N/A	0.721** (0.346)	0.309 (0.339)
C(Preprint_Server)Protocols.io	-0.047 (0.043)	-0.040 (0.042)
C(Preprint_Server)PsyArXiv	0.181*** (0.039)	0.139*** (0.038)
C(Preprint_Server)Research Square	0.191*** (0.029)	0.194*** (0.028)
Constant	726.760*** (2.754)	751.681*** (2.695)

Observations	120,195	120,195
R ²	0.413	
Adjusted R ²	0.413	
Residual Std. Error (df = 120172)	0.772	0.721
F Statistic	3,838.518*** (df = 22; 120172)	

Note:

*p<0.1; **p<0.05; ***p<0.01

Table 10: Results for the full model with dependent variable as citation data for 1 year from publication.

	<i>Dependent variable:</i>	
	n_cit_1_log	
	<i>OLS</i>	<i>robust linear</i>
	(1)	(2)
n_authors_log	0.131*** (0.004)	0.109*** (0.003)
n_references_tot_log	0.113*** (0.004)	0.117*** (0.004)
p_year	0.020*** (0.001)	0.015*** (0.001)
p_month	-0.013*** (0.001)	-0.013*** (0.001)
h_index_mean_log	0.052*** (0.003)	0.050*** (0.002)
C(is_plos)True	0.129*** (0.007)	0.149*** (0.006)
C(is_plos_one)True	-0.303*** (0.006)	-0.302*** (0.005)
C(Data_Shared)True	0.015 (0.025)	0.015 (0.023)
C(Data_Location)Online	0.007 (0.025)	0.008 (0.023)
C(Data_Location)Supplementary Information	0.003 (0.025)	0.005 (0.023)
C(Repositories_data_bool)True	0.012* (0.006)	0.011* (0.006)
C(Code_Shared)True	0.084 (0.081)	0.075 (0.075)
C(Code_Location)Online	-0.091 (0.081)	-0.092 (0.075)
C(Code_Location)Supplementary Information	-0.101	-0.087

	(0.082)	(0.075)
C(Preprint_Match)True	0.133*** (0.004)	0.099*** (0.004)
C(division_1)True	0.047*** (0.004)	0.049*** (0.004)
C(division_2)True	0.019*** (0.005)	0.029*** (0.004)
C(division_3)True	0.003 (0.005)	-0.001 (0.004)
C(division_4)True	-0.015** (0.007)	-0.007 (0.006)
C(division_5)True	-0.037*** (0.008)	-0.030*** (0.007)
C(division_6)True	-0.041*** (0.008)	-0.044*** (0.007)
C(division_7)True	0.006 (0.009)	0.0003 (0.008)
C(division_8)True	-0.073*** (0.009)	-0.054*** (0.008)
C(division_9)True	-0.119*** (0.009)	-0.104*** (0.008)
C(division_10)True	0.045*** (0.010)	0.056*** (0.009)
C(division_11)True	-0.026** (0.012)	-0.019* (0.011)
C(division_12)True	-0.056*** (0.012)	-0.057*** (0.011)
C(division_13)True	-0.037** (0.014)	-0.030** (0.013)
C(division_14)True	-0.072*** (0.016)	-0.062*** (0.015)
C(division_15)True	-0.024 (0.018)	-0.012 (0.017)
C(division_16)True	0.010 (0.019)	0.0002 (0.018)
C(division_17)True	-0.056*** (0.020)	-0.045** (0.018)
C(division_18)True	-0.032*** (0.011)	-0.033*** (0.010)

Constant	−40.101*** (2.385)	−29.688*** (2.198)
Observations	106,733	106,733
R ²	0.120	
Adjusted R ²	0.120	
Residual Std. Error (df = 106699)	0.542	0.522
F Statistic	442.350*** (df = 33; 106699)	
Note:	*p<0.1; **p<0.05; ***p<0.01	

References

- [1] John Willinsky. “The Unacknowledged Convergence of Open Source, Open Access, and Open Science”. In: *First Monday* (Aug. 2005). ISSN: 1396-0466. DOI: 10.5210/fm.v10i8.1265. (Visited on 04/11/2024).
- [2] Nathaniel Tkacz. *Wikipedia and the Politics of Openness*. Chicago: University of Chicago Press, 2014.
- [3] Samuel A. Moore. “A Genealogy of Open Access: Negotiations between Openness and Access to Research”. In: *Revue française des sciences de l’information et de la communication* 11 (Aug. 2017). ISSN: 2263-0856. DOI: 10.4000/rfsic.3220. (Visited on 04/11/2024).
- [4] Ruben Vicente-Saez and Clara Martinez-Fuentes. “Open Science Now: A Systematic Literature Review for an Integrated Definition”. In: *Journal of Business Research* 88 (July 2018), pp. 428–436. ISSN: 01482963. DOI: 10.1016/j.jbusres.2017.12.043. (Visited on 02/27/2024).
- [5] Sabina Leonelli. *Philosophy of Open Science*. Cambridge University Press, Sept. 2023. ISBN: 978-1-00-941636-8 978-1-00-941639-9. DOI: 10.1017/9781009416368. (Visited on 04/11/2024).
- [6] UNESCO. *UNESCO Recommendation on Open Science*. Tech. rep. Paris: UNESCO, 2021. DOI: 10.54677/MNMH8546. (Visited on 04/11/2024).
- [7] Stylianos Serghiou et al. “Assessment of Transparency Indicators across the Biomedical Literature: How Open Is Open?” In: *PLOS Biology* 19.3 (Mar. 2021). Ed. by Lisa Bero, e3001107. ISSN: 1545-7885. DOI: 10.1371/journal.pbio.3001107. (Visited on 04/11/2024).
- [8] Joe Menke et al. “Establishing Institutional Scores with the Rigor and Transparency Index: Large-scale Analysis of Scientific Reporting Quality”. In: *Journal of Medical Internet Research* 24.6 (June 2022), e37324. ISSN: 1438-8871. DOI: 10.2196/37324. (Visited on 04/11/2024).
- [9] Samuel G. Robson et al. “Promoting Open Science: A Holistic Approach to Changing Behaviour”. In: *Collabra: Psychology* 7.1 (Dec. 2021), p. 30137. ISSN: 2474-7394. DOI: 10.1525/collabra.30137. (Visited on 04/11/2024).
- [10] Kristijan Armeni et al. “Towards Wide-Scale Adoption of Open Science Practices: The Role of Open Science Communities”. In: *Science and Public Policy* 48.5 (Oct. 2021), pp. 605–611. ISSN: 0302-3427, 1471-5430. DOI: 10.1093/scipol/scab039. (Visited on 12/20/2023).
- [11] Iain Hrynaszkiewicz and Lauren Cadwallader. *A Survey of Funders’ and Institutions’ Needs for Understanding Researchers’ Open Research Practices*. Preprint. Open Science Framework, Sept. 2021. DOI: 10.31219/osf.io/z4py9. (Visited on 04/11/2024).
- [12] Iain Hrynaszkiewicz and Veronique Kiermer. “PLOS Open Science Indicators Principles and Definitions”. In: (2022). DOI: 10.6084/m9.figshare.21640889.v1. (Visited on 12/27/2023).
- [13] Giovanni Colavizza et al. “The Citation Advantage of Linking Publications to Research Data”. In: *PLOS ONE* 15.4 (Apr. 2020). Ed. by Jelte M. Wicherts, e0230416. ISSN: 1932-6203. DOI: 10.1371/journal.pone.0230416. (Visited on 12/15/2023).
- [14] Thomas Klebel et al. “PathOS Deliverable 1.2: Scoping Review of Open Science Impact”. In: (Feb. 2024). (Visited on 03/27/2024).
- [15] Chun-Kai Huang et al. “Open Access Research Outputs Receive More Diverse Citations”. In: *Scientometrics* (Jan. 2024). ISSN: 0138-9130, 1588-2861. DOI: 10.1007/s11192-023-04894-0. (Visited on 03/09/2024).
- [16] Teresa Schultz. “All the Research That’s Fit to Print: Open Access and the News Media”. In: *Quantitative Science Studies* 2.3 (Nov. 2021), pp. 828–844. ISSN: 2641-3337. DOI: 10.1162/qss_a_00139. (Visited on 03/26/2024).
- [17] Michael Woelfle, Piero Olliaro, and Matthew H. Todd. “Open Science Is a Research Accelerator”. In: *Nature Chemistry* 3.10 (Oct. 2011), pp. 745–748. ISSN: 1755-4330, 1755-4349. DOI: 10.1038/nchem.1149. (Visited on 02/27/2024).

- [18] Lonni Besançon et al. “Open Science Saves Lives: Lessons from the COVID-19 Pandemic”. In: *BMC Medical Research Methodology* 21.1 (June 2021), p. 117. ISSN: 1471-2288. DOI: 10.1186/s12874-021-01304-y. (Visited on 02/27/2024).
- [19] Open Science Collaboration. “Estimating the Reproducibility of Psychological Science”. In: *Science* 349.6251 (Aug. 2015), aac4716. ISSN: 0036-8075, 1095-9203. DOI: 10.1126/science.aac4716. (Visited on 02/27/2024).
- [20] Nicki Lisa Cole et al. *The Societal Impact of Open Science—a Scoping Review*. Preprint. SocArXiv, Feb. 2024. DOI: 10.31235/osf.io/tqrwg. (Visited on 03/09/2024).
- [21] Michael J. Fell. “The Economic Impacts of Open Science: A Rapid Evidence Assessment”. In: *Publications* 7.3 (July 2019), p. 46. ISSN: 2304-6775. DOI: 10.3390/publications7030046. (Visited on 03/27/2024).
- [22] Directorate-General for Research and Innovation (European Commission) and PwC EU Services. *Cost-Benefit Analysis for FAIR Research Data: Cost of Not Having FAIR Research Data*. Tech. rep. Luxembourg: Publications Office of the European Union, 2018. (Visited on 03/27/2024).
- [23] Thomas J. Hostler. “The Invisible Workload of Open Research”. In: *Journal of Trial and Error* (May 2023). ISSN: 2667-1204. DOI: 10.36850/mr5. (Visited on 03/09/2024).
- [24] Tony Ross-Hellauer et al. “Dynamics of Cumulative Advantage and Threats to Equity in Open Science: A Scoping Review”. In: *Royal Society Open Science* 9.1 (Jan. 2022), p. 211032. ISSN: 2054-5703. DOI: 10.1098/rsos.211032. (Visited on 03/28/2024).
- [25] Sabina Leonelli. “Open Science and Epistemic Diversity: Friends or Foes?” In: *Philosophy of Science* 89.5 (Dec. 2022), pp. 991–1001. ISSN: 0031-8248, 1539-767X. DOI: 10.1017/psa.2022.45. (Visited on 03/26/2024).
- [26] Yulin Yu and Daniel M. Romero. *Does the Use of Unusual Combinations of Datasets Contribute to Greater Scientific Impact?* Feb. 2024. DOI: 10.48550/arXiv.2402.05024. arXiv: 2402.05024 [cs, econ, q-fin]. (Visited on 02/08/2024).
- [27] Heather A. Piwowar, Roger S. Day, and Douglas B. Fridsma. “Sharing Detailed Research Data Is Associated with Increased Citation Rate”. In: *PLOS ONE* 2.3 (Mar. 2007). Ed. by John Ioannidis, e308. ISSN: 1932-6203. DOI: 10.1371/journal.pone.0000308. (Visited on 03/09/2024).
- [28] Heather A. Piwowar and Todd J. Vision. “Data Reuse and the Open Data Citation Advantage”. In: *PeerJ* 1 (Oct. 2013), e175. ISSN: 2167-8359. DOI: 10.7717/peerj.175. (Visited on 03/09/2024).
- [29] Edwin A. Henneken and Alberto Accomazzi. *Linking to Data: Effect on Citation Rates in Astronomy*. Nov. 2011. DOI: 10.48550/arxiv.1111.3618. arXiv: 1111.3618 [astro-ph]. (Visited on 03/09/2024).
- [30] Garret Christensen et al. “A Study of the Impact of Data Sharing on Article Citations Using Journal Policies as a Natural Experiment”. In: *PLOS ONE* 14.12 (Dec. 2019). Ed. by Florian Naudet, e0225883. ISSN: 1932-6203. DOI: 10.1371/journal.pone.0225883. (Visited on 03/09/2024).
- [31] Lisa M. Federer. “Long-Term Availability of Data Associated with Articles in PLOS ONE”. In: *PLOS ONE* 17.8 (Aug. 2022). Ed. by Jelte M. Wicherts, e0272845. ISSN: 1932-6203. DOI: 10.1371/journal.pone.0272845. (Visited on 03/09/2024).
- [32] Leho Tedersoo et al. “Data Sharing Practices and Data Availability upon Request Differ across Scientific Disciplines”. In: *Scientific Data* 8.1 (July 2021), p. 192. ISSN: 2052-4463. DOI: 10.1038/s41597-021-00981-0. (Visited on 03/09/2024).
- [33] David Stuart et al. *Practical Challenges for Researchers in Data Sharing*. Tech. rep. Springer Nature, 2018. (Visited on 03/27/2024).
- [34] Patrick Vandewalle. “Code Sharing Is Associated with Research Impact in Image Processing”. In: *Computing in Science & Engineering* 14.4 (July 2012), pp. 42–47. ISSN: 1521-9615. DOI: 10.1109/MCSE.2012.63. (Visited on 03/19/2024).
- [35] Šimon Kucharský, Bobby Lee Houtkoop, and Ingmar Visser. *Code Sharing in Psychological Methods and Statistics: An Overview and Associations with Conventional and Alternative Research Metrics*. Feb. 2020. DOI: 10.31219/osf.io/daews. (Visited on 03/28/2024).
- [36] Hancheng Cao et al. *The Rise of Open Science: Tracking the Evolution and Perceived Value of Data and Methods Link-Sharing Practices*. 2023. DOI: 10.48550/arxiv.2310.03193. (Visited on 03/19/2024).
- [37] Donghyun Kang, TaeYoung Kang, and Junkyu Jang. “Papers with Code or without Code? Impact of GitHub Repository Usability on the Diffusion of Machine Learning Research”. In: *Information Processing & Management* 60.6 (Nov. 2023), p. 103477. ISSN: 03064573. DOI: 10.1016/j.ipm.2023.103477. (Visited on 03/20/2024).
- [38] Erin C McKiernan et al. “How Open Science Helps Researchers Succeed”. In: *eLife* 5 (July 2016), e16800. ISSN: 2050-084X. DOI: 10.7554/eLife.16800. (Visited on 02/27/2024).

- [39] Darwin Y Fu and Jacob J Hughey. “Releasing a Preprint Is Associated with More Attention and Citations for the Peer-Reviewed Article”. In: *eLife* 8 (Dec. 2019), e52646. ISSN: 2050-084X. DOI: 10.7554/eLife.52646. (Visited on 12/20/2023).
- [40] Nicholas Fraser et al. “The Relationship between bioRxiv Preprints, Citations and Altmetrics”. In: *Quantitative Science Studies* (Apr. 2020), pp. 1–21. ISSN: 2641-3337. DOI: 10.1162/qss_a_00043. (Visited on 12/20/2023).
- [41] Boya Xie, Zhihong Shen, and Kuansan Wang. *Is Preprint the Future of Science? A Thirty Year Journey of Online Preprint Services*. Feb. 2021. arXiv: 2102.09066 [cs]. (Visited on 12/20/2023).
- [42] Nicholas Fraser et al. “The Evolving Role of Preprints in the Dissemination of COVID-19 Research and Their Impact on the Science Communication Landscape”. In: *PLOS Biology* 19.4 (Apr. 2021). Ed. by Ulrich Dirnagl, e3000959. ISSN: 1545-7885. DOI: 10.1371/journal.pbio.3000959. (Visited on 12/20/2023).
- [43] Alice Fleerackers et al. “Unreviewed Science in the News: The Evolution of Preprint Media Coverage from 2014–2021”. In: *Quantitative Science Studies* (Jan. 2024), pp. 1–20. ISSN: 2641-3337. DOI: 10.1162/qss_a_00282. (Visited on 03/19/2024).
- [44] Narmin Rzayeva et al. “The Experiences of COVID-19 Preprint Authors: A Survey of Researchers about Publishing and Receiving Feedback on Their Work during the Pandemic”. In: *PeerJ* 11 (Aug. 2023), e15864. ISSN: 2167-8359. DOI: 10.7717/peerj.15864. (Visited on 03/19/2024).
- [45] Clarissa França Dias Carneiro et al. “Characterization of Comments about bioRxiv and medRxiv Preprints”. In: *JAMA Network Open* 6.8 (Aug. 2023), e2331410. ISSN: 2574-3805. DOI: 10.1001/jamanetworkopen.2023.31410. (Visited on 03/19/2024).
- [46] Sarvenaz Sarabipour et al. “On the Value of Preprints: An Early Career Researcher Perspective”. In: *PLOS Biology* 17.2 (Feb. 2019), e3000151. ISSN: 1545-7885. DOI: 10.1371/journal.pbio.3000151. (Visited on 12/20/2023).
- [47] Jesse F. Wolf et al. “Preprinting Is Positively Associated with Early Career Researcher Status in Ecology and Evolution”. In: *Ecology and Evolution* 11.20 (Oct. 2021), pp. 13624–13632. ISSN: 2045-7758, 2045-7758. DOI: 10.1002/ece3.8106. (Visited on 12/20/2023).
- [48] Mauricia Davidson et al. “No Evidence of Important Difference in Summary Treatment Effects between COVID-19 Preprints and Peer-Reviewed Publications: A Meta-Epidemiological Study”. In: *Journal of Clinical Epidemiology* 162 (Oct. 2023), pp. 90–97. ISSN: 08954356. DOI: 10.1016/j.jclinepi.2023.08.011. (Visited on 03/19/2024).
- [49] Dena Zeraatkar et al. “Consistency of Covid-19 Trial Preprints with Published Reports and Impact for Decision Making: Retrospective Review”. In: *BMJ Medicine* 1.1 (Oct. 2022), e000309. ISSN: 2754-0413. DOI: 10.1136/bmjmed-2022-000309. (Visited on 03/19/2024).
- [50] Michele Avissar-Whiting. “Downstream Retraction of Preprinted Research in the Life and Medical Sciences”. In: *PLOS ONE* 17.5 (May 2022). Ed. by Frederique Lisacek, e0267971. ISSN: 1932-6203. DOI: 10.1371/journal.pone.0267971. (Visited on 03/19/2024).
- [51] Public Library of Science. *PLOS Open Science Indicators (Version 5)*. 2023. DOI: 10.6084/m9.figshare.21687686.v5. (Visited on 12/15/2023).
- [52] Bethesda (MD) National Library of Medicine. *PMC Open Access Subset*. Oct. 2023.
- [53] Yassine Gargouri et al. “Self-Selected or Mandated, Open Access Increases Citation Impact for Higher Quality Research”. In: *PLOS ONE* 5.10 (Oct. 2010). Ed. by Robert P. Futrelle, e13636. ISSN: 1932-6203. DOI: 10.1371/journal.pone.0013636. (Visited on 12/15/2023).
- [54] Alfredo Yegros-Yegros, Ismael Rafols, and Pablo D’Este. “Does Interdisciplinary Research Lead to Higher Citation Impact? The Different Effect of Proximal and Distal Interdisciplinarity”. In: *PLOS ONE* 10.8 (Aug. 2015). Ed. by Wolfgang Glanzel, e0135095. ISSN: 1932-6203. DOI: 10.1371/journal.pone.0135095. (Visited on 12/15/2023).
- [55] Jian Wang, Reinhilde Veugelers, and Paula Stephan. “Bias against Novelty in Science: A Cautionary Tale for Users of Bibliometric Indicators”. In: *Research Policy* 46.8 (Oct. 2017), pp. 1416–1436. ISSN: 00487333. DOI: 10.1016/j.respol.2017.06.006. (Visited on 12/15/2023).
- [56] Vedran Sekara et al. “The Chaperone Effect in Scientific Publishing”. In: *Proceedings of the National Academy of Sciences* 115.50 (Dec. 2018), pp. 12603–12607. ISSN: 0027-8424, 1091-6490. DOI: 10.1073/pnas.1800471115. (Visited on 12/15/2023).
- [57] Vette I. Torvik and Neil R. Smalheiser. “Author Name Disambiguation in MEDLINE”. In: *ACM Transactions on Knowledge Discovery from Data* 3.3 (July 2009), pp. 1–29. ISSN: 1556-4681, 1556-472X. DOI: 10.1145/1552303.1552304. (Visited on 12/15/2023).
- [58] Z. Lu. “PubMed and beyond: A Survey of Web Tools for Searching Biomedical Literature”. In: *Database* 2011.0 (Jan. 2011), baq036–baq036. ISSN: 1758-0463. DOI: 10.1093/database/baq036. (Visited on 12/15/2023).

- [59] Anderson A. Ferreira, Marcos André Gonçalves, and Alberto H.F. Laender. “A Brief Survey of Automatic Methods for Author Name Disambiguation”. In: *ACM SIGMOD Record* 41.2 (Aug. 2012), pp. 15–26. ISSN: 0163-5808. DOI: 10.1145/2350036.2350040. (Visited on 12/15/2023).
- [60] Wanli Liu et al. “Author Name Disambiguation for PubMed”. In: *Journal of the Association for Information Science and Technology* 65.4 (Apr. 2014), pp. 765–781. ISSN: 2330-1635, 2330-1643. DOI: 10.1002/asi.23063. (Visited on 12/15/2023).
- [61] Jin G Zheng et al. “Entity Linking for Biomedical Literature”. In: *BMC Medical Informatics and Decision Making* 15.S1 (Dec. 2015), S4. ISSN: 1472-6947. DOI: 10.1186/1472-6947-15-S1-S4. (Visited on 12/15/2023).
- [62] Andreas Strotmann and Dangzhi Zhao. “Author Name Disambiguation: What Difference Does It Make in Author-based Citation Analysis?” In: *Journal of the American Society for Information Science and Technology* 63.9 (Sept. 2012), pp. 1820–1833. ISSN: 1532-2882, 1532-2890. DOI: 10.1002/asi.22695. (Visited on 12/15/2023).
- [63] Jinseok Kim and Jana Diesner. “Distortive Effects of Initial-based Name Disambiguation on Measurements of Large-scale Coauthorship Networks”. In: *Journal of the Association for Information Science and Technology* 67.6 (June 2016), pp. 1446–1461. ISSN: 2330-1635, 2330-1643. DOI: 10.1002/asi.23489. (Visited on 12/15/2023).
- [64] Lauren Cadwallader et al. “Advancing Code Sharing in the Computational Biology Community”. In: *PLOS Computational Biology* 18.6 (June 2022), e1010193. ISSN: 1553-7358. DOI: 10.1371/journal.pcbi.1010193. (Visited on 04/16/2024).
- [65] Emily Escamilla et al. “The Rise of GitHub in Scholarly Publications”. In: (2022). DOI: 10.48550/arxiv.2208.04895. (Visited on 04/12/2024).