

Evidence-based steps toward a culture for replicability and reproducibility

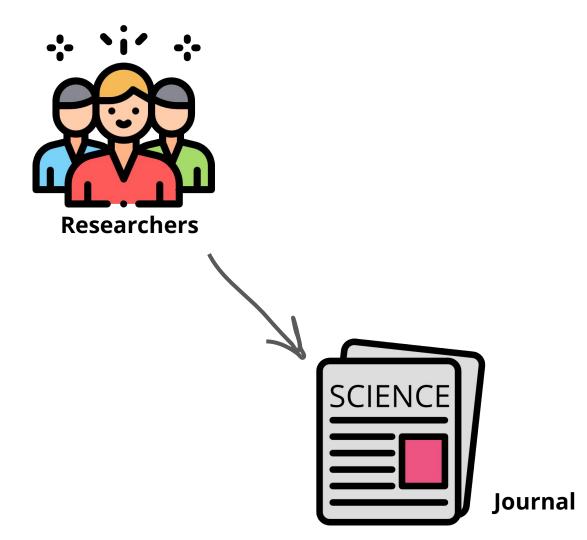
Ana Trisovic

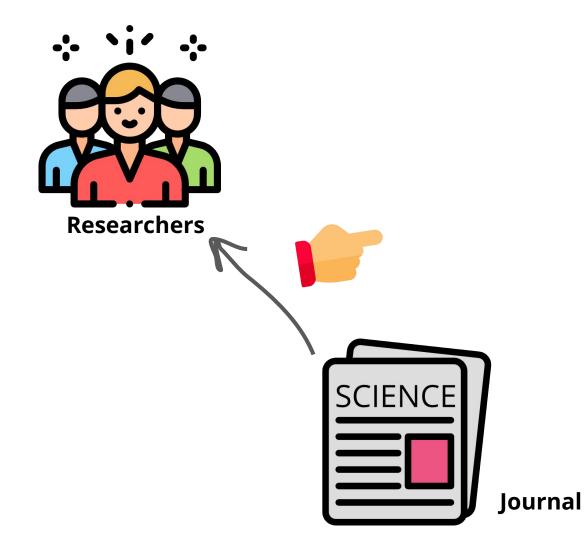
Harvard University

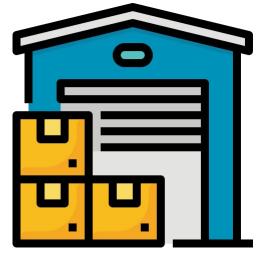




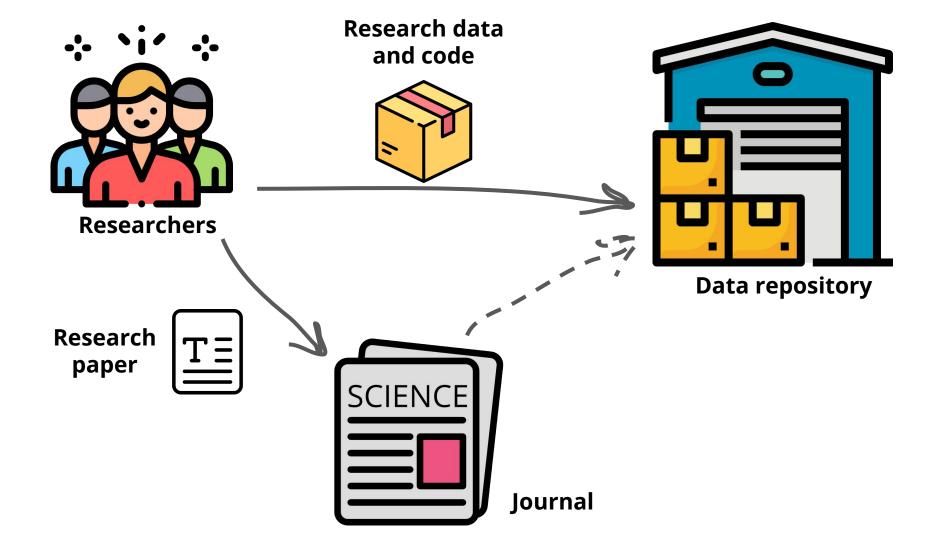








Data repository

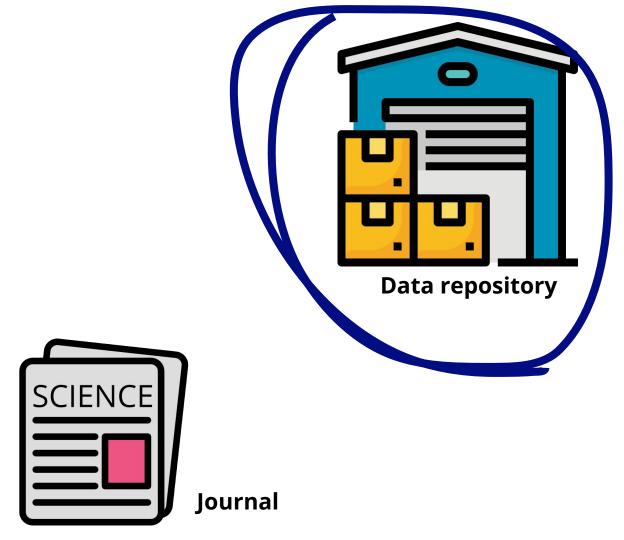


Presentation agenda

- 1. Dataverse research data repository
- 2. A large-scale study on code quality and execution
- 3. Results and discussion
- 4. What can researchers do
- 5. What can repositories do
- 6. What can journals do

Who are we?







- A free and open-source software platform to archive, share, and cite research data
 - Focus on data sharing and making data available
- Provides data repository software that can be installed at institutions
 - Supports research communities for entire countries (NO, NL)

70 institutions around the globe run Dataverse installations as their official data repository



Harvard Dataverse

https://dataverse.harvard.edu

		8 Harvard Dataverse				
HARVARD Dataverse		Add Data -	Search +	About User Guide	Support Sign Up	Log In
Deposit and share your data. Ge credit.	et academic	Organize datasets and gather metric own repository.	s in your	Publishing your data Dataverse!	is easy on Harvard	
Harvard Dataverse is a repository for res Deposit data and code here.	A dataverse is a container for all your datasets metadata.	s, files, and	Learn about getting starte repository here.	ed creating your own data	verse	
Add a dataset 🕂		Add a dataverse +		Getting started 🗗		
Find data across research fields,	, preview meta	data, and download files				

Data sharing

- Stand-alone or institutional account for depositing data (Dataverse, GitHub, Google, ORCID, University credentials)
- Individuals, institutes or journals may have own dataverse collections

HARVARD Dataverse		Add Data 👻	Search -	About	User Guide	Support
Host Dataverse 😡	Changing the host dataverse will clear any fiel Harvard Dataverse	ds you may have	entered data i	nto.		
*Asterisks indicate required fields						
Citation Metadata 🛧						
Title * 🕢	Enter title					
Author * 🕢	Add "Replication Data for" to Title Name *	Affiliation 3				
	Identifier Scheme 💿	Harvard Unive	ərsity			
Contact * 😔	Name 🕢	Affiliation 🕢				
	Trisovic, Ana E-mail * anatrisovic@fas.harvard.edu	Harvard Unive	ərsity		+	
Description * 🕢	This field supports only certain HTML tags. Text * €				+	
	Date 🕢					
	YYYY-MM-DD					
Subject * 🕢	Select				•	
Keyword 😔	Term 🕄	Vocabulary 😏			+	
	Vocabulary URL @ Enter full URL, starting with http://					
Related Publication 🕢	Citation 😡				+	

Data sharing

- Stand-alone or institutional account for depositing data (Dataverse, GitHub, Google, ORCID, University credentials)
- Individuals, institutes or journals may have own dataverse collections

Uataverse		Add Data 👻	Search +	About	User Guide	Support
Host Dataverse 9	Changing the host dataverse will clear any fie	lds you may have	entered data i	into.		
	Harvard Dataverse					
*Asterisks indicate required fields						
Citation Metadata A						
Title * 🕢	Enter title					
	Add "Replication Data for" to Title					
Author * 😧	Name * 📀	Affiliation 📀				
	Trisovic, Ana	Harvard Univ	ersity		+	
	Identifier Scheme 🕄	Identifier 🕢				
	Select •					
Contact * 🕄	Name 📀	Affiliation 🕢				
	Trisovic, Ana	Harvard Univ	ersity		+	
	E-mail * 🕄				_	J
	anatrisovic@fas.harvard.edu					
Description * 🕢	This field supports only certain HTML tags.					
manage agence - company Theorem	Text * ()				+	
	Date 🕢					
	YYYY-MM-DD					
Subject * 🕢	Select				•	
Keyword 🕄	Term 🕢	Vocabulary 😔				
· · · · · · · · · · · · · · · · · · ·					+	
	Vocabulary URL 🕢					1
	Enter full URL, starting with http://					
Related Publication @	Citation ()					
					+	
					T	

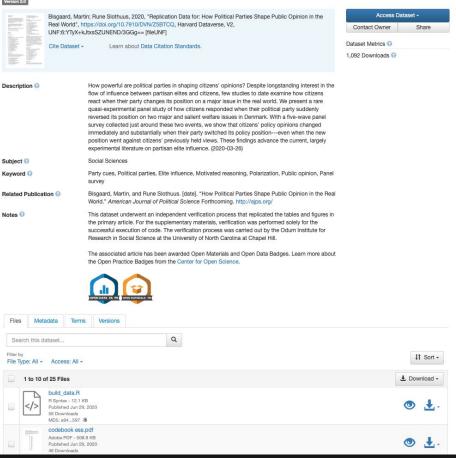
Replication dataset

 Replication dataset - a bundle of data, code and other files needed to reproduce a published study



Harvard Dataverse > American Journal of Political Science (AJPS) Dataverse >

Replication Data for: How Political Parties Shape Public Opinion in the Real World



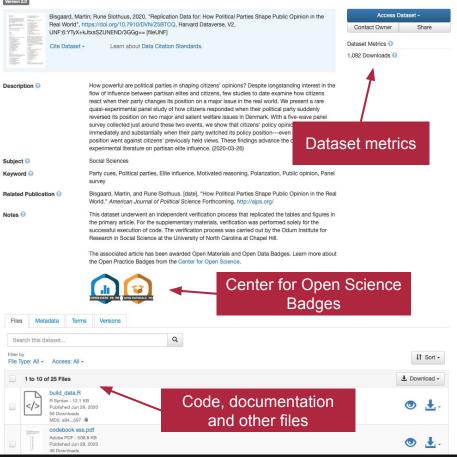
Replication dataset

 Replication dataset - a bundle of data, code and other files needed to reproduce a published study



Harvard Dataverse > American Journal of Political Science (AJPS) Dataverse >

Replication Data for: How Political Parties Shape Public Opinion in the Real World



Summary

- Dataverse data repositories have versatile support for data sharing
- Research data and code are shared in a "replication dataset" that often belong to a journal or institutional collection

How reusable are our replication datasets?



A large-scale study on research code quality and re-execution

A large-scale study on research code quality and

reusable replication datasets

code quality and re-execution

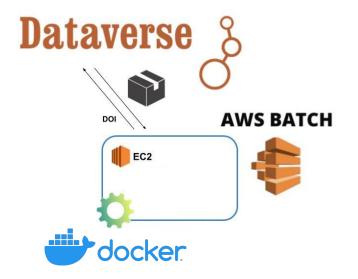
A large-scale study on research code quality and

re-execution

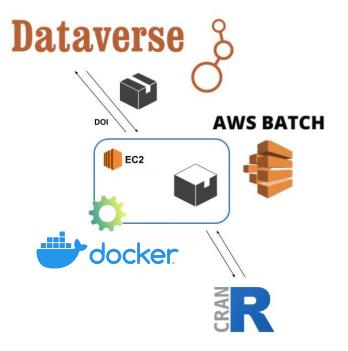
reusable replication datasets

not reusable replication datasets

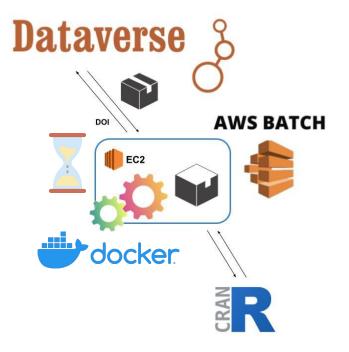
- 1. Replication dataset is retrieved from Harvard Dataverse to AWS
- 2. We collect data on the content, code, install used libraries etc.
- We attempt code re-execution for an allocated time of 1h per file and 5h in total
- 4. The re-execution result and other collected data are passed to the backend database for analysis



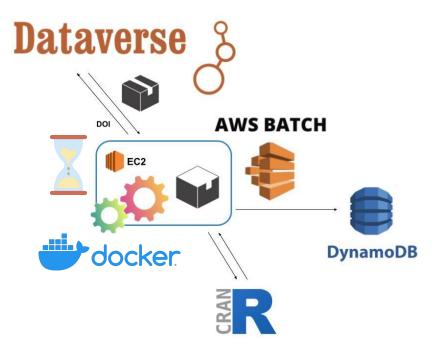
- 1. Replication dataset is retrieved from Harvard Dataverse to AWS
- 2. We collect data on the content, code, install used libraries etc.
- We attempt code re-execution for an allocated time of 1h per file and 5h in total
- 4. The re-execution result and other collected data are passed to the backend database for analysis



- 1. Replication dataset is retrieved from Harvard Dataverse to AWS
- 2. We collect data on the content, code, install used libraries etc.
- We attempt code re-execution for an allocated time of 1h per file and 5h in total
- 4. The re-execution result and other collected data are passed to the backend database for analysis



- 1. Replication dataset is retrieved from Harvard Dataverse to AWS
- 2. We collect data on the content, code, install used libraries etc.
- We attempt code re-execution for an allocated time of 1h per file and 5h in total
- 4. The re-execution result and other collected data are passed to the backend database for analysis



Results (basic properties)

Retrieved 2109 publicly available replication datasets containing 9078 R files

Over 94% of the datasets belonged to social sciences

Results (basic properties)

Retrieved 2109 publicly available replication datasets containing 9078 R files

Over 94% of the datasets belonged to social sciences

> Dataset size (median): 3.2 MB



Number of files (median): 8 (typically less than 15)

Results (basic properties)

Retrieved 2109 publicly available replication datasets containing 9078 R files

Over 94% of the datasets belonged to social sciences

> Dataset size (median): 3.2 MB

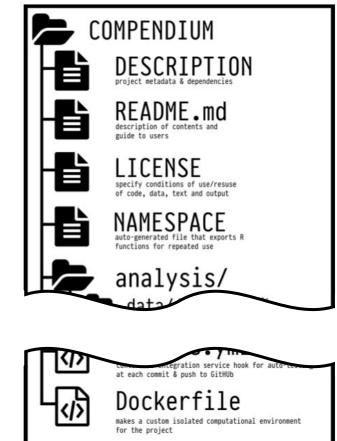
Number of files (median): 8 (typically less than 15)



- File name length: 10-20 characters
 - Documentation present in 57% of the datasets
- Comments comprise 20% of the shared R code

Presence of conventional files

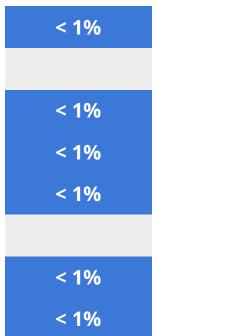
DESCRIPTION **README.md** LICENSE NAMESPACE Dockerfile **R** Markdown .Rproj install.R

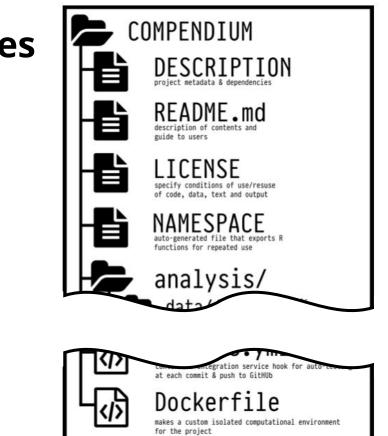


Marwick, Ben, Carl Boettiger, and Lincoln Mullen. "Packaging data analytical work reproducibly using R (and friends)." The American Statistician 72.1 (2018)

Presence of conventional files

DESCRIPTION **README.md** LICENSE NAMESPACE Dockerfile **R** Markdown .Rproj install.R

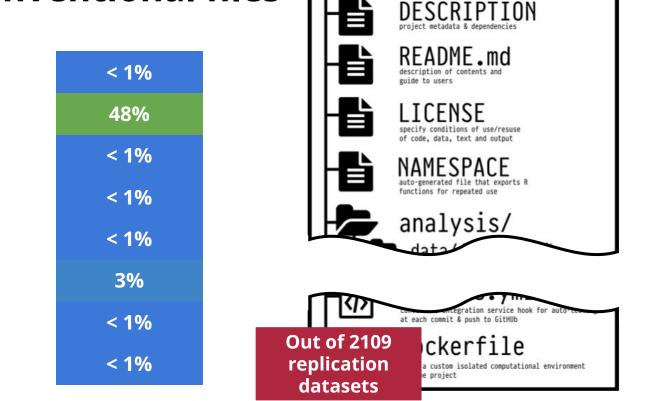




Marwick, Ben, Carl Boettiger, and Lincoln Mullen. "Packaging data analytical work reproducibly using R (and friends)." The American Statistician 72.1 (2018)

Presence of conventional files

DESCRIPTION **README.md** LICENSE NAMESPACE Dockerfile **R** Markdown .Rproj install.R



COMPENDIUM

Marwick, Ben, Carl Boettiger, and Lincoln Mullen. "Packaging data analytical work reproducibly using R (and friends)." The American Statistician 72.1 (2018)

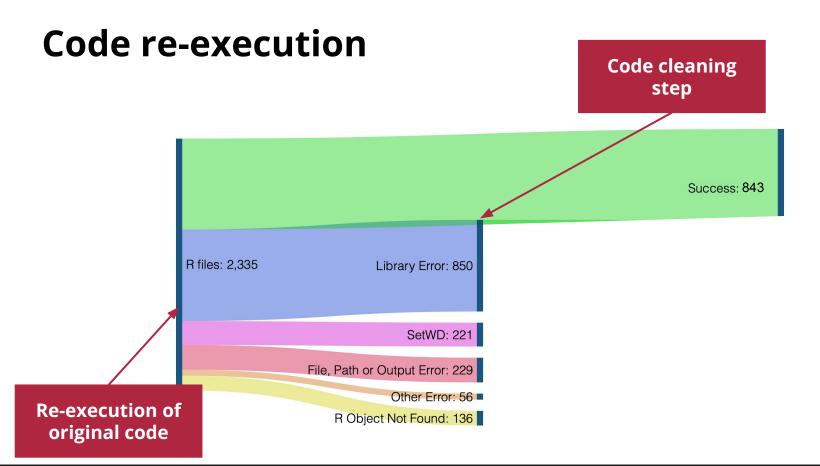
Most used libraries in research code:

- Data visualization and plotting (ggplot2, lattice)
- Data wrangling and display in a tabular form (xtable)
- 3. Data import and export (foreign, dplyr,
 - plyr, reshape2)
- 4. Statistical analysis (stargazer, MASS,

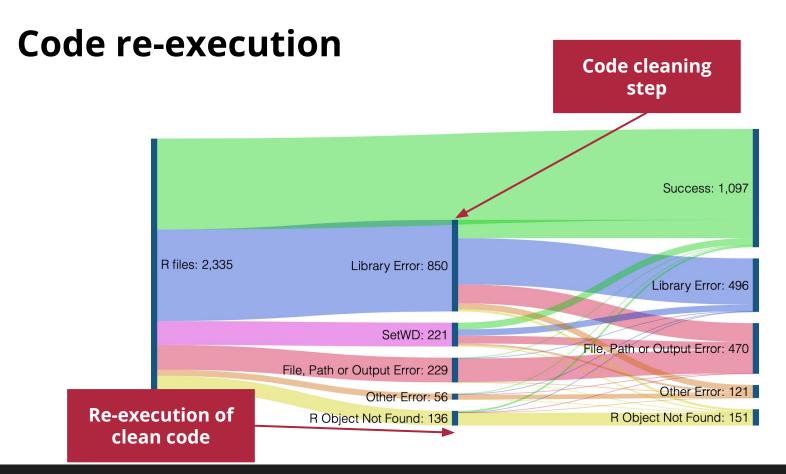
lmetest, car)

In the research code, no libraries detected for:

- Code testing (runit, testthat, tinytest, unitizer)
- 2. Provenance tracking (provR,
 - provenance, RDTlite, provTraceR)
- 3. Runtime environment management
 - (packrat, pacman)
- 4. Workflow libraries (workflowR,
 - workflows, drake)



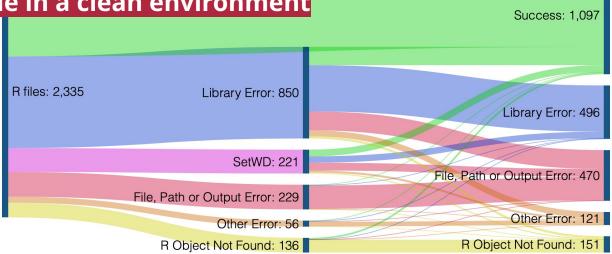
Trisovic, Ana, et al. "A large-scale study on research code quality and execution." arXiv preprint arXiv:2103.12793 (2021).



Trisovic, Ana, et al. "A large-scale study on research code quality and execution." arXiv preprint arXiv:2103.12793 (2021).

Code re-execution

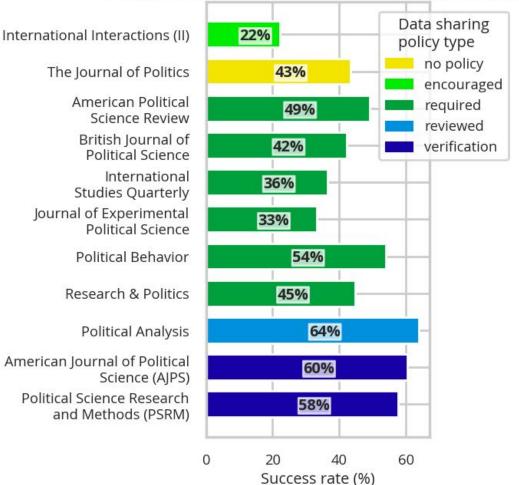
Many code errors can be avoided by capturing library dependencies and testing code in a clean environment



Trisovic, Ana, et al. "A large-scale study on research code quality and execution." arXiv preprint arXiv:2103.12793 (2021).

Portion of replication datasets with re-executable code files



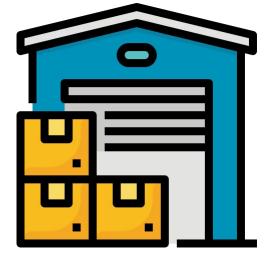


Journal average: 47% Total average: 45%

Summary

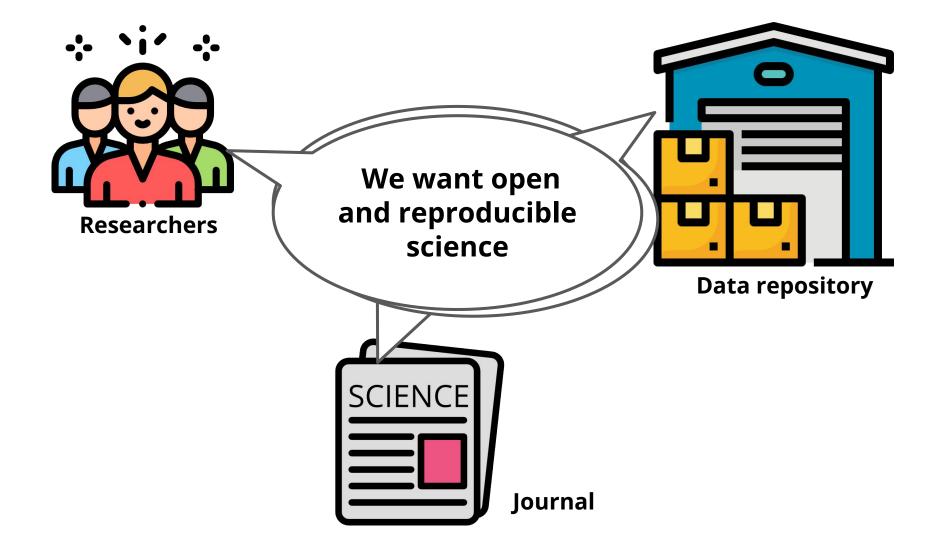
- No libraries for code testing, provenance or workflows
- Simple code cleaning resulted in substantial improvement in re-execution
- Re-execution correlates with journals' data sharing policy strictness





Data repository





What can researchers do?



Library versions should be captured by, minimally, the output of sessionInfo() from the researcher's R session, or a DESCRIPTION file, or install.R, or by using the renv package to track the libraries and their version number.
 When referring to data, code or other files, use relative file paths, as full paths will cause an error when the code is

executed on other systems.

- Library versions should be captured by, minimally, the output of sessionInfo() from the researcher's R session, or a DESCRIPTION file, or install.R, or by using the renv package to track the libraries and their version number.
- 2. When referring to data, code or other files, use relative file paths, as full paths will cause an error when the code is executed on other systems.

- 3. Workflow capture and managementment methods such as R Markdown, targets (or drake) will help to automate your code and specify the correct execution sequence.
- 4. Use Docker to document your runtime environment in a machine-readable format, and to ensure others can recreate your computing environment.

- 3. Workflow capture and managementment methods such as R Markdown, targets (or drake) will help to automate your code and specify the correct execution sequence.
- 4. Use Docker to document your runtime environment in a machine-readable format, and to ensure others can recreate your computing environment.



What can repositories do?



1. Create and maintain documentation on adequate deposit of research code.

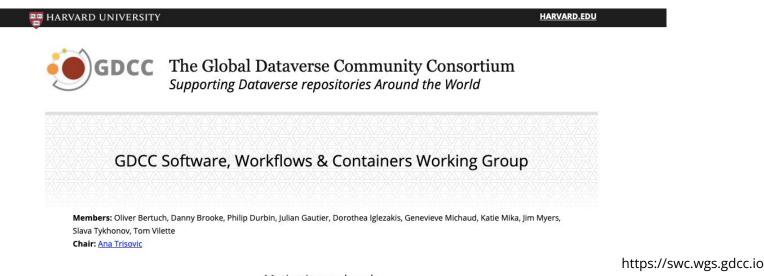
O Dataverse Project About + Community Best Practices + Software + Contact	
Search User Guide	Research Code Code files - such as Stata, R, MATLAB, or Python files or scripts - have become a frequent addition to the research data deposited in Dataverse repositories. Research code is typically developed by few researchers with the primary goal of obtaining results, while its
Account Creation + Mana	gement reproducibility and reuse aspects are sometimes overlooked. Because several independent studies reported issues trying to rerun research code, please consider the following guidelines if your dataset contains code.
Finding and Using Data Dataverse Collection Man	The following are general guidelines applicable to all programming languages.
Dataset + File Manageme	Create a README text file in the top-level directory to introduce your project. It should answer questions that reviewers or
Tabular Data File Ingest	README template for social science replication packages.
Data Exploration Guide	Depending on the number of files in your dataset, consider having data and code in distinct directories, each of which should
Appendix	 have some documentation like a README. Consider adding a license to your source code. You can do that by creating a LICENSE file in the dataset or by specifying the
Admin Guide	license(s) in the README or directly in the code. Find out more about code licenses at the Open Source Initiative webpage.
API Guide	If possible, use free and open-source file formats and software to make your research outputs more reusable and accessible.
Installation Guide	 Consider testing your code in a clean environment before sharing it, as it could help you identify missing files or other errors. For example, your code should use relative file paths instead of absolute (or full) file paths, as they can cause an execution error.
Developer Guide	Consider providing notes (in the README) on the expected code outputs or adding tests in the code, which would ensure that its
Style Guide	functionality is intact. guides.dataverse.or

2. Integrations with reproducibility platforms such as CodeOcean, WholeTale, Jupyter Binder and Renku will facilitate environment capture and encapsulation of research code.



Trisovic, Ana, et al. "Advancing computational reproducibility in the Dataverse data repository platform." P-RECS'20.

3. An internal working group will help identify community-wide problems, prioritize them, and implement solutions.



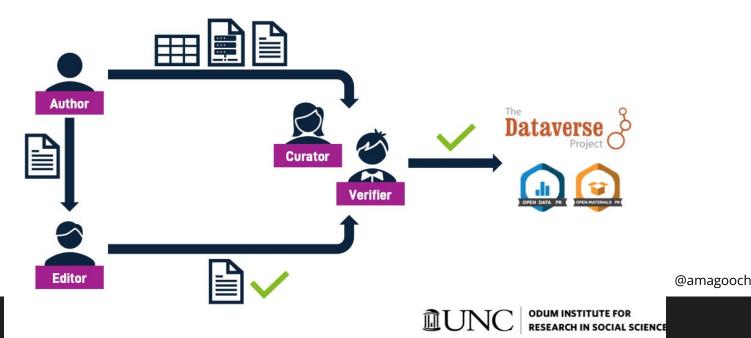
Motivation and goals

The SWC group aims to support research software, workflows, and container dissemination for reproducibility and reuse. The group discusses the necessary metadata, file formats, tools, and infrastructure necessary to incorporate these resources in data repositories

What can journals do?



1. Encourage a simple review of all deposited material if a code verification is infeasible



2. Create reproducibility checklist or templates for authors



A template README for social science replication packages.

The template README provided on this website is in a form that follows best practices as defined by a number of data editors at social science journals.

https://social-science-data-editors.github.io

Template README and Guidance

INSTRUCTIONS: This README suggests structure and content that have been approved by various journals, see Endorsers. It is available as Markdown/txt, Word, LaTeX, and PDF. In practice, there are many variations and complications, and authors should feel free to adapt to their needs. All instructions can (should) be removed from the final README (in Markdown, remove lines starting with > INSTRUCTIONS). Please ensure that a PDF is submitted in addition to the chosen native format.

Overview

INSTRUCTIONS: The typical README in social science journals serves the purpose of guiding a reader through the available material and a route to replicating the results in the research paper. Start by providing a brief overview of the available material and a brief guide as to how to proceed from beginning to end.

Example: The code in this replication package constructs the analysis file from the three data sources (Ruggles et al, 2018; Inglehart et al, 2019; BEA, 2016) using Stata and Julia. Two main files run all of the code to generate the data for the 15 figures and 3 tables in the paper. The replicator should expect the code to run for about 14 hours.

Data Availability and Provenance Statements

@larsvil

3. Integrations with reproducibility platforms

Home > arXiv updates > Instant access to code, for any arXiv paper



🍪 eLife

Instant access to code, for any arXiv paper

[IIII] Papers With Code

partners with



eLife and Stencila announce roadmap for bringing reproducible publishing to more authors

The next phase of the Executable Research Article project will focus on reducing barriers to the authoring and publication of reproducible research papers.

Conclusions

Conclusion

- We've seen evidence of both good and bad coding and dissemination practices (documentation, commenting, convention files rarely used)
- It is hard to re-execute "old" code and even harder to reuse it
 Curated replication datasets have higher re-execution rates.
 Things are looking up!
- Employing proposed recommendations would help researchers, repositories and journals contribute to research transparency and reproducibility.

Conclusion

- We've seen evidence of both good and bad coding and dissemination practices (documentation, commenting, convention files rarely used)
- It is hard to re-execute "old" code and even harder to reuse it
 - Curated replication datasets have higher re-execution rates.
 - Things are looking up!
- Employing proposed recommendations would help researchers, repositories and journals contribute to research transparency and reproducibility.

References

• This presentation was based on findings at Trisovic, Ana, et al. "A large-scale study on research code quality and execution." arXiv:2103.12793 (2021). A large-scale study on research code quality and execution

Ana Trisovic^{1*}, Matthew K. Lau², Thomas Pasquier³, Mercè Crosas¹

March 25, 2021

1. Institute for Quantitative Social Science, Harvard University, Cambridge, MA, USA

2. CAS Key Laboratory of Forest Ecology and Management, Institute of Applied Ecology, Chinese Academy of Sciences, Shenyang, China

3. Department of Computer Science, University of British Columbia, Vancouver, BC, Canada

*corresponding author(s): Ana Trisovic (anatrisovic@g.harvard.edu)

Abstract

This article presents a study on the quality and execution of research code from publicly-available replication datasets at the Harvard Dataverse repository. Research code is typically created by a group of scientists and published together with academic papers to facilitate research transparency and reproducibility. For this study, we define ten questions to address aspects impacting research reproducibility and reuse. First, we retrieve and analyze more than 2000 replication datasets with over 9000 unique R files published from 2010 to 2020. Second, we execute the code in a clean runtime environment to assess its ease of reuse. Common coding errors were identified, and some of them were solved with automatic code cleaning to aid code execution. We find that 74% of R files crashed in the initial execution, while 56% crashed when code cleaning was applied, showing that many errors can be prevented with good coding practices. We also analyze the replication datasets from journals' collections and discuss the impact of the journal policy strictness on the code re-execution rate. Finally, based on our results, we propose a set of recommendations for code dissemination aimed at researchers, journals, and repositories.

1 Introduction

Researchers increasingly publish their data and code to enable scientific transparency, reproducibility, reuse, or compliance with funding bodies, journals, and academic institutions [1]. Reusing data and code should propel new research and save researchers' time, but in practice, it is often easier to write new code



Thank you!

Email: anatrisovic@g.harvard.edu GitHub & Twitter: atrisovic Dataverse Project: https://dataverse.org/contact