

The Dataverse Project - Data sharing, Reproducibility, Research, Development and the Community

CNSTAT Expert meeting on Guidance on Data Sharing for NIA Longitudinal Studies May 10, 2021

Ana Trisovic, Harvard University on behalf of the Dataverse Project team



- 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)
- Developed at Harvard's Institute for Quantitative Social Science (IQSS) with contributions from the Dataverse community
 - 122 contributors to the software



- 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)
- Developed at Harvard's Institute for Quantitative Social Science (IQSS) with contributions from the Dataverse community
 - 122 contributors to the software

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



Metrics



https://dataverse.org/metrics (29 Dataverse installations)

Presentation agenda

- 1. Data sharing at Dataverse repositories
- 2. Documentation and guidelines
- 3. Reproducibility and reuse
- 4. Open questions, research and development

Data sharing

Data sharing

- Stand-alone or institutional account for depositing data
 - Dataverse installation account
 - GitHub, Google, ORCID, University credentials
- Two approaches for data sharing:
 - Via UI in the web browser
 - Via API (command line or Dataverse Software clients)

HARVARD Dataverse		Add Data 👻	Search +	About	User Guide	Suppor
Host Dataverse 💿	Changing the host dataverse will clear any fie	lds you may have	entered data	into.		
	Harvard Dataverse					
Asterisks indicate required fields						
Citation Metadata 🔨						
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 * 😧					
	anatrisovic@fas.harvard.edu					
Description * 🕢	This field supports only certain HTML tags.				+	
	Text * 😡					
	Date 🕢					
	YYYY-MM-DD					
Subject * 🕢	Select				•	
Keyword 😧	Term 🕄	Vocabulary 😔				
	Variabilitati URI O				+	
	Vocabulary URL 🕄					
Poloted Bublication	Citation ()					
Related Publication 🕢	Citation 9				+	
						ļ

Dataverse Software Clients

A client is an application that connects to and communicates with the remote Dataverse installation to transfer or manage data through the API.



All clients are open-source and available on GitHub.

Dataverse collection and datasets

- A dataverse collection is a collection of datasets and/or other collections
- Individuals, institutes or journals may have own dataverse collections





Container for your Datasets and/or Dataverses*

Dataverse collection and datasets

- A dataverse collection is a collection of datasets and/or other collections
- Individuals, institutes or journals may have own dataverse collections



Container for your Datasets and/or Dataverses*

		Henry A. Murray Research Archive at Harvard University
Murray Res	earch Archive Dataverse (Harvard Unive	arsity) Home
Harvard Datave	erse >	
The Henry A. collection cor	THE UNIVERSITY OF TEXAS Adoptio	on Project Dataverse (The University of Texas at Austin) Home Page
All "restricted	Harvard Dataverse >	
More informa		
<	Data from the Texas Adoption Proje	ct. Includes data from the main sample of 300 adoptive families and from an auxiliary sample
	Search this dataverse	Q Advanced Search
	Otataverses (0)	1 to 1 of 1 Result
	Datasets (1)	Texas Adoption Project
	Files (5)	Apr 14, 2010
	Publication Year	Joseph Horn; John Loehlin, 2010, "Texas Adoption Project", https://doi.org/ UNF:5:ov2C/MUSnhND+HKHhDinxQ== [fileUNF]
	2010 (1) Author Name John Loehlin (1)	This study presents data from two samples from the Texas Adoption Project. The main from a Texas home for unwed mothers between 1963 and 1971. Included are ability and children in the

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



Support for multiple metadata standards



Tabular files and variable metadata



Tabular files and variable metadata

Access File -				ndicators as Social Pressure	in Internatio	onal Relations
oad Options 🛓 nal File Format (Stata Binary) Delimited	Kelley, Judith;			ber: Indicators as Social Pressure in Internat	Varia	ble search
a Format	Q date	×	7 Results Download \$	Chart View	Table View	
ble Metadata	ID	Name	Label			
File Citation	1754727	crim1date		Variabl	e crimdate1:	
Explorer Data = ,	= 17547149	crimdate1		Values		Categories
	= 17547168	crimdate2		Summ	ary Statistics	N
	17547244	crimdate3				5977
		unctocdate		Maxin	num	2123 2012
	17547265			Minim	um	1967
	17547283	protocoldate				2006.4731470637444
						2007
	17547148	tier_date		155 variables		6.89046957527742

Tabular files and variable metadata

Replication Data for: Politics by Number: Indicators as Social Pressure in International Rel kelley_simmons_ajps_2014_replication.tab

Kelley, Judith; Simmons, Beth, 2018, "Replication Data for: Politics by Number: Indicators as Social Pressure in International Relations", https://doi.org/10.7910 UNF:6:IEYJGdv4VpPallzdj0qJNg== [fileUNF]





Extensive variable metadata (descriptive statistics) automatically derived from tabular data file in DDI format

Summary	Statistics
---------	------------

Cases	N
	5977
	2123
Maximum	2012
Minimum	1967
	2006.4731470637444
	2007
	6.89046957527742

Data sharing

Summary

- Dataverse data repository software
- Versatile support for data sharing via UI and the clients
- Extensive support in the repository for data, variables and metadata

Documentation and guidelines

Documentation and guidelines

The Dataverse Project team maintains an extensive set of guidelines for repository managers, developers and users.

Dataverse Community meetings

- projects.ig.harvard.edu/dcm2021
- lune 15-17, 2021



guides.dataverse.org

Monitoring

Documentation and guidelines

The Dataverse Project team maintains an extensive set of guidelines for repository managers, developers and users.

Dataverse Community meetings

- projects.iq.harvard.edu/dcm2021
- June 15-17, 2021



Monitoring

New guidelines

& Dataverse Project About -Community Best Practices -Software -Contact **Research Code** Search Code files - such as Stata, R. MATLAB, or Python files or scripts - have become a frequent addition Dataverse repositories. Research code is typically developed by few researchers with the primary User Guide reproducibility and reuse aspects are sometimes overlooked. Because several independent studie Account Creation + Management research code, please consider the following guidelines if your dataset contains code. Finding and Using Data The following are general guidelines applicable to all programming languages. Dataverse Collection Management Create a README text file in the top-level directory to introduce your project. It should ansi Dataset + File Management reusers would likely have, such as how to install and use your code. If in doubt, consider us 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 have some documentation like a README. Appendix Consider adding a license to your source code. You can do that by creating a LICENSE file Admin Guide license(s) in the README or directly in the code. Find out more about code licenses at the API Guide If possible, use free and open-source file formats and software to make your research outp Consider testing your code in a clean environment before sharing it, as it could help you ide Installation Guide example, your code should use relative file paths instead of absolute (or full) file paths, as ti **Developer Guide** Consider providing notes (in the README) on the expected code outputs or adding tests in functionality is intact. Style Guide Capturing code dependencies will help other researchers recreate the necessary runtime environr able to run correctly (or at all). One option is to use platforms such as Whole Tale, Jupyter Binder reproducibility. Have a look at Dataverse Integrations for more information. Another option is to us capture, which is often supported through the programming language. Here are a few examples: If you are using the conda package manager, you can export your environment with the cor *upcoming environment.yml. For more information, see the official documentation. Python has multiple conventions for capturing its dependencies, but probably the best-kno requirements.txt file, which is created using the command pip freeze > requirements.txt

with pip is explained in the official documentation.

guides.dataverse.org

Documentation and guidelines

- Each Dataverse installation may have local documentation and support
- At Harvard, user training is organized by the Harvard Library and IQSS Data Curation Services

HARVARD Dataverse Support	Getting Started V	Curation Services	Events	Policies & Governance 🔻	Go to Har
GETTING STARTED For Researchers For Journals For Organizations	Repository!	YOUR data is ea	container fo	the Harvard Data	g datasets. T
Contact Us	etc. We recommend collection. Files up to 2.5GB in a processed to access metadata search and	you first create a Data any format are accepted their metadata and thus	averse colle . Some form be able to p	er, research group, an entir ction and then add datasets ats (RData, SPSS, STATA, CS rovide more features for these you will need to:	to that Data

- Log in with the following options:
- · Your institutional affiliation login, e.g. HarvardKey
- Username/email
- ORCID
- GitHub
- Google

Documentation and guidelines

- Each Dataverse installation may have local documentation and support
- At Harvard, user training is organized by the Harvard Library and IQSS Data **Curation Services**
- Video guides on YouTube



HARVARD	Getting Started V	Curation Services	Events	Policies & Governance V	Go to Har	
	Getting Started Curation Services Events Policies & Governance Go to Har HOME / Getting Started Publishing your data is easy on the Harvard Dataverse Repository! A Dataverse collection is a customizable container for organizing and showcasing datasets. T Dataverse collection can be for an individual researcher, research group, an entire department etc. We recommend you first create a Dataverse collection and then add datasets to that Dataverse collection.					
	 metadata search and To deposit data in the Log in with the follow 	re-formatting). ne Harvard Dataverse	Repository y	rovide more features for these	type of files (

Documentation and guidelines

Summary

- Documentation and guidelines (global) are actively maintained
- Institutional guidelines, support and training

- Reproducibility: "obtaining consistent computational results using the same input data, steps, code, and conditions of analysis"
- Conducted large scale analysis to examine execution and quality of R and Python code
- We retrieved replication datasets from Harvard Dataverse repository and R code was re-executed in a pre-installed Docker container



National Academies of Sciences, Engineering, and Medicine. Reproducibility and replicability in science. National Academies Press, 2019. Trisovic, Ana, et al. "A large-scale study on research code quality and execution." arXiv preprint arXiv:2103.12793 (2021).

- Reproducibility: "obtaining consistent computational results using the same input data, steps, code, and conditions of analysis"
- Conducted large scale analysis to examine execution and quality of R and Python code
- We retrieved replication datasets from Harvard Dataverse repository and R code was re-executed in a pre-installed Docker container



National Academies of Sciences, Engineering, and Medicine. Reproducibility and replicability in science. National Academies Press, 2019. Trisovic, Ana, et al. "A large-scale study on research code quality and execution." arXiv preprint arXiv:2103.12793 (2021).



Trisovic, Ana, et al. "Repository Approaches to Improving Quality of Shared Data and Code." Data 6.2 (2021): 15.



Trisovic, Ana, et al. "Repository Approaches to Improving Quality of Shared Data and Code." Data 6.2 (2021): 15.



Trisovic, Ana, et al. "Repository Approaches to Improving Quality of Shared Data and Code." Data 6.2 (2021): 15.



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

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





Integration with reproducibility platforms

- New cloud tools have emerged to support collaborative work and research reproducibility by capturing code dependencies inside a web browser
- Dataverse software integration with reproducibility platforms allows:
 - Importing new research data with necessary libraries
 - Exporting and reusing the existing data





Solution of the second seco



CODE OCEAN Dashboard E	plore Learn		• Ø
ඹ 🕲 Combo Synergy of Two Drugs	Deta	ils Code Interface 🔛	😪 🐼 Run
Source Files <	main.R	> Resu	ults
	require(RSQLite) require(plotly)	Q Search	
🞧 main.R		> 요 Run 90665	0 O Q
run.sh	imake.plot<-function(drugA, drugB, cell){ guery<-paste('select drugA_name, drugA_conc, drugB_name, drugB_conc, x2x0, HSA, Bliss from		0 O O
	query<-paste('select drugA_mame, drugA_conc, drugB_mame, drugB_conc, x2x0, HSA, Bliss from drugA, '', '', drugB, '') and drugB_mame in ('', drugA, '', '', drugB, '') and cell_line cell, '' order by drugA_conc, drugB_conc', sege'')	* > _& Run 9655826	000
	data<-dbGetQuery(db,query)	> Run 9650940	9 0 Q
	if (dim(data)[1]>0) (> Run 3563692	90Q
	<pre>drug8_conc<-unique(data[,*drug8_conc*])</pre>	> 🖉 Run 3363841	000
	b.conc<-rep(c(0,1,2,3),4)	> Run 2826088	000
	<pre>observed.data<-matrix(data[,*x2x0*], nrow=4, ncol=4, byrow=TRUE)</pre>	> S Published Result	00
DrugCombodb IIS.76 MB	<pre>5 add_trace(x+b.conc, y=a.conc, z=da add_trace(x+b.conc, y=a.conc, z=da add_trace(x+b.conc, y=a.conc, z=da layout(tile = pate(drugA, "6 zaxis = list(tile = yaxis = list(tile = yaxis = list(tile = yaxis = list(tile = aspectratio=list(x=n, htmlwidgets::saveWidget(as.widget(else { write(* Sorry, This combination</pre> 5-FU & ABT-888 Combination	• 1	Observed Data Predicted (HSA) Predicted (Bliss)

In-brows of the Co platform

Figure credit: tech.cornell.edu/news
Integration with reproducibility platforms

- New cloud tools have emerged to support collaborative work and research reproducibility by capturing code dependencies inside a web browser
- Dataverse software integration with reproducibility platforms allows:
 - Importing new research data with necessary libraries
 - Exporting and reusing the existing data





Solution



Example: Jupyter Binder and Harvard Dataverse repository

HARVARD Add Data -About User Guide Support Ana Trisovic -Search -Dataverse **Replication Data for: Repository approaches to** improving quality of shared data and code Version 4.1 Sbinder Trisovic, Ana, 2020, "Replication Data for: Repository approaches to Access Dataset improving quality of shared data and code", https://doi.org/10.7910/DVN/E Edit Dataset -3LC5, Harvard Dataverse, V4 Link Dataset Cite Dataset -Learn about Data Citation Standards. Contact Owner Share Dataset Metrics 22 Downloads 🕣 Description 🕤 This is supplementary, ata to the article "Repository approaches to improving quality of shared data and code," and in particular, its first section on completeness of research code. Run this code on Julyter Binder here: 🚱 launch binder (2020-09-27) Subject 🕄 Computer and Information Science Files Metadata Terms Versions Search this dataset... Q Find + Upload Files

Example: Jupyter Binder and Harvard Dataverse repository



Turn a Git repo into a collection of interactive

notebooks

Have a repository full of Jupyter notebooks? With Binder, open those notebooks in an executable environment, making your code immediately reproducible by anyone, anywhere.

 HARVARD
Dataverse

Version 4.1

Search this dataset...

Add Data - Search - About User Guide Suppo

Replication Data for: Repository approaches to improving quality of shared data and code

Q Find

Trisovic, Ana, 2020, "Replication Data for: Repository approaches to Access Da New to Binder? Get started with a Zero-to-Binder tutorial in Julia. Python or R. improving quality of shared data and code", https://doi.org/10.7910/DVN/EA Edit Data **3LC5**, Harvard Dataverse, V4 Link Date Cite Dataset -Learn about Data Citation Standards Build and launch a repository Contact Owner Dataverse DOI (10.7910/DVN/TJCLKP) Dataset Metrics Dataverse DOI - Dataverse DOI (10.7910/DVI/TJCLKP) 22 Downloads 🕣 Description 🕤 This is supplementary data to the article "Repository approaches to improving guality of shared data and code," and in particular, its first Gi ref (branch, tag, or commit) Path to a notebook file (optional) section on completeness of research code. Run this code on Jupyter Binder here: <a>[g] launch binder (2020-09-27) AD Path to a notebook file (optional) File -Subject 🕄 Computer and Information Science mybinder.org Files Metadata Versions Terms

Upload Files

Build and launch a repository

Dataverse DOI (10 7910/DVN/TICLKP)

Example: Jupyter Binder and ranch, tag, or commit) **Harvard Dataverse** repository Copy the URL below and share your Binder with others: https://mybinder.org/v2/dataverse/10.7910/DVN/EA3LC5/

Expand to see the text below, paste it into your README to show a binder badge: S launch binder

Jupyter Binder creates runtime environment for research code



Path to a notebook file (optional)

Path to a notebook file (optional)

File -

自

8 - 8 < > 🗉 0 C đ A hub-binder.mybinder.ovh 0 ð Start Page **Example: Jupyter Binder and** Visit repo Copy Binder link Not Trusted Python 3 O Harvard Dataverse CH + = La Download 🔕 🙆 🗹 Dataverse % Binder Memory: 203.7 MB / 2 GB Code repository Analysis of Python results In [1]: import numpy as np import matplotlib.pyplot as plt import seaborn as sns from matplotlib import rc import pandas as pd # plot styles sns.set style('whitegrid') sns.set_style({'font.family': 'Times New Roman'}) In [2]: df = pd.read_csv("python-study-data.csv", index_col=0) In [3]: df.head() Out [3]: doi result3 filename result2 list of all size TypeError TypeError: coercina to invalid file: simulation_output.txt;format_cdc_data.sh;evalu... 27499 0 doi:10.7910/DVN/8TB7GO el preprocessing.py Unicode: need None string or .. SyntaxError: TypeError: Anyone can explore unsupported Missing simulation output.txt:format cdc data.sh:evalu... 274991 doi:10.7910/DVN/8TB7GO ei preprocessing india.pv operand type(s) parentheses in call to 'p ... for +: data, make direct AttributeError SyntaxError: 'NoneType Missing 2 doi:10.7910/DVN/8TB7GO ei_preprocessing_ipums_census_acs_samples.py simulation_output.txt;format_cdc_data.sh;evalu... 274991 changes in the code object has no parentheses in attri ... call to 'p... or data AttributeError SyntaxError 'NoneType Missing 3 doi:10.7910/DVN/8TB7GO ei_preprocessing_ipums_full_census.py simulation_output.txt;format_cdc_data.sh;evalu... 274991 object has no parentheses in attri. call to 'p ... TypeError TypeError: coercina to doi:10.7910/DVN/8TB7GO invalid file: simulation output.txt;format cdc data.sh;evalu... 274991 ei preprocessing race.pv Unicode: need None string or .. In [91]: df[df.result2.isnull()] Out[91]: doi result3 filename result2 list of all size

Reproducibility versus reuse



Reproducibility versus reuse



Reproducibility versus reuse



Harvard Dataverse > Murray Research Archive Dataverse > Early Head Start Research and Evaluation Project Dataverse >

Research and prototype projects

Research reproducibility test at upload

Ana Trisovic

- 1 Project reference: sJV8vKMPYjutkCOBQmWR
- 2 downloading file: https://firebasestorage.googleapis.com/v0/b/re3-fb.appspot.com/o/rep
- 3 downloading file: https://firebasestorage.googleapis.com/v0/b/re3-fb.appspot.com/o/rep
- 4 downloading file: https://firebasestorage.googleapis.com/v0/b/re3-fb.appspot.com/o/re
- 5 downloading file: https://firebasestorage.googleapis.com/v0/b/re3-fb.appspot.com/o/rep
- 6 downloading file: https://firebasestorage.googleapis.com/v0/b/re3-fb.appspot.com/o/rep
- 7 downloading file: https://firebasestorage.googleapis.com/v0/b/re3-fb.appspot.com/o/rep
- 8 downloading file: https://firebasestorage.googleapis.com/vo/u/.2 fb.appspot.com/o/rep
- 9 FILENAME. ButlerHomola_Excludability_Appendix.R
- 10 f_SULT: ['Error in library("stargazer") : there is no package called 'stargaze.']
 11 FILENAME: ButlerHomola_Excludability_Analysis.R
- 12 RESULT: ['Error in library("stargazer") : there is no package called 'stargazer'

Research and prototype projects

Research reproducibility test at upload

Ana Trisovic

1 Project reference: sJV8vKMPYjutkCOBQmWR 2 downloading file: https://firebasestorage.googleapis.com/v0/b/re3-fb.appspot.com/o/rep 3 downloading file: https://firebasestorage.googleapis.com/v0/b/re3-fb.appspot.com/o/rep 4 downloading file: https://firebasestorage.googleapis.com/v0/b/re3-fb.appspot.com/o/rep 5 downloading file: https://firebasestorage.googleapis.com/v0/b/re3-fb.appspot.com/o/rep 6 downloading file: https://firebasestorage.googleapis.com/v0/b/re3-fb.appspot.com/o/rep 7 downloading file: https://firebasestorage.googleapis.com/v0/b/re3-fb.appspot.com/o/rep 8 downloading file: https://firebasestorage.googleapis.com/v0/b/re3-fb.appspot.com/o/rep 8 downloading file: https://firebasestorage.googleapis.com/v0/b/re3-fb.appspot.com/o/rep 9 fILENA". ButlerHomola_Excludability_Appendix.R 10 fSULT: ['Error in library("stargazer") : there is no package called 'stargaze.'] 13 FILENAME: ButlerHomola_Excludability_Analysis.R 12 RESULT: ['Error in library("stargazer") : there is no package called 'stargazer'] 13

Code readability assessment



Bahaidarah, Layan, et al. "Toward Reusable Science with Readable Code and Reproducibility." (to appear)

Reproducibility and reuse

Summary

- A large-scale study sheds light on reproducibility challenges
- Integrations with cloud tools facilitate software capture
- Further development focuses on reuse

Open questions, research and development

Open questions, research and development

- DDI Cross Domain Integration (DDI-CDI)
 - Planning to implement cross-domain metadata to facilitate merging data files
 - Controlled vocabularies for each variable
- Will incorporate software metadata
 - Enabling licencing, attribution and software dependency
- Global Dataverse Community Consortium (GDCC) facilitates Dataverse community efforts, and supports Dataverse repositories around the world

 GDCC working groups are formed to tackle the open questions

Open questions, research and development

- DDI Cross Domain Integration (DDI-CDI)
 - Planning to implement cross-domain metadata to facilitate merging data files
 - Controlled vocabularies for each variable
- Will incorporate software metadata
 - Enabling licencing, attribution and software dependency
- Global Dataverse Community Consortium (GDCC) facilitates Dataverse community efforts, and supports Dataverse repositories around the world
 GDCC working groups are formed to tackle the open questions

Software, workflows and containers working group

- SWC (<u>swc.gdcc.io</u>) identifies requirements and undertakes developments to support research software and reproducibility at the Dataverse Project
 - We welcome new members!
 - One solution for containers:



Conclusion

- The Dataverse Project creates **research data repository** software
- Ongoing support and ongoing enhancements for descriptive metadata for datasets and variable metadata for tabular files
- Research is conducted to:
 - examine common practices,
 - identify shortcomings of existing approaches and
 - new software developments
- Further development include supporting dissemination of emerging computational components, such as **workflows** and **containers**, with specific metadata, tools and infrastructure

Conclusion

- The Dataverse Project creates **research data repository** software
- Ongoing support and ongoing enhancements for descriptive metadata for datasets and variable metadata for tabular files
- Research is conducted to:
 - examine common practices,
 - identify shortcomings of existing approaches and
 - new software developments
- Further development include supporting dissemination of emerging computational components, such as **workflows** and **containers**, with specific metadata, tools and infrastructure



Thank you!

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