MAKING SOCIAL SCIENCE MORE REPRODUCIBLE

(Starting with Data Curation and Code Review)

Limor Peer Associate Director for Research, Institution for Social and Policy Studies Data and Research Specialist, Office of the Provost Yale University

GESIS Lecture Series

March 13, 2018



INSTITUTION FOR SOCIAL AND POLICY STUDIES

ISPS was founded in 1968 as an interdisciplinary center to support social science and public policy research at Yale University



Field Experiments Initiative





Biodiversity Field Experiment (University of Minnesota's Cedar Creek Ecosystem Science Reserve, with permission)

In 2001, the Institution for Social and Policy Studies announced an initiative designed to encourage field experimentation in the social sciences at Yale.

The term 'field experiment' refers to fully-randomized research designs in which observations found in a naturalistic setting – voters, patients, welfare recipients, community organizations, government entities, and the like – are assigned to treatment and control conditions. Recent examples of this kind of research at ISPS include randomized studies of voter mobilization, peer counseling of homeless people, campaign activities in Africa, and the persuasiveness of televised campaign advertisements.







ISPS DATA ARCHIVE

Since 2011 Specialized community Open access Website integration

90 studies 1,400 files 15 GB

Author	Discipline		Keywords			Area of Study		
- Any - 🗸	- Any -	~	- Any -	~		- Any -	~	
Location	Year		Research	design	• • •			
- Any -	-Year	~	- Any -	~		SEAR		
TITLE				AUTHOR(S)			YEAR ARCHIVED	~
Did Shy Trump Supporters Bias the 2 representative List Experiment	016 Polls? Evidence from a Nat	tionally-		Alexander Coppock			2017	
Trading Barriers: Firms, Immigration	and the Remaking of Globaliza	ition		Margaret E. Peters			2017	
The Majority-Minority Divide in Atti from Mumbai	tudes Toward Internal Migrati	on: Evidend	<u>ce</u>	Nikhar Gaikwad and (Gareth N	Vellis	2017	
Chocolate Scents and Product Sales: Bookstore and Café	A Randomized Controlled Tria	al in a Cana	<u>dian</u>	Mary C. McGrath, Pet Shotwell	ter M. A	ronow, Vivien	2017	

OUTLINE

Making social science more reproducible: Why, What?

Defining terms

Curating for reproducibility: What? How?

Curation Tool: Yale Application for Research Data (YARD)

OUTLINE

Making social science more reproducible: Why, What?

Defining terms

Curating for reproducibility: What? How?

Curation Tool: Yale Application for Research Data (YARD)

E.G., REINHART & ROGOFF

Free exchange The 90% question



A seminal analysis of the relationship between debt and growth comes under attack

Print edition | Finance and economics > Apr 20th 2013



GOVERNMENT indebtedness matters. Default and financial panic are the stuff of finance-minister nightmares. Government borrowing can crowd out private investment, dragging growth down. Yet economists have struggled to specify when a country needs to worry about its debt load. In a 2010 paper Carmen Reinhart, now a professor at Harvard Kennedy School, and Kenneth Rogoff, an economist at Harvard University, seemed to provide an answer. They argued that GDP growth slows to a snail's pace once government-debt levels exceed 90% of GDP.

https://www.economist.com/news/financeand-economics/21576362-seminal-analysisrelationship-between-debt-and-growth-comesunder

E.G., REINHART & ROGOFF

In their Excel spreadsheet, Reinhart and Rogoff had not selected the entire row when averaging growth figures: they omitted data from Australia, Austria, Belgium, Canada and Denmark...

University of Massachusetts Amherst researchers uncovered the mistake after they obtained the actual spreadsheet that Reinhart and Rogoff used for their calculations.

[The authors] have acknowledged making a spreadsheet calculation mistake in a 2010 research paper... But the authors stand by their conclusion that higher government debt is associated with slower economic growth.

E.G., PIKETTY

The data underpinning Professor Piketty's 577-page tome, which has dominated best-seller lists in recent weeks, contain a series of errors that skew his findings, according to a Financial Times investigation.

Piketty findings undercut by errors



Economist and author Thomas Piketty

Chris Giles in London MAY 23, 2014

257

Thomas Piketty's book, 'Capital in the Twenty-First Century', has been the publishing sensation of the year. Its thesis of rising inequality tapped into the zeitgeist and electrified the post-financial crisis public policy debate.

E.G., PIKETTY - CONTINUED

"On the plus side, he published his code ... on the negative side, it appears that Piketty's code contains mistakes, fudging and other problems. ... Simply put, spreadsheets are good for quick and dirty work, but they are not designed for serious and reliable work. ... Spreadsheets make code review difficult. The code is hidden away in dozens if not hundreds of little cells. If you are not reviewing your code carefully and if you make it difficult for others to review it, how do expect it to be reliable?" -- Daniel Lemire

http://lemire.me/blog/archives/2014/05/23/you-shouldnt-use-a-spreadsheet-for-important-work-i-mean-it/

"it's a red herring in the Piketty discussion, except insofar as both [Piketty and Reinhart & Rogoff] are examples that help flesh out **standards and guidelines for data/code release** in economics." -- Victoria Stodden

http://blog.stodden.net/2014/05/27/mistakes-piketty/

SCIENCE ADVANCES WHEN SCIENTIFIC CLAIMS ARE SUBJECTED TO SCRUTINY

One of the core principles of the scientific process is that other scientists are able to repeat your experiment and either confirm or refute your results.

This is referred to as reproducibility or replication.



Repeat After Me, Maki Naro https://thenib.com/repeat-after-me

1.5

Individual Studies

OUTLINE

Making social science more reproducible: Why, What?

Defining terms

Curating for reproducibility: What? How?

Curation Tool: Yale Application for Research Data (YARD)

REPRODUCIBILITY DEFINITION

computational reproducibility empirical reproducibility VALIDATION

REPEATABILITY REPLICABILITY

methodological reproducibility

VERIFICATION **REPRODUCIBILITY**

statistical reproducibility

REPRODUCIBILITY DEFINITION

Reproducibility: calculation of quantitative scientific results by independent scientists using the original datasets and methods

Stodden, V. (Ed.), Leisch, F. (Ed.), Peng, R.D. (Ed.). (2014). *Implementing Reproducible Research*. New York: Chapman and Hall/CRC.

Note:

- 1. cf. **Replication**: implementing experiments and collecting new data for analysis by other researchers to validate the science
- 2. Different disciplines may interpret these terms differently

SHARING DATA IS ESSENTIAL

- ★ To reproduce or to verify research
- **★** To make the results of publicly funded research available to the public
- ★ To enable others to ask new questions of extant data
- ★ To advance the state of research and innovation

Borgman, C. L. (2012). The conundrum of sharing research data. *Journal of the American Society for Information Science and Technology, 63*(6), 1059-1078. http://doi.org/10.1002/asi.22634

SHARING DATA IS REQUIRED







laura and john arnold foundation™







LIMOR PEER, YALE UNIVERSITY

BUT...

SHARING DATA IS NOT SUFFICIENT

The most commonly reported problems associated with [replication] attempts were the lack of... data and code, followed by insufficient documentation.

Janz, N., Werfel, S., Wykstra S. (2014). Replication in political science graduate courses: an untapped resource? *Monkey Cage.* https://www.washingtonpost.com/news/monkey-cage





Christie Bahlai @cbahlai 21h Aaaannnd: unexplained missing data. Must consult with data creator before proceeding. Will write stats code while I wait. #otherpeoplesdata

Details

SHARING DATA WITH CARE

Because there are more ways to share data, and because the scholarly landscape supports and encourages that, there is a proliferation of data files on many different types of systems that **do not meet the criterion of quality**...

Peer, L., Green, A., & Stephenson, E. (2014). Committing to data quality review. *International Journal of Digital Curation*, 9(1). http://doi.org/10.2218/ijdc.v9i1.317

REPRODUCIBILITY STANDARD (1)

Reproducibility refers to the ability of a researcher to duplicate the results of a prior study using the same materials as were used by the original investigator. That is, a second researcher might use the same raw data to build the same analysis files and implement the same statistical analysis in an attempt to yield the same results. Reproducibility is a minimum necessary condition for a finding to be believable and informative.

K. Bollen, J. T. Cacioppo, R. Kaplan, J. Krosnick, J. L. Olds (2015). *Social, Behavioral, and Economic Sciences Perspectives on Robust and Reliable Science,* National Science Foundation.

REPRODUCIBILITY STANDARD (2)

If a paper makes, or implies, claims that require software, those claims should be backed up.... We are after repeatability, which is simply **the act of checking the claims made in the paper**... Repeatability should become a standard feature of the dissemination of research results.

Krishnamurthi, S., Vitek, J. (2015). *Communications of the ACM*, Vol. 58 No. 3, Pages 34-36, https://doi.org/10.1145/2658987

REPRODUCIBILITY STANDARD (3)

The standard of reproducibility calls for the **data and the computer code** used to analyze the data be made available to others. This standard falls short of full replication because the same data are analyzed again, rather than analyzing independently collected data... However, under this standard, limited exploration of the data and the analysis code is possible and **may be sufficient to verify the quality of the scientific claims**.

Peng, Roger D. (2011). Reproducible research in computational science. Science. Vol. 334, Issue 6060, pp. 1226-1227 https://doi.org/10.1126/science.1213847

REPRODUCIBILITY STANDARD (4)

The replication standard holds that sufficient information exists with which to understand, evaluate, and build upon a prior work if a third party could replicate the results without any additional information from the author.

King, G. (1995). Replication, replication. *PS: Political Science & Politics, 28*(3), 444–452. http://doi.org/10.2307/420301

REPRODUCIBILITY STANDARD (5)

Could the published computational findings be **reproduced on an independent system by using the data and code** provided?

Stodden, V., McNutt, M., Bailey, D. H., Deelman, E., Gil, Y., Hanson, B., ... Taufer, M. (2016). Enhancing reproducibility for computational methods. *Science*, *354*(6317), 1240–1241. https://doi.org/10.1126/science.aah6168

REPRODUCIBILITY STANDARD DEMANDS DATA QUALITY

A set of measures that determine if data are **independently understandable for informed reuse**.

Peer, L., Green, A., & Stephenson, E. (2014). Committing to data quality review. *International Journal of Digital Curation, 9*(1). http://doi.org/10.2218/ijdc.v9i1.317

OUTLINE

Making social science more reproducible: Why, What?

Defining terms

Curating for reproducibility: What? How?

Curation Tool: Yale Application for Research Data (YARD)

CURATION DEFINITION

The active and ongoing management of data through its life cycle of interest and usefulness to scholarship, science, and education. Data curation activities enable data discovery and retrieval, maintain its quality, add value, and provide for reuse over time, and this new field includes authentication, archiving, management, preservation, retrieval, and representation.

The University of Illinois' Graduate School of Library and Information Science, https://ischool.illinois.edu/articles/2012/04/palmer-discusses-big-data-u-i-news-bureau







Roger Peng (Johns Hopkins University)

SOFTWARE & CODE

"Code and software are now an integral part of research data generation, usage and analysis"

"Deep intellectual contributions now encoded in Software" "Code as first-class research object"

Nature editorial, 2015 http://www.nature.com/articles/sdata20 154

Stodden, 2014

https://web.stanford.edu/~vcs/talks /OpenConNov152014-STODDEN.pdf Mozilla Science Lab in collaboration with figshare and github https://science.mozilla.org/blog/code-as-

a-research-object-a-new-project

Science Code Manifesto

Manifesto Discussion Endorse Resources About

Software is a cornerstone of science. Without software, twenty-first century science would be impossible. Without better software, science cannot progress.

But the culture and institutions of science have not yet adjusted to this reality. We need to reform them to address this challenge, by adopting these five principles:

Code	All source code written specifically to process data for a published paper must be available to the reviewers and readers of the paper.
Copyright	The copyright ownership and license of any released source code must be clearly stated.
Citation	Researchers who use or adapt science source code in their research must credit the code's creators in resulting publications.
Credit	Software contributions must be included in systems of scientific assessment, credit, and recognition.
Curation	Source code must remain available, linked to related materials, for the useful lifetime of the publication.

CODE SHARING REQUIREMENTS

Journals: <u>Stodden et al</u> (2013) find 30% increase in journals requiring code.

- Do you want to publish here?
 - Nature group
 - Science
 - PLoS
- AJPS (see APSA DART)
- •And the list is growing...



- Assign persistent identifier
- Create study citation and study-level metadata record
- Record file size details
- Check for presence of all files
- Verify content of files matches expected format
- Create non-proprietary versions of files
- Implement migration strategy for file formats



- Confirm presence of comprehensive descriptive information necessary for informed reuse
 - Data definitions
 - Variable construction
 - Methodology
 - Sampling information
 - Original data source citation
 - Analysis software version
- ✓ Link to related research products



- Check for undocumented variable and value information
- Examine data for inconsistencies and errors
 - Discrepancies in number of observations
 - Out-of-range or wild codes
 - Undefined null values
- Review data for confidentiality issues



- Convert absolute file paths to relative file paths
- Check code for presence of non-executable comments that document analysis processes
- Identify packages required to execute code
- Execute code to ensure code is error-free
- Compare code output to findings presented in article

CURATING FOR REPRODUCIBILITY THE CURE CONSORTIUM

Supporting research data curation and code review for the purpose of facilitating the digital preservation of the evidence base necessary for future understanding, evaluation, and replication of scientific claims.

https://cure.web.unc.edu/



Establish Standards

Share Practices

Promote Data Quality Review



CURATING FOR REPRODUCIBILITY THE CURE CONSORTIUM





CU RATING for RE PRODUCIBILITY

https://cure.web.unc.edu/

CURATING FOR REPRODUCIBILITY MODELS OF PRACTICE



- Institution for Social and Policy Studies (ISPS)
 Aligning Data Curation Workflows with Data Quality Review
- **2. Cornell Institute for Social and Economic Research** Providing Data Curation and Reproduction of Results (R²) Services
- **3. Odum Institute for Research in Social Science** Enforcing Journal Data Replication Policies

OUTLINE

Making social science more reproducible: Why, What?

Defining terms

Curating for reproducibility: How?

Curation Tool: Yale Application for Research Data (YARD)



CURATING FOR REPRODUCIBILITY

Institution for Social and Policy Studies (ISPS)

Aligning Data Curation Workflows with Data Quality Review

ISPS data archivists developed a data curation workflow that implements the ideals of scientific reproducibility and transparency



MISSING LABELS

"We are missing labels for the following variables: _n1, _n0, V1 and V0."

Archive staff

"Here are the labels: _n1 is the number of observations in the treated strata before matching _n0 is the number of observations in the comparison strata before matching v1 = turnout for treated observations <u>v0 = turnout</u> for comparison observations

... this reminds me that I needed to include the .ado code in the Matching Code folder. I just did that and updated the readme file. Boy, the things your forget about after not thinking about something for two years!"

Researcher

CODE BEFORE

🛃 Do-f	file Editor - CogBiasRept	
File	Edit View Project Tools	
D P2		
		- *
	jbiaskepi Untitled.do	•
1	set more off	
3		
4	log using "CogBiasRepl.log", replace	
5		
7	constudy I con	
8	use "studyldata.dta", clear	
9		
10	**Randomization Check	
12	x1: miggit treatment ideology partisanship age education income remain black white	
13	**Measure Descriptives	
14	sum VaccineOpinion	
15	sum fear_flu	
16	**Table 1	
18	*Column 1	
19	reg VaccineOpinion bc bic comp_combine	E
20	lincom bc-bic	
21	*Column 2	
22	reg vaccineOpinion bc blc comp_combine drear_iu ixbc_i ixblc_i ixcbc_i	
24	*Column 3	
25	reg VaccineOpinion bc bic comp_combine dangry_flu anxbc_f anxbic_f anxcbc_f	
26	*Column 4	
27	reg VaccineOpinion bc bic comp_combine drear_flu IXbc_f IXblc_f Ixcbc_f dangry_flu anxbc_f anxbic_f anxbic_f	
29	**Mediation Analysis	
30	sgmediation VaccineOpinion if comp_combine==1, mv(NewVacPers) iv(dfear_flu) cv(StockVacPers)	
31	sgmediation VaccineOpinion if comp_combine==1, mv(StockVacPers) iv(dfear_flu) cv(NewVacPers)	
32		
34	***Study 2***	
35		
36	use "study2data.dta", clear	
37	**Pandomization Check	
39	mlogit arg treatment education white rep dem age income female feardisp	
40		
41	**Measure Descriptives	
42	alpha nervous frightened afraid jittery	
44	sum anxiety	
45	sum BanGaySupport	
46		
47	**Manipulation Check	
49	cost analoy, by (learphotos)	*
	"	E.
Ready		Line: 34, Col: 13 CAP NUM OVR

LIMOR PEER, YALE UNIVERSITY

44

IS THE CODE WELL-ANNOTATED?

🕑 Do-fi	ile Editor - Antonio Contractor a bles_do	
<u>F</u> ile	Edit View Project Tools	
DR		
	Renarcos Rosanz 2 lables_do Untitled.do	• ×
1	****STUDY 1****	<u>^</u>
2	Attribute to researchers. Insert fileseth to access data in W2	
4	*****For mediation analysis, download the somediation command using "findit somediation."******	
5		
6	**Randomization Check	
7	xi: mlogit treatment ideology partisanship age education income female black white	
8		
9	**Measure Descriptives	
110	sum vaccineopinion	
12		
13	**Table 1	
14	*Column 1	
15	reg VaccineOpinion bc bic comp_combine	
16	lincom bc-bic	
17	*Column 2	E
18	reg vaccineupinion bc bic comp combine drear ilu ixbc i ixbic i ixcbc i	
20	Column 3	
21	reg VaccineOpinion bc bic comp combine dangry flu anxbc f anxbic f anxcbc f	
22	*Column 4	
23	reg VaccineOpinion bc bic comp_combine dfear_flu fxbc_f fxbic_f fxcbc_f dangry_flu anxbc_f anxcbc_f	
24		
25	**Mediation Analysis	
20	symediation VaccineOpinion if comp_combine==1, mV(NewVacPers) IV(dreat_IIU) cV(StocKVaCPers)	
28	Symediation vaccineopinion in comp_combinei, mv(Stockvacreis) iv(diear_iid) cv(Newvacreis)	
29		
30	****STUDY 2****	
31	responses for the second	
32	****Note to researchers: insert filepath to access data in "Arceneaux_AJFS_2012_study_2_dta"****	
33	**Dendering for Charle	
35	"Annual and treatment education white ren dem age income female feardign	
36	moget day_ordenants candedten where rep dam age income remark restarby	
37	**Measure Descriptives	
38	alpha nervous frightened afraid jittery	
39	egen anxiety=rowmean(nervous frightened afraid jittery)	
40	sum anxiety	
41	sum sancaysupport	
43	**Manipulation Check	
44	ttest anxiety, by(fearphotos)	
45		+
•	III .	F 1
Ready		Line 12 Col: 0 CAP NUM OVR

DOES THE CODE FULLY EXECUTE?

Example 1:	File E	Edit D	Data	Graphics	Statistics	User	Window	Help
	6	a [• •	g • g	B 💷	00	
	Review	T 4 >	× [.					
	# Comn	nand _	·	/*Creat	e varial	hlee r	read in	regressionst/
	1 doedit	t "C	;	. gen mcl	L mcend=(ores c O	iscu in	regressions-/
	2 do "C; 2 cave "(\Us	-1	-	-			
	J Save (C. \			replace	e mcl_	_mcend=1	if mccain==1 & endorse==1
				(12 real	cnanges	made)	1	
				. gen mcl	L_obend=	0		
/*Create variables used in regressions*/								
gen mcl_mcend=0				(36 real	changes	e mci_ made)	_opena=1	1 11 mccain==1 & endorse==5
<pre>replace mci_mcend=1 11 mccain==1 & endorse==1 gen mcl obend=0</pre>						,		
replace mcl_obend=1 if mccain==1 & endorse==5			ŀ	. gen obl	L_mcend=	0		
gen obl_mcend=0					replace	e obl	mcend=1	if mccain==0 & endorse==1
gen differ=0				(15 real	changes	made)		
replace differ=1 if mcl_obe==1 obl_mce==1								
				. gen dif	ffer=0			
					replace	e diff	fer=1 if	mcl obe==1 obl mce==1
			((51 real	changes	made)		
(

CHECK THE CODE AGAINST THE PUBLISHED PAPER

Example 2:

LIMOR PEER, YALE UNIVERSITY

. /*Table	2*/														
	prob	it ir	itere	st_in_le	tter	mccai	n								
Iteration	0:	log	like	lihood =	-57	.30569	2								
Iteration	1:	log	like	lihood =	- 5	6.1950	5								
Iteration	2:	log	like	lihood =	-56	5.19382	7								
Iteration	3:	log	like	lihood =	-56	5.19382	7								
														_	
Probit re	gress:	ion							Numk	ber of	obs	=		(100	
									LR d	chi2(1))	=		2.22	
									Prok	> ch:	i2	=	0	.1359	1
Log likel:	ihood	= -5	56.19	3827					Psei	ido R2		=	0	.0194	
interest_	in_let	tter		Coef.	3	std. Er	r.	z		P> z		[958	Conf	. Int	erval]
							_								
	mc	cain		4061989		273681	8	1.4	8	0.138		130	2077	. 9	426054
	_(cons	IC	.8557124		.2009	7	-4.2	6	0.000		-1.24	9606	4	618184
-							_								
		n	nfx												
Marginal (effect	ts af	fter	probit											
Y :	= Pr(:	inter	rest_	in_lette	r) (predic	t)								
	= .2	55694	496												
variable		dy/	/dx	Std. E	rr.	z		P> z	[95%	c.ı	.]		Х	
					_				_						-
mccain*		13045	522	.087	06	1.50		0.134	0	040174	.3	01079	•	.49	1

(*) dy/dx is for discrete change of dummy variable from 0 to 1

(table)

Butler and Schofield		365
Table 2. The Effect of Candida	ite Support on Interest in	Letter to the Editor
Dependent Variable = Interested in Publishing Letter Independent Variable	Coefficient (Standard Error) [Change in Probability]	Coefficient (Standard Error) [Change in Probability]
Pro-McCain letter Circulation (in units of 10,000) Unemployment rate in	0.41 (0.27) [13.0%]	0.58** (0.29) [16.3%] -0.028** (0.013) [-16.4% -0.009 (0.13) [-0.3%]
metro area Intercept N Pseudo R ² Log likelihood	-0.86** (0.20)	-0.23 (0.75) 100 .09

Note: The dependent variable is a binary variable that takes the value of 1 if the newspaper either tried to contact the alias for verification purposes or if it published the letter and 0 otherwise. Standard errors are given in parentheses. The estimated predicted probabilities are given in brackets. For the binary variables, the predicted probabilities report the change in the predicted probability when the value of the variable goes from 0 to 1 while holding other variables constant. For the continuous variables, the predicted probabilities report the change in predicted probability when increasing the value of that variable from the mean value to one standard deviation above the mean.

*p < .10. **p < .05.

CURATION TOOL: YARD YALE APPLICATION FOR RESEARCH DATA

 A Web application that allows Depositors, Curators, and Administrators to submit, review, process, and publish data.



Production and code release in 2017-2018

CURATION TOOL: REQUIREMENTS

- Curation workflow automation, integration, management, and tracking
- Integrate and capture DDI metadata production with data and code review and cleaning
- Version control for metadata and files
- Preservation metadata and formats
- ✓ Secure storage and access
- ✓ Preference for open source solutions
- Push out relevant information to pre-determined destinations
 - i.e., a user, the archive administrators, a Web based dissemination system, or preservation systems
- Fit into repository and research workflows

CURATION TOOL: YARD YALE APPLICATION FOR RESEARCH DATA

Log ii	n
Yale ISP	<u>r</u> S
Log in to th password.	e ISPS Data Curation Tool with your username and
Don't have Create an a	a ISPS Data Curation Tool account? account.
Email	
Password	
	Remember me
	Log in
Forgot you	r password?

LIMOR PEER, YALE UNIVERSITY

https://docs.colectica.com/curation/technical-documentation/ddi-mapping/

SOME FINAL THOUGHTS...

IF YOU'RE RECEIVING DATA

When do you engage in the research process?

- ✓ What type of quality assurance do you offer for the data in your repository?
 - Do you review files? (bits, formats)
 - Do you review data content? (validation, missing info, confidentiality)
 - Do you review documentation, metadata? (comprehensiveness, links)
 - Do you review code / analytic methods? (execute code, replicate results)
- ✓ Data repository / archive is the first data re-user



IF YOU'RE RECEIVING DATA

When do you engage in the research process?

- ✓ What type of quality assurance do you offer for the data in your repository?
 - Do you review files? (bits, formats)
 - Do you review data content? (validation, missing info, confidentiality)
 - Do you review documentation, metadata? (comprehensiveness, links)
 - Do you review code / analytic methods? (execute code, replicate results)
- ✓ Data repository / archive is the first data re-user

Can curators get involved in the research process earlier?



IF YOU'RE PROVIDING DATA

Do you establish a practice and a culture of transparency?

- Capture your process.
- Check your code, re-check, and check again. Then cross-check.
- ✓ Leave little (no?) room for mistakes: Use workflow tools.
- ✓ Show (don't tell): Use packaging tools.
- ✓ Aim for long-term re-use.



IF YOU'RE PROVIDING DATA

Do you establish a practice and a culture of transparency?

- Capture your process.
- ✓ Check your code, re-check, and check again. Then cross-check.
- ✓ Leave little (no?) room for mistakes: Use workflow tools.
- ✓ Show (don't tell): Use packaging tools.
- ✓ Aim for long-term re-use.

Can we embed curation practices in the research workflow?



THANK YOU!

limor.peer@yale.edu

@l peer

http://isps.yale.edu/

