# Reproducibility and transparency in academia, and implications for statistical agencies

Lars Vilhuber

2025-10-09

# Follow along



larsvilhuber.github.io/transparency-statistical-agencies/ (HTML zipped, PDF)

## Disclaimer

# Goals of my talk

What is the state of reproducibility and transparency in academic economics?

# What are the benefits of reproducibility and transparency?

# Increasing broad consensus in academia

# What are the implications for statistical agencies?

# AEA Journals

## American Economic Review

The *American Economic Review* is a general-interest economics journal. Established in 1911, the *AER* is among the nation's oldest and most respected scholarly journals in economics.

## American Economic Review: Insights

*AER: Insights* is designed to be a top-tier, general-interest economics journal publishing papers of the same quality and importance as those in the *AER*, but devoted to publishing papers with important insights that can be conveyed succinctly.

## Journal of Economic Literature

The *Journal of Economic Literature (JEL)*, first published in 1969, is designed to help economists keep abreast of and synthesize the vast flow of literature.

## Journal of Economic Perspectives

The *Journal of Economic Perspectives (JEP)* fills the gap between the general interest press and academic economics journals.

## American Economic Journal: Applied Economics

*American Economic Journal: Applied Economics* publishes papers covering a range of topics in applied economics, with a focus on empirical microeconomic issues.

## American Economic Journal: Economic Policy

*American Economic Journal: Economic Policy* publishes papers covering a range of topics, the common theme being the role of economic policy in economic outcomes.

## American Economic Journal: Macroeconomics

*American Economic Journal: Macroeconomics* focuses on studies of aggregate fluctuations and growth, and the role of policy in that context.

## American Economic Journal: Microeconomics

*American Economic Journal: Microeconomics* publishes papers focusing on microeconomic theory; industrial organization; and the microeconomic aspects of international trade, political economy, and finance.

# What is a replication package?

- AEA Data and Code Availability policy

- Data and Code Availability Standard DCAS v1.0

- AEA Data and Code Repository

# AEA policy



Membership   About AEA   Log In

**AMERICAN ECONOMIC ASSOCIATION**

Journals   Annual Meeting   Careers   Resources   EconLit   Committees   Ethics/Ombuds

Home  ›  Journals  ›  AEA Data and Code Policies and Guidance  ›  Data and Code Availability Policy

## Journals

American Economic Review

AER: Insights

AEJ: Applied Economics

AEJ: Economic Policy

AEJ: Macroeconomics

AEJ: Microeconomics

Journal of Economic Literature

Journal of Economic Perspectives

## Data and Code Availability Policy

**It is the policy of the American Economic Association to publish papers only if the data and code used in the analysis are clearly and precisely documented and access to the data and code is nonexclusive to the authors.**

**Authors of accepted papers that contain empirical work, simulations, or experimental work must provide, prior to acceptance, information about the data, programs, and other details of the computations sufficient to permit replication, as well as information about access to data and programs.**

The Editor should be notified at the time of submission if access to the data used in a paper is restricted or limited or if, for some other reason, the requirements above cannot be met.

If data or programs cannot be published in an openly accessible trusted data repository, authors must commit to preserving data and code for a period of no less than five years following publication of the manuscript and to providing reasonable assistance to requests for clarification and replication.

# Tenets of the Policy

- **Transparency**
- **Completeness**
- **Preservation**

# Transparency

- Provenance of the *data*

- Processing of the data, from raw data to results (code)

It is the policy of the American Economic Association to publish papers only if the **data** used in the analysis are **clearly and precisely documented** and **access** to the data and code is **clearly and precisely documented** and is non-exclusive to the authors.

# Completeness

- All data needs to be identified and and access described

- All code needs to be described and provided

Authors ... must provide, prior to acceptance, the **data, programs, and other details** of the computations **sufficient** to permit replication

# Preservation

- All data needs to be preserved for future replicators
  - Ideally, within the replication package, subject to ToU, for convenience
  - Otherwise, in a **trusted repository**

# Preservation

- Code must be in a trusted repository
    - Usually, within the replication package
    - Websites, Github, are ***not acceptable***

# Historically



AER 2011 thanks to Stefano Dellavigna

# Modern preservation

OPEN**ICPSR**

Find Data    Share Data    Repositories

Find Data / Data and Code for: "Indirect Savings from Public Procurement Centralization" / Indirect-Effects-Centralization-main

## Data and Code for: "Indirect Savings from Public Procurement Centralization"

**Principal Investigator(s):** ⍰ Clarissa Lotti, Lear; Arieda Muço, Central European University; Giancarlo Spagnolo, Site - Stockholm School of Economics; Tommaso Valletti, Imperial College London

**Version:** ⍰ V1

**AMERICAN ECONOMIC ASSOCIATION**

| Name ⊡ | File Type ⊡ | Size ⊡ | Last Modified ⊡ |
|---|---|---|---|
| 📁 code | | | 06/18/2024 01:15:PM |
| 📁 data | | | 06/18/2024 01:16:PM |
| 📁 output | | | 06/18/2024 01:14:PM |
| 📄 CITATION.CFF | text/plain | 862 bytes | 06/18/2024 09:14:AM |
| 📄 LICENSE.txt | text/plain | 1.2 KB | 06/18/2024 09:14:AM |
| 📄 README.md | text/x-web-markdown | 6 KB | 06/18/2024 09:14:AM |
| 📄 main.sh | application/x-sh | 2.4 KB | 06/18/2024 09:14:AM |

⊕ DOWNLOAD THIS FOLDER

### Usage Metrics ⍰

**Overall Project Metrics**

| 14 | 3 | 3 |
|---|---|---|
| Views | Downloads | Publications |

**Folder/File-Level Metrics**

| 0 | 0 |
|---|---|
| Views | Downloads |

Download Detailed Metrics

# Side note: Government

- Data are often confidential
  - Are they preserved? (**NARA**, otherwise)
  - Are they accessible to others? (**FSRDC, NORC**, etc.)
- Code is sometimes deemed "confidential"
  - We will return to this topic!

# Exceptions to the Policy

None

●●●

... there is a grey zone:

- When data do not belong to researcher, no control over preservation, access!

- Sometimes, ToU prevent researcher from revealing metadata (name of company, location)

# Transparency again

- However:
  - No exception for need to **describe** access (own and other)
  - No exception for need to fully **describe** processing (possibly with redacted code)

# Enforcement of the AEA Policy

# Reproducibility?

# Reproducibility

"Reproducibility" refers to the ability of a researcher to duplicate the results of a prior study using the same materials and procedures as were used by the original investigator." [1]

# Testing for ...

- **Transparency**
- **Completeness**

through **reproducibility**

# Criteria: Transparency?

- Can a reasonable person understand the description of **acquisition of data** and **processing** via code?

# Criteria: Completeness?

- Do the provided materials allow to reproduce all the **tables** and **figures** in the paper?

# Who is the target person?



Student replicators

# Who is the target person?

Over the past 6 years, over **170** ***undergraduate*** students have been involved in verifying these articles.

- Economics, biostatistics, sociology

- Typically recruited in sophomore or junior year, but will consider freshmen through master's students

# Who is the target person?

- **You** (in 4 years, between prepping 2 new courses, an R&R, a new child, and tenure coming up in 2 years)

- **Your RA** (in 4 years, because you are… see above)

- Your **future readers** who will cite you (in 4-10 years, who may want to extend or replicate your study, but won't if it is too complex)

# Tracing inputs from outputs

# Credibility

## Loss in the Time of Cholera: Long-Run Impact of a Disease Epidemic on the Urban Landscape†

By Attila Ambrus, Erica Field, and Robert Gonzalez*

How do geographically concentrated income shocks influence the long-run spatial distribution of poverty within a city? We examine the impact on housing prices of a cholera epidemic in one neighborhood of nineteenth century London. Ten years after the epidemic, housing prices are significantly lower just inside the catchment area of the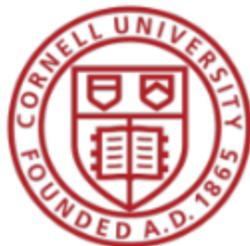 water pump that transmitted the disease. Moreover, differences in housing prices persist over the following 160 years. We make sense of these patterns by building a model of a rental market with frictions in which poor tenants exert a negative externality on their neighbors. This showcases how a locally concentrated income shock can persistently change the tenant composition of a block. (JEL D62, O18, R21, R31)

Indeed, it is the peculiar nature of epidemic disease to create terrible urban carnage and leave almost no trace on the infrastructure of the city.
—Steven Johnson, *The Ghost Map*

Can disease exert a permanent effect on the geography of urban poverty? While it is well understood that illness is impoverishing, because health shocks have no direct impact on infrastructure or land, it is not obvious that epidemics which affect a small number of residents would leave an economic footprint on a city. As the quote above illustrates, a common presumption is that residential migration will preserve the spatial distribution of income in the long run, erasing such shocks from the map over time. In this manner, idiosyncratic income shocks to households should not lead to lasting pockets of poverty in a city. Yet, in reality, spatial discontinuities in urban land values are frequently observed and do not always appear related to discrete changes in local amenities.
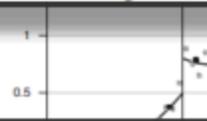
We examine this question in the context of a cholera epidemic that hit a single urban parish of London in 1854. Over the course of one month, 660 residents living

2

ics 126 (1): 145–205.

Lagunoff, Roger, and Akihiko Matsui. 1997. "Asynchronous Choice in Repeated Coordination Games." *Econometrica* 65 (6): 1467–77.

Lalive, Rafael. 2008. "How Do Extended Benefits Affect Unemployment Duration? A Regression Discontinuity Approach." *Journal of Econometrics* 142 (2): 785–806.

**Land Registry.** 2014. "Price Paid Data." http://bit.ly/1HNQAiA (accessed December 19, 2014).

Lee, David S. 2008. "Randomized Experiments from Non-Random Selection in U.S. House Elections." *Journal of Econometrics* 142 (2): 675–97.

Lee, David S., and Thomas Lemieux. 2010. "Regression Discontinuity Designs in Economics." *Journal of Economic Literature* 48 (2): 281–355.

Lee, Sanghoon, and Jeffrey Lin. 2018. "Natural Amenities, Neighbourhood Dynamics, and Persistence in the Spatial Distribution of Income." *Review of Economic Studies* 85 (1): 663–94.

LonRes. 2015. "LonRes: Rental Price Archives." Access provided by Greater London Properties.

the names of the primary occupant at each address ...rds.
...ended in 1963. Hence, for the years 1995–2013 we ... from the Land Registry of England (Land Registry ... property address as well as the sale price and date of ...ntal prices of all properties rented within the Soho area ...May 2015 from the *LonRes* data archives, the primary ...ntral London and only available to verified real estate ...2015, we obtain house value estimates from Zoopla, ...website.[15] We digitized all valuations and addresses ...d addresses above by matching them to housing maps from the relevant time period (for historic records) or using Google's geocoder tool (for current house records).

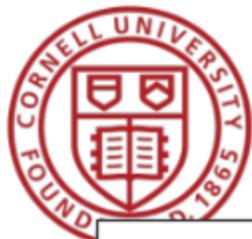To assess the spatial distribution of cholera deaths, we map the total number of deaths by house using the Cholera Inquiry Committee's 1855 map (Cholera Inquiry ...t of a scientific ...local chaplain ...h all residences ...eal disease that ...resulting map, ...r of deaths per ...eath certificates ...which records ...om individuals ...ure from maps

*...long-run spatial distribution of poverty within a city? We examine the impact on housing prices of a cholera epidemic in one neighborhood of nineteenth century London. Ten years after the epidemic, housing prices are significantly lower just inside the catchment area of the water pump that transmitted the disease. Moreover, dif-*

Calonico, Sebastian, Matias D. Cattaneo, and Rocío Titiunik. 2015. "Optimal Data-Driven Regression Discontinuity Plots." *Journal of the American Statistical Association* 110 (512): 1753–69.

Card, David, Alexandre Mas, and Jesse Rothstein. 2008. "Tipping and the Dynamics of Segregation." *Quarterly Journal of Economics* 123 (1): 177–218.

**Cholera Inquiry Committee.** 1855. *Report on the Cholera Outbreak in the Parish of St. James, Westminster, during the Autumn of 1854.* London: J. Churchill.

Conley, T. G. 1999. "GMM Estimation with Cross Sectional Dependence." *Journal of Econometrics* 92 (1): 1–45.

Conley, Timothy G. 2008. "Spatial Econometrics." Unpublished.

Dell, Melissa. 2010. "The Persistent Effects of Peru's Mining Mita." *Econometrica* 78 (6): 1863–1903.

https://doi.org/10.1257/aer.20190759

# Data and Code for: Loss in the Time of Cholera: Long-run Impact of a Disease Epidemic on the Urban Landscape

**Principal Investigator(s):** ❷ Attila Ambrus, Duke University; Erica Field, Duke University; Robert Gonzalez, University of South Carolina

**AMERICAN ECONOMIC ASSOCIATION**

**Version:** ❷ V2

## Do-files, input Data, and Output Figures and Tables

NOTE: Master do-file (Master.do) provides all Tables and Figures

| Do-file | Input datasets | Output |
|---|---|---|
| Table_summary_stats.do | houses_1853_final.dta | Table 1<br>Table B1 |
| Table_deaths.do | Merged_1853_1864_data.dta | Table 2 |
| Table_main_results.do | Merged_1853_1864_data.dta<br>Merged_1846_1894_data.dta<br>houses_1936_final.dta | Table 3 |
| Table_moved.do | Merged_1853_1864_data.dta | Table 4 |
| Table_migration.do | Merged_1853_1864_data.dta | Table 5 |
| Table_census.do | Data_census.dta | Table 6 |
| Table_Booth_data.do | final_booth_RG.dta | Table 7 |
| Table_current_results.do | houses_current_final.dta<br>current_rentals_final.dta | Table 8 |
| Fig_RD_plots.do | Merged_1853_1864_data.dta<br>Merged_1846_1894_data.dta<br>houses_1936_final.dta<br>Data_census.dta<br>final_booth_RG.dta<br>houses_current_final.dta<br>current_rentals_final.dta | Figure 2<br>Figure 3<br>Figure B1<br>Figure B2<br>Figure B3<br>Figure B4<br>Figure B5 |
| Fig_variance_grid.do | grid_house_final | Figure 4 |
| Table_fuzzy_iv.do | Merged_1853_1864_data.dta | Table B2 |

st Modified ❏

/02/2019 02:23:PM

/21/2019 10:47:AM

of Cholera: Long-run Impact of
ation [publisher], 2020. Ann
20-01-31. https://doi.org

un spatial distribution of
one neighborhood of 19th
t inside the catchment area of
persist over the following 160
tions in which poor tenants

# Data and Code for: Loss in the Time of Cholera: Long-run Impact of a Disease Epidemic on the Urban Landscape

**Principal Investigator(s):** ❓ Attila Ambrus, Duke Unive... University; Robert Gonzalez, University of South Carolin...

**Version:** ❓ V2

**Do-files, input Data, and Output Figures and Tables**

NOTE: Master do-file (Master.do) provides all Tables and Figures

| Do-file | Input datasets | |
|---|---|---|
| Table_summary_stats.do | houses_1853_final.dta | |
| Table_deaths.do | Merged_1853_1864_data.dta | |
| Table_main_results.do | Merged_1853_1864_data.dta | |
| | Merged_1846_1894_data.dta | |
| | houses_1936_final.dta | |
| Table_moved.do | Merged_1853_1864_data.dta | |
| Table_migration.do | Merged_1853_1864_data.dta | |
| Table_census.do | Data_census.dta | |
| Table_Booth_data.do | final_booth_RG.dta | |
| Table_current_results.do | houses_current_final.dta | Table 8 |
| | current_rentals_final.dta | |
| Fig_RD_plots.do | Merged_1853_1864_data.dta | Figure 2 |
| | Merged_1846_1894_data.dta | Figure 3 |
| | houses_1936_final.dta | Figure B1 |
| | Data_census.dta | Figure B2 |
| | final_booth_RG.dta | Figure B3 |
| | houses_current_final.dta | Figure B4 |
| | current_rentals_final.dta | Figure B5 |
| Fig_variance_grid.do | grid_house_final | Figure 4 |
| Table_fuzzy_iv.do | Merged_1853_1864_data.dta | Table B2 |

| Name ▾ | File Type ▲ |
|---|---|
| 📊 Data_census.dta | application/ |
| 📊 Merged_1846_1894_data.dta | application/ |
| 📊 Merged_1853_1864_data.dta | application/ |
| 📊 | plication/ |

| Name ▾ | File Ty |
|---|---|
| 📁 mccrary-s-ado | |
| 📁 spatial_HAC | |
| 📊 Fig_RD_plots.do | text/x- |
| 📊 Fig_bandwidth_sensitivity.do | text/x- |
| 📊 Fig_pre-trends.do | text/x- |

...one neighborhood of 19th
...t inside the catchment area of
...ersist over the following 160
...tions in which poor tenants

# Data and Code for: Loss in the Time of Cholera: Long-run Impact of a Disease Epidemic on the Urban Landscape

**File Type**

application/

application/

Do-files, i

NOTE: Ma

| Do-file |
|---|
| Table_su |
| Table_de |
| Table_ma |
| Table_mo |
| Table_mi |
| Table_ce |
| Table_Bo |
| Table_cu |

application/

data.dta

data.dta

**File Ty**

plication/

Fig_RD_p

itivity.do

text/x-

text/x-

text/x-

Fig_varia

| Table_fuzzy_iv.do | Merged_1853_1864_data.dta | Table B2 |

neighborhood of 19th
inside the catchment area of
persist over the following 160
tions in which poor tenants

```stata
1   *===============================================================
2   * Purpose: Do-file creates RDplots
3   * Outcome:
4   * Figure 2:  Cholera Deaths and BSP Boundary (1854)
5   * Figure 3:  RD plots for Main Outcomes (in logs)
6   * Figure B1: Covariate RD Plots (1853)
7   * Figure B2: Histogram and Density of Forcing Variable (Distance to BSP boundary)
8   * Figure B3: RD Plots for Residential Mobility Outcome
9   * Figure B4: RD Plots for House Occupancy Outcomes
10  * Figure B5: RD Plots for Socioeconomic Outcomes
11  *===============================================================
12
13  clear all
14  set more off
15
16
17
18  *****************************************************************
19  *  Figure 2a, 2b: Cholera Deaths and BSP Boundary (1854)
20  *****************************************************************
21  * RD Program
22  capture program drop myrdplot
23  program define myrdplot
24  args outcome
25
26      * large sample
27      local width = 20
28      local hwidth = 10
29      local limit = 100 - `width'
30      local gr_limit = `limit'+`width'
31      local gr_width = `gr_limit'/4
```

# Reproducibility



Find Data / Data and Code for: Loss in the Time of Cholera: Long-run Impact of a Disease Epidemic on the Urban Landscape

## Data and Code for: Loss in the Time of Cholera: Long-run Impact of a Disease Epidemic on the Urban Landscape

**Principal Investigator(s):** ❓ Attila Ambrus, Duke University; Erica Field, Duke University; Robert Gonzalez, University of South Carolina

**Version:** ❓ V2

3

**Version Title:** ❓ Corrected author information

| Name ⊡ | File Type ⊡ | Size ⊡ | Last Modified ⊡ |
|---|---|---|---|
| 📁 aer_replication | | | 09/02/2019 02:23:PM |
| 📄 README.pdf | application/pdf | 587 KB | 08/21/2019 10:47:AM |

**Project Citation:**

Ambrus, Attila, Field, Erica, and Gonzalez, Robert. Data and Code for: Loss in the Time of Cholera: Long-run Impact of a Disease Epidemic on the Urban Landscape. Nashville, TN: American Economic Association [publisher], 2020. Ann Arbor, MI: Inter-university Consortium for Political and Social Research [distributor], 2020-01-31. https://doi.org/10.3886/E111523V2

## Project Description

**Summary:** ❓ How do geographically concentrated income shocks influence the long-run spatial distribution of poverty within a

**⊕ DOWNLOAD THIS PROJECT**

**Usage Metrics** ❓

**Overall Project Metrics**

| 597 | 155 | 1 |
|---|---|---|
| Views | Downloads | Publications |

Download Detailed Metrics

**Published Versions**

AMERICAN ECONOMIC ASSOCIATION

Credibility

# Reproducibility in Economics and beyond

Social Science
Data Editors

# Social Science Data Editors

Improving reproducibility in the social and economic sciences

---

## Data and Code Availability Standard

DCAS The Data and Code Availability Standard (DCAS) is a standard for sharing research code and data, endorsed by leading journals in social sciences. See https://datacodestandard.org/ for more information.

DOI  10.5281/zenodo.7436134

DCAS

Data and Code Availability Standard

About    Journals

| | Data | |
|---|---|---|
| 1 | Data Availability Statement | A Data Availability Statement is provided with detailed enough information such that an independent researcher can replicate the steps needed to access the original data, including any limitations and the expected monetary and time cost of data access. |
| 2 | Raw data | Raw data used in the research (primary data collected by the author and secondary data not otherwise available) is made publicly accessible. Exceptions are explained under Rule 1. |

# Data Editors

- American Economic Association (8)
- Econometric Society (3)
- Canadian Journal of Economics (1)
- Royal Economic Society (2)
- Western Economic Association International (1)
- European Economic Association (1)
- Review of Economic Studies (1)
- **Journal of the European Economic Association** (1)
- **Journal of Political Economy** (3)

## ƆCAS
Data and Code Availability Standard

## Journals

The following journals endorse the Data and Code Availability Standard.

1. American Economic Journal: Applied Economics 🐦 🐘
2. American Economic Journal: Economic Policy 🐦 🐘
3. American Economic Journal: Macroeconomics 🐦 🐘
4. American Economic Journal: Microeconomics 🐦 🐘
5. American Economic Review 🐦 🐘
6. American Economic Review: Insights 🐦 🐘
7. Canadian Journal of Economics 🐦
8. Econometrica 🐦
9. Econometrics Journal
10. Economic Inquiry 🐦
11. Economic Journal 🐦
12. Journal of Economic Literature 🐦 🐘
13. Journal of Economic Perspectives 🐦 🐘
14. Journal of the European Economic Association 🐦
15. Quantitative Economics 🐦
16. Review of Economic Studies 🐦
17. Theoretical Economics 🐦

# Common policies

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

# Elsewhere: Political Science



APSR



AJPS

# Elsewhere: Sociology



sociological science

Articles    For Authors

Home › Reproducibility Policy

## Reproducibility Policy

Over the last decade, we have witnessed a crisis in science in which many admired research studies
non-replicable. Researchers increasingly recognize that publication itself does not imply that findings
questioned the credibility of social science research. In order to advance the credibility of sociologica
has adopted a reproducibility policy.

Starting with submissions received after April 1, 2023, authors of articles relying on statistical or com
required to deposit replication packages as a condition of publication in *Sociological Science*. Replic
the statistical code and — when legally and ethically possible — the data required to fully reproduce
policy, Sociological Science hopes other high-impact journals in Sociology will follow suit in setting s
published work.

In addition to depositing replication packages, papers relying on experimental methods must adhere
registration requirements outlined in the journal's Policy on Findings from Experimental Data below.

Under many legitimate circumstances, data cannot legally or ethically be made available to readers.
data available, they must explain why in the main text of the paper. In such cases, making code and
required, unless doing so would violate legal or ethical constraints.

Researchers using qualitative data, such as interviews or participant observation data, are not requi
package. We encourage authors to make qualitative data available when possible, and urge them to
as interview protocols or coding schemes can be shared.

4

# Trust in Government Statistics

# United Nations

Fundamental Principles of Official Statistics, **Principle 3**:

> ***Accountability and Transparency*** To facilitate a correct interpretation of the data, the statistical agencies are to present information according to scientific standards on the sources, methods and procedures of the statistics [5]

# National Academies

Principles and Practices for a Federal Statistical Agency, **Principle 2**:

> Credibility among Data Users A federal statistical agency must have credibility with those who use its data and information [6]

# OMB

"... flow of objective, **credible** statistics to support the decisions of individuals, households, governments, businesses, and other organizations."

# OMB

Statistical Policy Directive No. 1, 4:

"Any **loss of trust** in the integrity of the Federal statistical system and its products could lessen respondent cooperation with Federal statistical surveys, decrease the quality of statistical system products, and foster uncertainty about the validity of measures our Nation uses to monitor and assess its performance and progress."
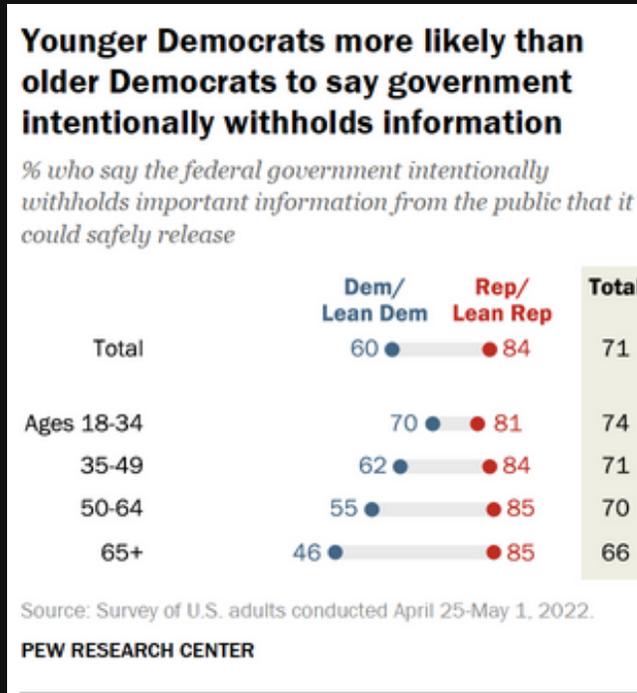
# Agency efforts

## U.S. Census Bureau



U.S. Census Trust and Safety Center

# Agency efforts



Responding to the Pandemic with Trusted Economic Analysis

In Fiscal Year (FY) 2020, the nation faced one of its greatest challenges that impacted nearly every facet of day-to-day life. As the coronavirus pandemic spread to nearly every country in the world, agricultural and manufacturing supply chains experienced dramatic challenges. Stay-at-home orders became common changing many people's eating habits, from where they bought and consumed their food, to how much it cost, and whether they had enough to eat.

ERS Annual Report 2020



Collaborating Across Government to Respond to Emerging Issues with Trusted Economic Information

Each year, ERS releases widely cited reports on familiar topics, such as food security and farm income, that help everyone from policymakers to consumers make better decisions related to agriculture and food. However, ERS also has another important role within USDA and the Federal government that often goes unseen.

ERS Annual Report 2021

# Joint Statement

Joint Statement on Commitment to Scientific Integrity and Transparency

- **Principle 2**: a Federal statistical agency must have credibility with those who use its data and information;

- **Principle 3**: a Federal statistical agency must have the trust of those whose information it obtains;
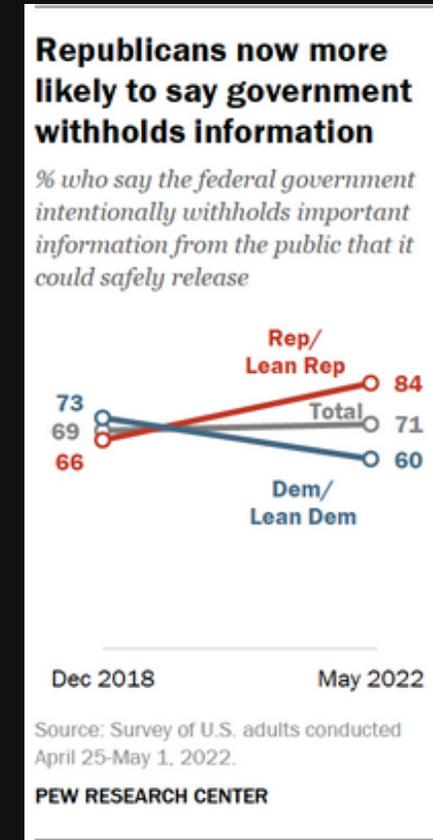


Joint Statement

# Waning trust

Pew Research



Democrats on withholding data



Republicans on withholding data

# Computational Reproducibility and Official Statistics

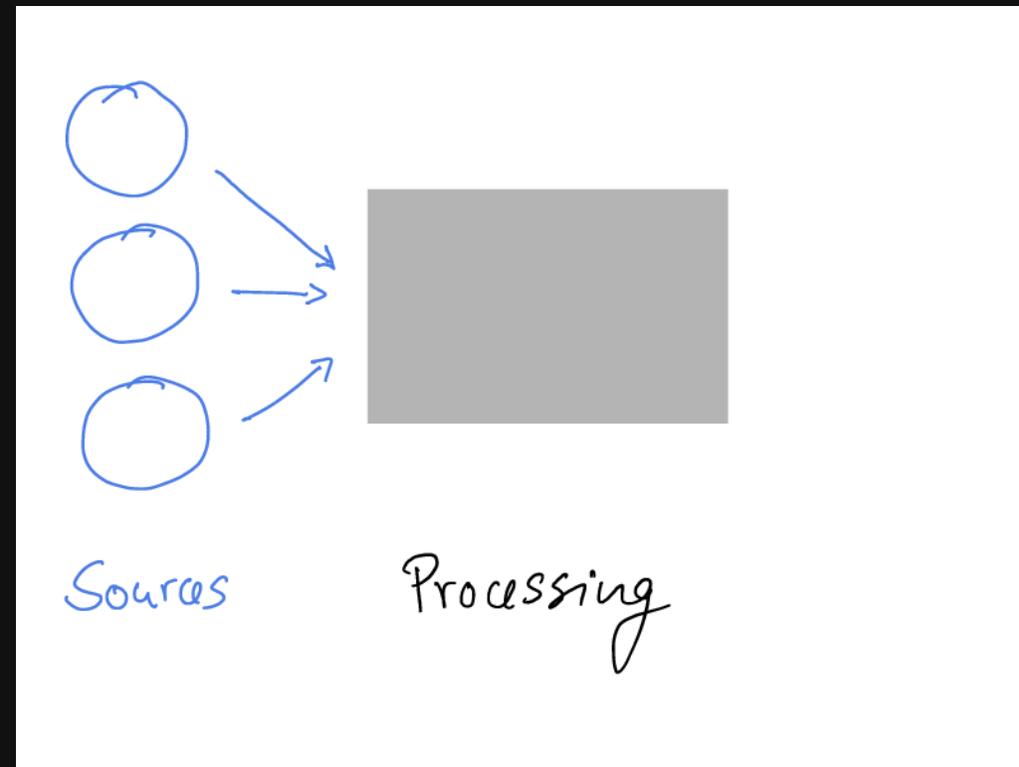Agencies do provide
**detailed information on sources**

- Surveys

- Administrative data

# Computational Reproducibility and Official Statistics

But: **Availability of "computing instructions"?**

- *Code* for cleaning, aggregation, imputation
- Including for *disclosure avoidance*



Sources    Processing

# Computational Reproducibility and Official Statistics

But: Availability of **reliable, trusted data archives**

- Of released data – ability to reproduce downstream uses
- Of source data – ability to reproduce released data

# The analogy

# The analogy



Sources      Processing

# The analogy

# Some principles from the academic world

Which are starting to be infused into the federal system

# FAIR Principles

FAIR:

- **F**indable
- **A**ccessible
- **I**nteroperable
- **R**eusable

# Data Citation Principles



## To make it findable,



Data Citation Principles

7

### 1. Importance

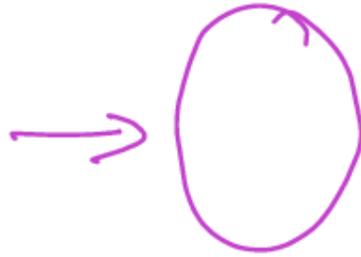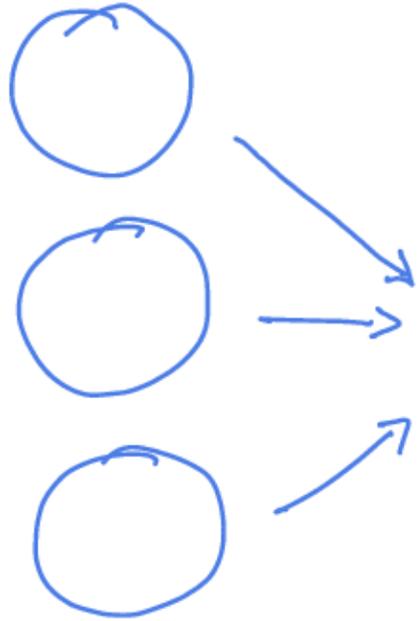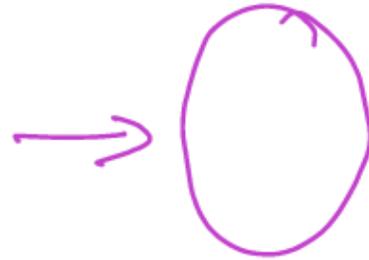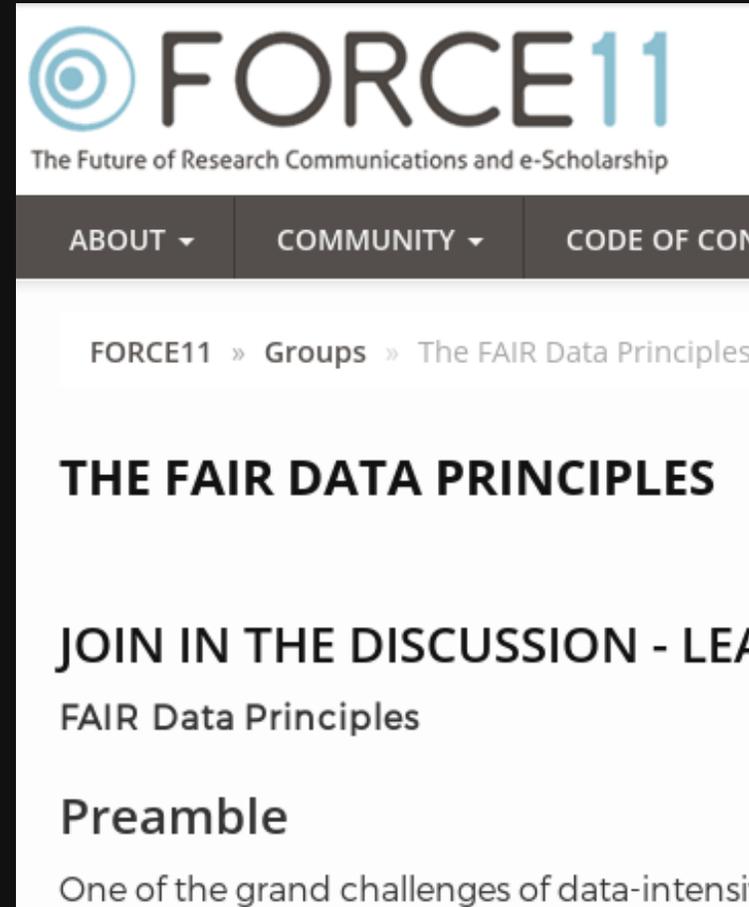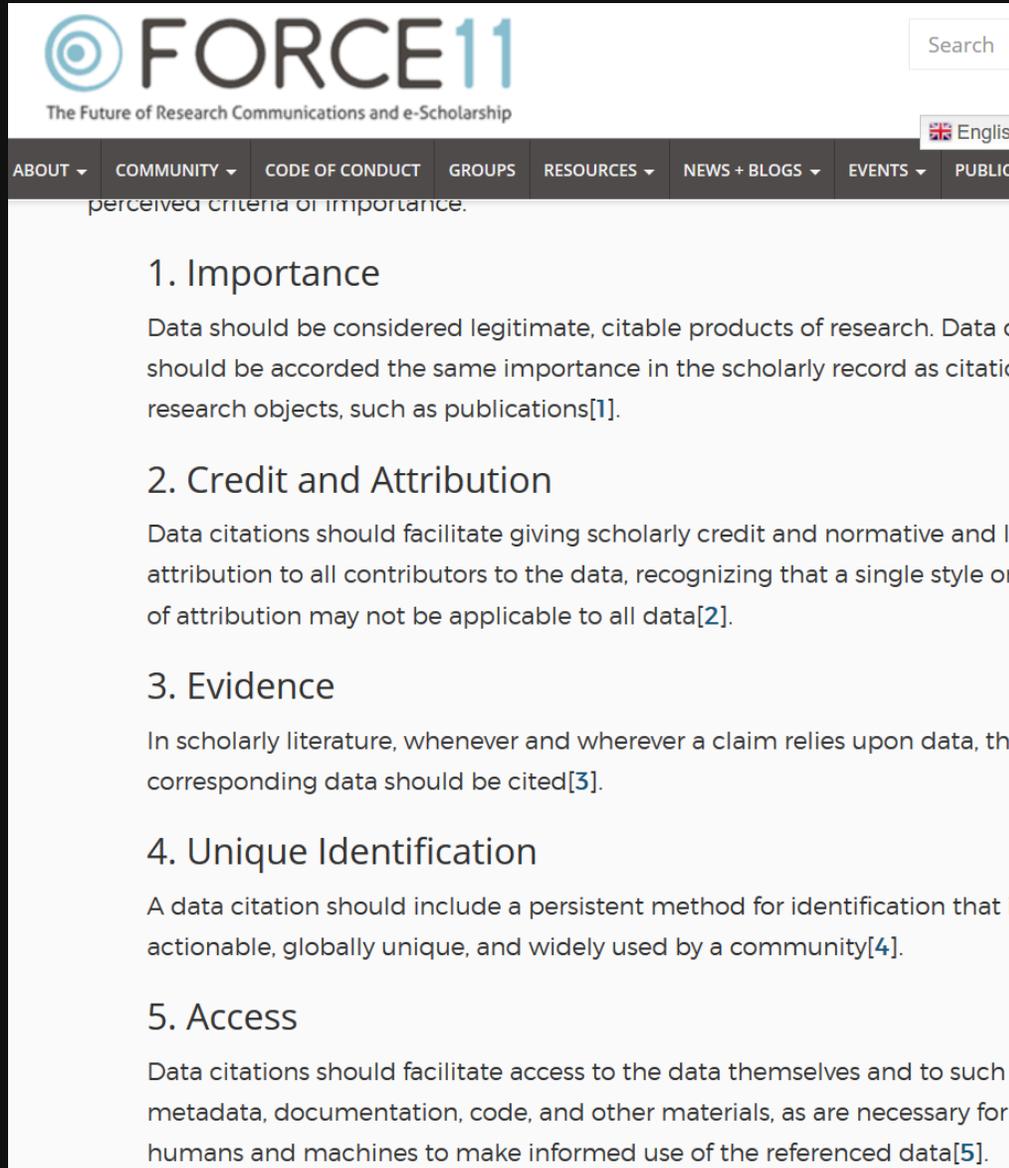Data should be considered legitimate, citable products of research. Data citations should be accorded the same importance in the scholarly record as citations of other research objects, such as publications[1].

### 2. Credit and Attribution

Data citations should facilitate giving scholarly credit and normative and legal attribution to all contributors to the data, recognizing that a single style or mechanism of attribution may not be applicable to all data[2].

### 3. Evidence

In scholarly literature, whenever and wherever a claim relies upon data, the corresponding data should be cited[3].

### 4. Unique Identification

A data citation should include a persistent method for identification that is actionable, globally unique, and widely used by a community[4].

### 5. Access

Data citations should facilitate access to the data themselves and to such associated metadata, documentation, code, and other materials, as are necessary for both humans and machines to make informed use of the referenced data[5].

# An example from ERS

# Website



Website

# Sources

**Data Sources**

The ARS and RRS are calculated using two data sources. The first is the 7.5 arc-second resolution, Global Multi-resolution Terrain Elevation Data 2010 (GMTED2010) from U.S. Department of the Interior, U.S. Geological Survey and U.S. Department of Defense, National Geospatial-Intelligence Agency. The second is the vintage 2010 census tract TIGER (Topographically Integrated Geographic Encoding and Referencing)/Line boundary files from the U.S. Department of Commerce, Bureau of the Census. Additionally, ESRI's ArcGIS StreetMap Premium 2021, North American Q3 road network data were used when creating the RRS. All road types in the dataset were used, including highways, arterial, collector, local, and semi-private roads. Finally, population, population density, land area, and rurality data for vintage 2010 census tracts are from USDA, ERS's Rural-Urban Commuting Area Codes data product.

High-level description of sources

# Sources

U.S. Bureau of the Census, U.S. Department of Commerce. (2012). *2010 census tract TIGER (Topographically Integrated Geographic Encoding and Referencing)/Line shapefile.*

Sources are cited!

# Methods

Data are available for the vintage 2010 census tracts within the 50 States and Washington, DC.

## Methods

Creating the Area Ruggedness Scale (ARS) and Road Ruggedness Scale (RRS) was a three-step process.

### Step 1: Computing a Grid Cell Terrain Ruggedness Index

The ruggedness measures are based on the Terrain Ruggedness Index (TRI) developed by the Riley et al. (1999) article, A Terrain Ruggedness Index that Quantifies Topographic Heterogeneity. The TRI is calculated using data from a digital elevation model (DEM), a detailed representation of the Earth's terrain at the scale of small, regularly spaced grid cells. A TRI value is computed for each grid cell by calculating "the sum change in elevation between [the given] grid cell and its eight neighbor cells," as illustrated below. Lower values indicate less change in elevation within the 3-by 3-grid-cell neighborhood, and higher values indicate areas with higher elevation differences.

Two sets of TRI values were calculated using grid cells that average 0.15 square miles in size, one including all territory (the Area TRI) and one including just those grid cells containing roads (the Road TRI). For the Area TRI, ruggedness values were calculated for all 263 million grid cells covering the United States. For the Roads TRI,

High-level description of methods, but no (obvious) code

# Methods

Esri. (2021). *ArcGIS StreetMap premium* (North America 2021 Release 3) [Data set]. Esri.

Evans, J.S., & Murphy, M.A. (2021, May 14). *Spatial analysis and modelling utilities*, version 1.3-7. CRAN – Package spatialEco.

Some methods - R code - is cited

# Own citation?

**Recommended Citation**

U.S. Department of Agriculture, Economic Research Service. *Area and Road Ruggedness Scales*, September 2023.

Own citation does not include a URL

# Findability?



Data.gov is not great

# Findability?



Not even close.

Google Dataset Search is worse

# Repeatability of downloads?

URL is

**https://ers.usda.gov/webdocs/DataFiles/107356/Ruggedne**
**v=6316.8**

- What will the 2020 tract-based data URL look like?

# Reproducibility?

- Most of the data inputs seem to be public data, or commercially available (ESRI)

- If the code were provided, others should be able to reproduce the analysis

# Opening up technical possibilities

# How can we know that a data source is reliably obtained?

# Consider the case of Gino



The New York Times

Account ⌄

## The Harvard Professor and the Bloggers

When Francesca Gino, a rising academic star, was accused of falsifying data — about how to stop dishonesty — it didn't just torch her career. It inflamed a crisis in behavioral science.

Francesca Gino

# The case of Gino

- Francesca Gino was a tenured professor at Harvard Business School, writing on honesty (!)

# The case of Gino

- Several articles were investigated by third parties (Data Colada, in particular [9]), and found to be problematic

# The case of Gino

- At least one of them had manipulated data **AFTER** it had been collected, **BEFORE** it had been analyzed.



Data manipulation



Results of manipulation

# Generic survey processing

# Generic survey processing

# Requiring transparency in academia



Generic survey processing

# Verifying transparency in academia



Generic survey processing

# Verification by journals

- **Provision** (publication of materials) provides transparency

- **Verification** (running the analysis again - computational reproducibility) compensates for *mistrust*/*absence of trust*

# Which journals again

- American Economic Association (8)

- Econometric Society (3)

- Canadian Journal of Economics (1)

- Royal Economic Society (2)

- Western Economic Association International (1)

- European Economic Association (1)

- Review of Economic Studies (1)

- Journal of the European Economic Association (1)

- Journal of Political Economy (3)

- American Journal of Political Science (1)

- American Political Science Review (1)

# Verification by others

- Pre-publication: cascad

- Post-publication: Data Colada, Institute for Replication



cascad



I4R

# Verification by institutions

- World Bank[10]



World Bank RRR

# Taking it a step further



Survey flow

# Taking it a step further

- Has been discussed by authors behind Data Colada

- Survey tool provider (Qualtrics, etc.) exports data, posts checksum

- Survey tool provider exports data only to institution directly into trusted repository, researchers obtain data from there (with privacy protections)

# Does not prevent all fraud

## Toronto researcher loses Ph.D.

**Exclusive: Psychology researcher loses PhD after allegedly using husband in study and making up data**

A psychology researcher already under fire for several questionable studies has had her PhD revoked by a university tribunal that found it likely she fabricated data in her thesis.

Ping Dong, who was a doctoral student at the University of Toronto from 2012 to 2017, had already earned retractions for two papers based on her thesis before the tri-

Ping Dong

Toronto case

## MIT student makes up firm data

AI

MIT disavows doctoral student paper on AI's productivity benefits

Anthony Ha — 12:30 PM PDT · May 17, 2025

MIT case

# Can Assurances be created for Statistical Agencies?



Survey flow

# Challenges for statistical agencies

- Documenting full transparency is hard
- Complex and legacy survey tools make the process harder
- Presence of (legitimate!) manual edits is an issue
- Production processes are long and complex
- Much of the code base is not open source

# How to document the full process?



Survey flow

# A sketch: Transparency Certified

https://transparency-certified.github.io/

TRAnsparency CErtified (TRACE)                    Home    Jobs    Specification    Infrastructure    About

## TRACE: Building trust in computational research

### A new approach to computational transparency and reproducibility

## Trusting computational research without repeating it

How can we trust the integrity of results from research that relies on computations without repeating them? By certifying the successful _original_ execution of a computational workflow that produced findings _in situ_. With certifications in hand, consumers of research can trust the transparency of results without necessarily repeating computations. Learn more

# Work in progress

- Working with cascad, several INEXDA members, and others

- Relying on external certification of data inputs (data catalogs with metadata, checksums)

# Wrapping it all up

# What is the state of reproducibility and transparency in academic economics?

- An increasing number of journals are not just **requiring** complete data, code, and transparent description, but also **verifying** that the code and data are correct.

- At the AEA: since 2019, reviewed around **3000** articles, ran code for about **2/3** of them.

# What are the benefits of reproducibility and transparency?

- Greater **trust** in the results

- Greater **ease** of building on results

- Greater **transparency** of the process, but also of the **provenance**

# Increasing broad consensus in academia

- FAIR principles

- Data Citation Principles

- Computational Reproducibility

# What are the implications for statistical agencies?

# Producers of statistical products

- May want to provide greater transparency into the **production process**.

- May need to do more for **long-term, unbiased preservation** of input data, output products, and code/software to link the two.

- May want to start with the low-hanging fruit: ***dashboards and fully public processes***.

# Coherence with stated principles

The emerging consensus is fully in line with the **decades-strong principles** of statistical agencies:

# Greater trust by the public?

- **Transparency** should be correlated with greater trust in the work of the statistical agencies

- But: Transparency can also lead to vulnerability through misinterpretation (no panacea)

# Thank you

# One more thing...

# That confidential code thing...

- IRS variable names

- File paths b/c your IT department said so

- Use of confidential **data** in code (`if name="Lars" then confid=2`)

# Solution

Don't do that.

# Solution

labordynamicsinstitute.github.io/reproducibility-confidential/

**Reproducibility when data are confidential**

Lars Vilhuber

2024-10-01

Also here.

# Appendix

# Secrets in the code

# What are secrets?

- API keys

- Login credentials for data access

- File paths (FSRDC!)

- Variable names (IRS!)

# Standard practice

Store secrets in environment variables or files that are not published.

# Some services are serious about this

**About secret scanning**

GitHub scans repositories for known types of secrets, to prevent fraudulent use of secrets that were committed accidentally.

Github secret scanning

# Where to store secrets

- **environment variables**

- "dot-env" files (Python), "Renviron" files (R)

- or some other clearly identified file in the project or home directory

# Environment variables

Typed interactively (here for Linux and Mac)

```
1  MYSECRET="dfad89ald"
2  CONFDATALOC="/path/to/irs/files"
```

(this is **not** recommended)

# Storing these in files

Same syntax used for contents of "dot-env" or "Renviron" files, and in fact `bash` or `zsh` startup files (`.bash_profile`, `.zshrc`)

# Using In R

Edit `.Renviron` (note the dot!) files:

```r
1  # Edit global (personal) Renviron
2  usethis::edit_r_environ()
3  # You can also consider creating project-specific settings:
4  usethis::edit_r_environ(scope = "project")
```

Use the variables defined in `.Renviron`:

```r
1  mysecret <- Sys.getenv('MYSECRET')
```

# Using In Python

## Loading regular environment variables:

```python
1  import os
2  mysecret = os.getenv("MYSECRET")   # will load environment variables
```

## Loading with dotenv

```python
1  from dotenv import load_dotenv
2  load_dotenv()   # take environment variables from project .env.
3  mysecret = os.getenv("MYSECRET")   # will load environment variables
```

# Using in Stata

Yes, this also works in Stata

```
1  // load from environment
2  global mysecret : env MYSECRET
3  display "$mysecret"  // don't actually do this in code
```

and via (what else) a user-written package for loading
from files:

```
1  net install doenv, from(https://github.com/vikjam/doenv/raw/master/)
2  doenv using ".env"
3  global mysecret "`r(MYSECRET)'"
4  display "$mysecret"
```

# Simplest solution

```stata
 1   //============ non-confidential parameters =========
 2   include "config.do"
 3
 4   //============ confidential parameters =============
 5   capture confirm file "$code/confidential/confparms.do"
 6   if _rc == 0 {
 7       // file exists
 8       include "$code/confidential/confparms.do"
 9   } else {
10       di in red "No confidential parameters found"
11   }
12   //============ end confidential parameters =========
```

# Confidential code?

# What is confidential code, you say?

- In the United States, some **variables on IRS databases** are considered super-top-secret. So you can't name that-variable-that-you-filled-out-on-your-Form-1040 in your analysis code of same data. (They are often referred to in jargon as "Title 26 variables").

# What is confidential code, you say?

- Your code contains the **random seed you used to anonymize** the sensitive identifiers. This might allow to reverse-engineer the anonymization, and is not a good idea to publish.

# What is confidential code, you say?

- You used a **look-up table hard-coded** in your Stata code to anonymize the sensitive identifiers (`replace anoncounty=1 if county="Tompkins, NY"`).

A **really bad idea**, but yes, you probably want to hide that.

# What is confidential code, you say?

- Your IT specialist or disclosure officer thinks publishing the **exact path** to your copy of the confidential 2010 Census data, e.g., "/data/census/2010", is a security risk and refuses to let that code through.

# What is confidential code, you say?

- You have adhered to disclosure rules, but for some reason, the precise minimum cell size is a confidential parameter.

# What is confidential code, you say?

So whether reasonable or not, **this is an issue**. How do you do that, without messing up the code, or spending hours redacting your code?

# Example

- This will serve as an example. None of this is specific to Stata, and the solutions for R, Python, Julia, Matlab, etc. are all quite similar.

- Assume that variables `q2f` and `q3e` are considered confidential by some rule, and that the minimum cell size `10` is also confidential.

```
1  set seed 12345
2  use q2f q3e county using "/data/economic/cmf2012/extract.dta", clear
3  gen logprofit = log(q2f)
4  by county: collapse (count)  n=q3e (mean) logprofit
5  drop if n<10
6  graph twoway n logprofit
```

# Example

Only one line that does not contain "confidential" information.

```
1  set seed 12345
2  use q2f q3e county using "/data/economic/cmf2012/extract.dta", clear
3  gen logprofit = log(q2f)
4  by county: collapse (count)  n=q3e (mean) logprofit
5  drop if n<10
6  graph twoway n logprofit
```

# Do not do this

A bad example, because literally making more work for you and for future replicators, is to manually redact the confidential information with text that is not legitimate code:

```
1  set seed NNNNN
2  use <removed vars> county using "<removed path>", clear
3  gen logprofit = log(XXXX)
4  by county: collapse (count)  n=XXXX (mean) logprofit
5  drop if n<XXXX
6  graph twoway n logprofit
```

The redacted program above will no longer run, and will be very tedious to un-redact if a subsequent replicator obtains legitimate access to the confidential data.

# Better

Simply replacing the confidential data with replacement that are valid placeholders in the programming language of your choice is already better. Here's the confidential version of the file:

```stata
 1  //============ confidential parameters =============
 2  global confseed    12345
 3  global confpath    "/data/economic/cmf2012"
 4  global confprofit  q2f
 5  global confemploy  q3e
 6  global confmincell 10
 7  //============ end confidential parameters =========
 8  set seed $confseed
 9  use $confprofit county using "${confpath}/extract.dta", clear
10  gen logprofit = log($confprofit)
11  by county: collapse (count)  n=$confemploy (mean) logprofit
12  drop if n<$confmincell
13  graph twoway n logprofit
```

# Better

and this could be the released file, part of the replication package:

```
 1  //=========== confidential parameters ============
 2  global confseed   XXXX      // a number
 3  global confpath    "XXXX"   // a path that will be communicated to you
 4  global confprofit  XXX      // Variable name for profit T26
 5  global confemploy  XXX      // Variable name for employment T26
 6  global confmincell XXX      // a number
 7  //=========== end confidential parameters =========
 8  set seed $confseed
 9  use $confprofit county using "${confpath}/extract.dta", clear
10  gen logprofit = log($confprofit)
11  by county: collapse (count)  n=$confemploy (mean) logprofit
12  drop if n<$confmincell
13  graph twoway n logprofit
```

While the code won't run as-is, it is easy to un-redact, regardless of how many times you reference the confidential values, e.g., q2f, anywhere in the code.

# Best

- Main file

- Conditional processing

- Separate file for confidential parameters which can simply be excluded from disclosure request

# Best

Main file `main.do`:

```stata
 1  //=========== confidential parameters =============
 2  capture confirm file "$code/confidential/confparms.do"
 3  if _rc == 0 {
 4      // file exists
 5      include "$code/confidential/confparms.do""
 6  } else {
 7      di in red "No confidential parameters found"
 8  }
 9  //=========== end confidential parameters =========
10
11  //=========== non-confidential parameters =========
12  global safepath "$rootdir/releasable"
13  cap mkdir "$safepath"
14
15  //=========== end parameters =====================
```

# Best

Main file `main.do` (continued)

```stata
 1  // ::::  Process only if confidential data is present
 2
 3  capture confirm  file "${confpath}/extract.dta"
 4  if _rc == 0 {
 5      set seed $confseed
 6      use $confprofit county using "${confpath}/extract.dta", clear
 7      gen logprofit = log($confprofit)
 8      by county: collapse (count)  n=$confemploy (mean) logprofit
 9      drop if n<$confmincell
10      save "${safepath}/figure1.dta", replace
11  } else { di in red "Skipping processing of confidential data" }
12
13  //=========== at this point, the data is releasable ======
14  // ::::  Process always
15
16  use "${safepath}/figure1.dta", clear
17  graph twoway n logprofit
18  graph export "${safepath}/figure1.pdf", replace
```

# Best

Auxiliary file `$code/confidential/confparms.do"` (not released)

```
1  //=========== confidential parameters ============
2  global confseed   12345
3  global confpath   "/data/economic/cmf2012"
4  global confprofit q2f
5  global confemploy q3e
6  global confmincell 10
7  //=========== end confidential parameters =========
```

# Best

Auxiliary file `$code/include/confparms_template.do` (this is released)

```
1  //=========== confidential parameters =============
2  // Copy this file to $code/confidential/confparms.do and edit
3  global confseed    XXXX     // a number
4  global confpath    "XXXX"   // a path that will be communicated to you
5  global confprofit  XXX      // Variable name for profit T26
6  global confemploy  XXX      // Variable name for employment T26
7  global confmincell XXX      // a number
8  //=========== end confidential parameters =========
```

# Best replication package

Thus, the replication package would have:

```
1  ...
2  code/main.do
3  README.md
4  include/confparms_template.do
5  releasable/figure1.dta
6  releasable/figure1.pdf
```

# Keeping on top of provenance

- Licenses
- Streamlining for reproducibility

# Licenses

# Where does the file come from?

- How can we describe this later to somebody?
  - Point and click is long to describe
  - What are the rights we have?

# What is a license?

A license (licence) is an official permission or permit to do, use, or own something (as well as the document of that permission or permit).[11] [12]

# Examples

- **Creative Commons licenses**, used for artistic products and data

- **Open Source licenses** (BSD, GPL, MIT, etc.), used for software (code)

# License applying to Geodist data

- CEPII GeoDist is under an "Etalab 2.0 license"

# Can we re-publish the file?

# Downloading via code

# Easiest:

**Stata**

```
1  use "$URL" , clear
```

# Why not?

- will it be there in two months? in 6 years?

- what if the internet connection is down?

# Easy:

## Stata

```
1  global URL "https://www.cepii.fr/distance/dist_cepii.dta"
2  copy "$URL" (outputfile), replace
```

## R

```
1  download.file(url="$URL",destfile="(outputfile)")
```

We will get to even better methods a bit later

# Creating a README

- Template README
  - Cite both dataset and working paper
  - Add data URL and time accessed (can you think of a way to automate this?)
  - Add a link to license (also: download and store the license)

# Link

Step 1: Stata, R [13]

# Wrapping it all up

# Wrapping up

- Public replication package contains intelligible code, omits confidential details (but provides template code), has detailed data provenance statements

- Confidential replication package contains all the same, plus the confidential code, is archived in the FSRDC

# Things to remember

- Use code to save figures and tables (`estout`, `graph export`, `regsave`)

- Create log files for each run (`stata -b do file.do` not fine-grained enough) link

# Things to remember

Run it all again, top to bottom!

# Things to remember

- When doing a disclosure review request, remember to request the **code**

- When outputting statistics, **consider the disclosure rules** - the less changes, the faster the output (in theory), but in particular fewer surprises

- Do not think "**nobody will ever read this code**" - somebody is very likely to!

# End

Now you wait for the replicators to show up!

# Footnotes

1.

Bollen et al. 2015. "Social, Behavioral, and Economic Sciences Perspectives on F
and Reliable Science." National Science Foundation.
https://www.nsf.gov/sbe/AC_Materials/SBE_Robust_and_Reliable_Research_Rep

2.

Ambrus, Attila, Erica Field, and Robert Gonzalez. 2020. "Loss in the Time of Cho
Long-Run Impact of a Disease Epidemic on the Urban Landscape." American Ec
Review, 110 (2): 475–525. https://doi.org/10.1257/aer.20190759

3.

Ambrus, Attila, Field, Erica, and Gonzalez, Robert. Data and Code for: Loss in th
of Cholera: Long-run Impact of a Disease Epidemic on the Urban Landscape. Na
TN: American Economic Association [publisher], 2020. Ann Arbor, MI: Inter-univ
Consortium for Political and Social Research [distributor], 2020-01-31.
https://doi.org/10.3886/E111523V2

4.