

Introduction to open-sourcing data

Matthew J. Salganik
Department of Sociology
Princeton University

Summer Institute in Computational Social Science
June 20, 2019

The Summer Institute in Computational Social Science is supported by grants from the Russell Sage Foundation and the Alfred P. Sloan Foundation.





https://www.youtube.com/watch?v=66oNv_DJuPc

Brief introduction into open-sourcing your data:

- ▶ Store your data in a simple format

Brief introduction into open-sourcing your data:

- ▶ Store your data in a simple format
- ▶ Provide documentation

Brief introduction into open-sourcing your data:

- ▶ Store your data in a simple format
- ▶ Provide documentation
- ▶ Beware of privacy

Store your data in a simple format

Store your data in a simple format

In this case .csv should be good.

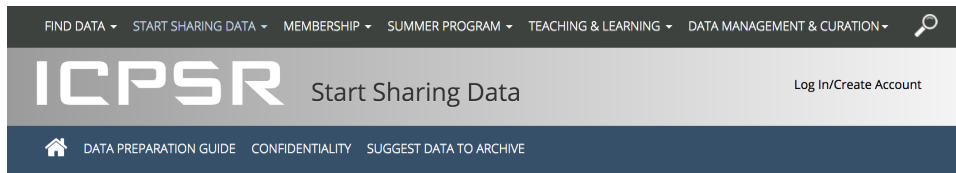
Provide documentation

Provide documentation

What would another researcher want to know?

- ▶ How and when was this data collected?
- ▶ What do the different variables describe?

Provide documentation (more details)



Data Preparation Guide

Introduction

1. Proposal Development and

Guide to Social Science Data Preparation and Archiving Phase 3: Data Collection and File Creation

Best Practices in Creating Metadata

<https://www.icpsr.umich.edu/icpsrweb/content/deposit/guide/chapter3docs.html>

Beware of privacy

Beware of privacy

Remove personally identifying information

NIST definition [\[edit \]](#)

The following data, often used for the express purpose of distinguishing individual identity, clearly classify as PII under the definition used by the [National Institute of Standards and Technology](#) (described in detail below):^[15]

- [Full name](#) (if not common)
- [Face](#) (sometimes)
- [Home address](#)
- [Email address](#) (if private from an association/club membership, etc.)
- [National identification number](#) (e.g., [Social Security number](#) in the U.S.)
- [Passport number](#)
- [Vehicle registration plate number](#)
- [Driver's license number](#)
- [Face, fingerprints, or handwriting](#)
- [Credit card numbers](#)
- [Digital identity](#)
- [Date of birth](#)
- [Birthplace](#)
- [Genetic information](#)
- [Telephone number](#)
- [Login name](#), [screen name](#), [nickname](#), or [handle](#)

https://en.wikipedia.org/wiki/Personally_identifiable_information

Privacy and Security

Myths and Fallacies of “Personally Identifiable Information”

<http://dx.doi.org/10.1145/1743546.1743558>

In this case, we recommend:

- ▶ Removing PII (name, email address, etc)
- ▶ Removing TurkID
- ▶ Coarsen age, geography, and race/ethnicity
- ▶ Coarsen timestamp
- ▶ Anything else?

For more on coarsening, see this code:

https://github.com/compsocialscience/summer-institute/blob/master/2018/materials/day4-surveys/mturk_data_cleaning.Rmd

For more about de-identification, see *Bit by Bit*, Sec 6.6.2 “Understanding and managing informational risk”

The 5Ws of data release:

The 5Ws of data release:

- ▶ Who: You

The 5Ws of data release:

- ▶ Who: You
- ▶ What: make your data available to other researchers in a responsible way

The 5Ws of data release:

- ▶ Who: You
- ▶ What: make your data available to other researchers in a responsible way
- ▶ Where: Dataverse, ICPSR, or an archival data repository

The 5Ws of data release:

- ▶ Who: You
- ▶ What: make your data available to other researchers in a responsible way
- ▶ Where: Dataverse, ICPSR, or an archival data repository (your laptop is not an archival data repository)

The 5Ws of data release:

- ▶ Who: You
- ▶ What: make your data available to other researchers in a responsible way
- ▶ Where: Dataverse, ICPSR, or an archival data repository (your laptop is not an archival data repository)
- ▶ When: when you publish your paper

The 5Ws of data release:

- ▶ Who: You
- ▶ What: make your data available to other researchers in a responsible way
- ▶ Where: Dataverse, ICPSR, or an archival data repository (your laptop is not an archival data repository)
- ▶ When: when you publish your paper
- ▶ Why: It is good for you and it is good for the world

In this case, you should archive your data here:

<https://github.com/compsocialscience/summer-institute/tree/master/2018/materials/day4-surveys/datasets>

When you start your projects next week

- ▶ plan to release your data
- ▶ plan to release your code

Questions