# Teaching Reproducibility and Replicability in Spatial Data Science

Where to Start and What to Do

Peter Kedron        University of California Santa Barbara
Joseph Holler        Middlebury College
Andrew Trgovac    Arizona State University

HEGSRR.github.io

**Dr. Peter Kedron**

Associate Professor
Department of Geography
UC Santa Barbara

Associate Director
Center for Spatial Studies & Data Science

**Dr. Joseph Holler**

Associate Professor
Department of Geography
Middlebury College

**Dr. Andrew Trgovac**

Teaching Assistant Professor
School of Geographical Science and
Urban Planning
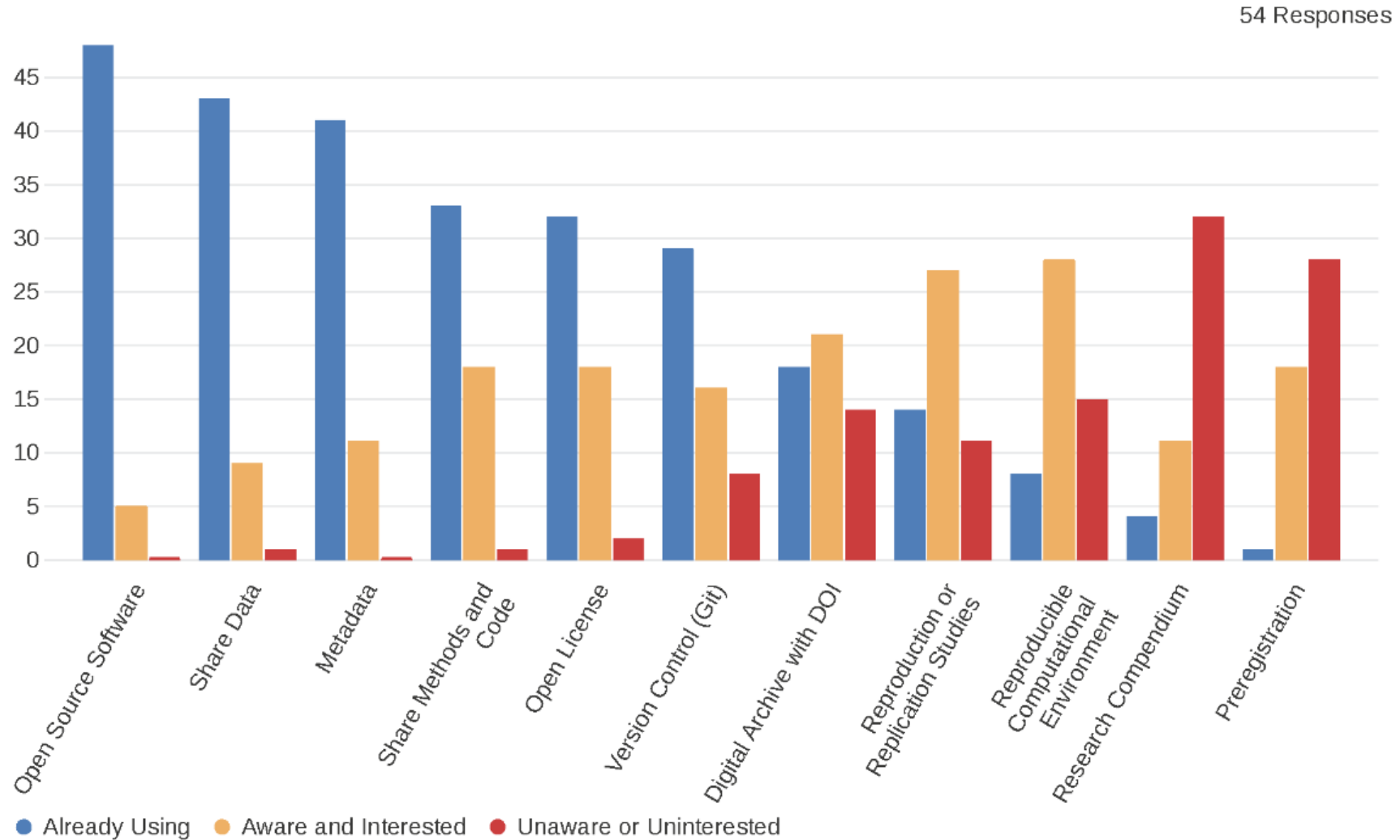Arizona State University

# Workshop Agenda

Introduce a competency rubric and a classroom model for R&R education in spatial data science

**Table 1.** Workshop Agenda and Presenters

| Time | Topic |
|---|---|
| 10 min | Introduction and Major Takeaways |
| 10 min | Competency Rubric |
| 10 min | Q&A |
| 05 min | Classroom Model |
| 05 min | Project Based Learning |
| 20 min | Implementing the Classroom Model |
| 10 min | Q&A |
| 05 min | Resources and Next Steps |

# Open Science Practices

More than sharing data and code



54 Responses

Legend: Already Using · Aware and Interested · Unaware or Uninterested

# A Working Definition

(Schmidt 2009, Gomez et al. 2010, Barba 2017, Christensen et al. 2019, NASEM 2019)

|  | Purpose | Data | Context | Procedure | Results |
|---|---|---|---|---|---|
| Reproduction | Internal Validity | Same | Same | Same | Same |
| Reanalysis | Internal Validity | Same | Same | Different | Similar |
| Replication | External Validity | Different | — | Similar | Similar |

# A Working Definition

**Veridical Spatial Data Science**

Principled inquiry to extract reliable and reproducible information from spatialtemporal data, with an ***enriched technical language*** to communicate and evaluate empirical evidence in the context of human decisions, domain knowledge, and geographic confounds; ***supported by a system of external validation and evidence accumulation based on the purposeful replication of findings across space and time***.

(Adapted from Kedron and Bardin 2021, Yu and Kumbier 2020)

# Key Ideas for Teaching Reproducibility and Replicability

What to take away from this workshop

**(1) Teach more than data and code sharing.**
Link R&R to the epistemology of science and open science

**(2) Attempt reproductions and replications with your students**
Pedagogically rich way to teach GIScience, spatial data science, and topical knowledge

**(3) Reproduction attempts create intrinsic and extrinsic rewards**
Reproduction attempts are rewarding for students, improve learning outcomes, produce publications

**(4) Use (and improve) our open educational materials**
Templates, past reproductions, teaching materials

# Reproducibility in Spatial Science

Competencies and Metrics for Curriculum Development and Learner Assessment

# Inspiration For This Work

Competencies of reproducible spatial data science

**2024**

O'Donnell, K. L., Aiello-Lammens, M., Bledsoe, E., Bowlick, F. J., Broughton, L., Calderon, O., Crispo, E., Emery, N., Farrell, K., Ngiramahoro, M., Patel, N., Paudel, S., Richardson, L., Soares, B. E., Supp, S., Weigel, E.

BEDE Network Data Science Skills Curriculum Map. Biological and Environmental Data Education (BEDE) Network

QUBES Educational Resources

https://qubeshub.org/publications/4859/alignments/1

# Structure of the Reproducibility Competency Rubric

We identify four competency categories and three levels of achievement

| Category | Skill | Foundational | Second Level | Third Level |
|---|---|---|---|---|
| **Science Context** | | | | |
| **Provenance** | | | | |
| **Project Organization and Sharing** | | | | |
| **Reproducible Code** | | | | |

# Skills and Student Learning Outcomes

What to take away from this workshop

| Category | Skill | Foundational | Second Level | Third Level |
|---|---|---|---|---|
| **Science Context** | 1. Evaluating prior research<br>2. Acknowledgement<br>3. Intellectual Property | | | |
| **Provenance** | 1. Creating Data<br>2. Using Data<br>3. Version Control<br>4. Documenting Versions | | | |
| **Project Organization and Sharing** | 1. Sharing<br>2. Storage<br>3. Detailing<br>4. Balancing Reproducibility with Compliance | | | |
| **Reproducible Code** | 1. Testing Code<br>2. Coding notebooks<br>3. Commenting<br>4. Documenting<br>5. Computational Environment | | | |

# Pedagogical Foundation

A systematic approach grounded in revised Bloom's Taxonomy

## Revised Bloom's

**Cognitive Process**     **Knowledge Types**

Remember                    Factual
Understand                  Conceptual
Apply                       Procedural
Analyze                     Metacognitive
Evaluate
Create

# Pedagogical Foundation

A systematic approach grounded in revised Bloom's Taxonomy

**Revised Bloom's**

| **Cognitive Process** | **Knowledge Types** | **Example** |
|---|---|---|
| Remember | Factual | **Category:** Provenance |
| Understand | Conceptual | **Skill:** Creating Data |
| Apply | Procedural | **Outcome:** Understand spatial data structures |
| Analyze | Metacognitive | |
| Evaluate | | **Blooms:** Understand, Conceptual |
| Create | | |

# Detailing

Specific learning outcomes associated with detailing project-level information

| Category | Skill | Foundational | Second Level | Third Level |
|---|---|---|---|---|
| **Science Context** | 1. **Evaluating research**<br>2. Acknowledgement<br>3. Intellectual Property | (F1) identify the research question and the type of research study (e.g., observational/experimental)<br><br>(F2) identify claims/conclusions of a study<br><br>(F3) deconstruct research workflow | (N1) evaluate internal validity of a study<br><br>(N2) critique a research workflow<br><br>(N3) plan to reproduce / replicate research | (T1) design solutions to address validity concerns<br><br>(T2) evaluate context of study in literature<br><br>(T3) attempt to replicate research in new context (e.g., location, population) |
| Bloom's Taxonomy | | (F1) Analyze, Conceptual<br>(F2) Analyze, Conceptual<br>(F3) Evaluate, Procedural | (S1) Evaluate, Conceptual<br>(S2) Evaluate, Procedural<br>(S3) Create, Metacognitive | (T1) Create, Metacognitive<br>(T2) Evaluate, Conceptuals<br>(T3) Apply, Procedural |

# Version Control

Specific learning outcomes associated with the version control skill

| Category | Skill | Foundational | Second Level | Third Level |
|---|---|---|---|---|
| **Provenance** | 1. Creating Data<br>2. Using Data<br>3. **Version Control**<br>4. Documenting Versions | (F1) understand what version control is<br><br>(F2) create a local record of changes for your data and/or code | (S1) identify if application has version control<br><br>(S2) use version control software tool (e.g., GitHub) individually to track workflows (e.g., push-pull or commit-add) | (T1) create a version control process<br><br>(T2) collaboratively use version control (e.g., use development branches or workflows that require pull requests) |
| Bloom's Taxonomy | | (F1) Understand, Factual<br>(F2) Apply, Procedural | (S1) Analyze, Conceptual<br>(S2) Apply, Procedural | (T1) Create, Procedural<br>(T2) Apply, Metacognitive |

# Detailing

Specific learning outcomes associated with detailing project-level information

| Category | Skill | Foundational | Second Level | Third Level |
|---|---|---|---|---|
| **Project Organization and Sharing** | 1. Sharing<br>2. Storage<br>**3. Detailing**<br>4. Balance Reproducibility with Compliance | (F1) sharing documented details about data and/or code, methods, or variables (metadata)<br><br>(F2) organized project directory structure | (N1) write a README file with project-level metadata and geographic extent<br><br>(N2) index project contents | (T1) Document data with ISO-standard geospatial metadata<br><br>(T2) Archive project with Dublin core metadata and detailed geographic extent |
| Bloom's Taxonomy | | (F1) Apply, Procedural<br>(F2) Apply, Procedural | (S1) Create, Factual<br>(S2) Analyze, Procedural | (T1) Create, Factual<br>(T2) Create, Factual |

# Teaching Reproducibility and Replicability

Doing Reproduction and Replications with Students

# Summary Introduction

How we reproduce studies with students: settings

- Advanced undergraduate Open GIScience
- Graduate spatial statistics / GIS
- Independent research / thesis
- Summer research assistants
- Special reading/study groups

# Summary Introduction

How we reproduce studies with students: process

Dissect → Plan → Register → Attempt → Compare → Share

# Summary Introduction

How we reproduce studies with students: toolbox

- Complete Reproduction Studies
- Template Research Compendium
- Template Analysis Plan
- Course websites
- Suggested studies
- Project-based learning

# Project Based Learning

Our pedagogical foundation

| Project Based Learning Component | Interpretation | Example Linkage to Reproductions and Replications |
|---|---|---|
| 1. Challenging Problem | *Meaningful question at appropriate level* | Critically evaluating the design and execution of prior studies |
| 2. Sustained Inquiry | *Extended process of questioning* | Iterative evaluation and revision of work |
| 3. Authenticity | *Real world context* | Engagement with peer-reviewed empirical research |
| 4. Student Voice & Choice | *Decisions & implementation done by students* | Students contribute to project selection and lead project design/revision |
| 5. Critique & Revision | *Give, receive & apply feedback* | Student to student consensus for replication design |
| 6. Reflection | *Active reflection on learning and effectiveness of decisions* | Collective assessment of design decisions, analyses & unexpected challenges |
| 7. Public Product | *Work shared beyond the classroom* | Share compendium of project design, data, code, results etc. for key decisions |

# Classroom Implementation I

Deconstruct a published study

Chakraborty, J. 2021.
Social inequities in the distribution of
COVID-19: An intra-categorical analysis of
people with disabilities in the U.S.
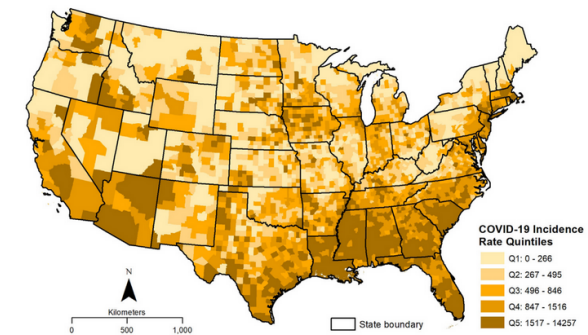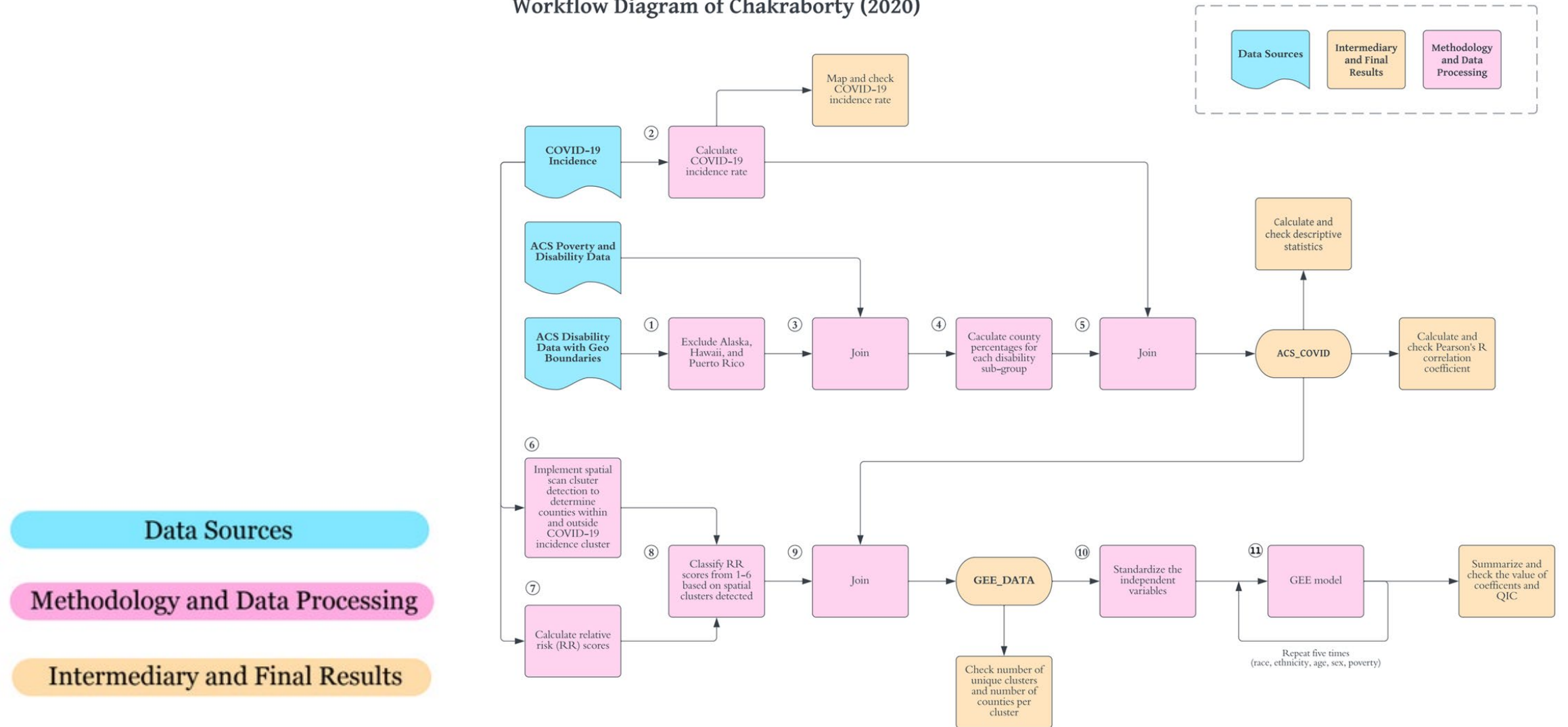*Disability and Health Journal* 14:1-5.
https://doi.org/10.1016/j.dhjo.2020.101007





Fig. 1. County level distribution of COVID-19 incidence rate (cases per 100,000 people) in the continental USA, August 1, 2020.



Data Sources

Methodology and Data Processing

Intermediary and Final Results

# Classroom Implementation I

Deconstruct a published study



Workflow Diagram of Chakraborty (2020)

# Classroom Implementation II

Plan a reproduction attempt with a research compendium template

1. Directory structure
2. Version Control
3. Metadata:
   - Study
   - Data Sources

# Classroom Implementation II

Plan a reproduction attempt with a research compendium template



I. Analysis Plan in Computational Notebook
⇒ PDF Preregistration

# Classroom Implementation III

Attempt a reproduction, addressing unexpected issues, and assessing outcomes

1. Analysis Plan in Computational Notebook
     ⇨ PDF Preregistration
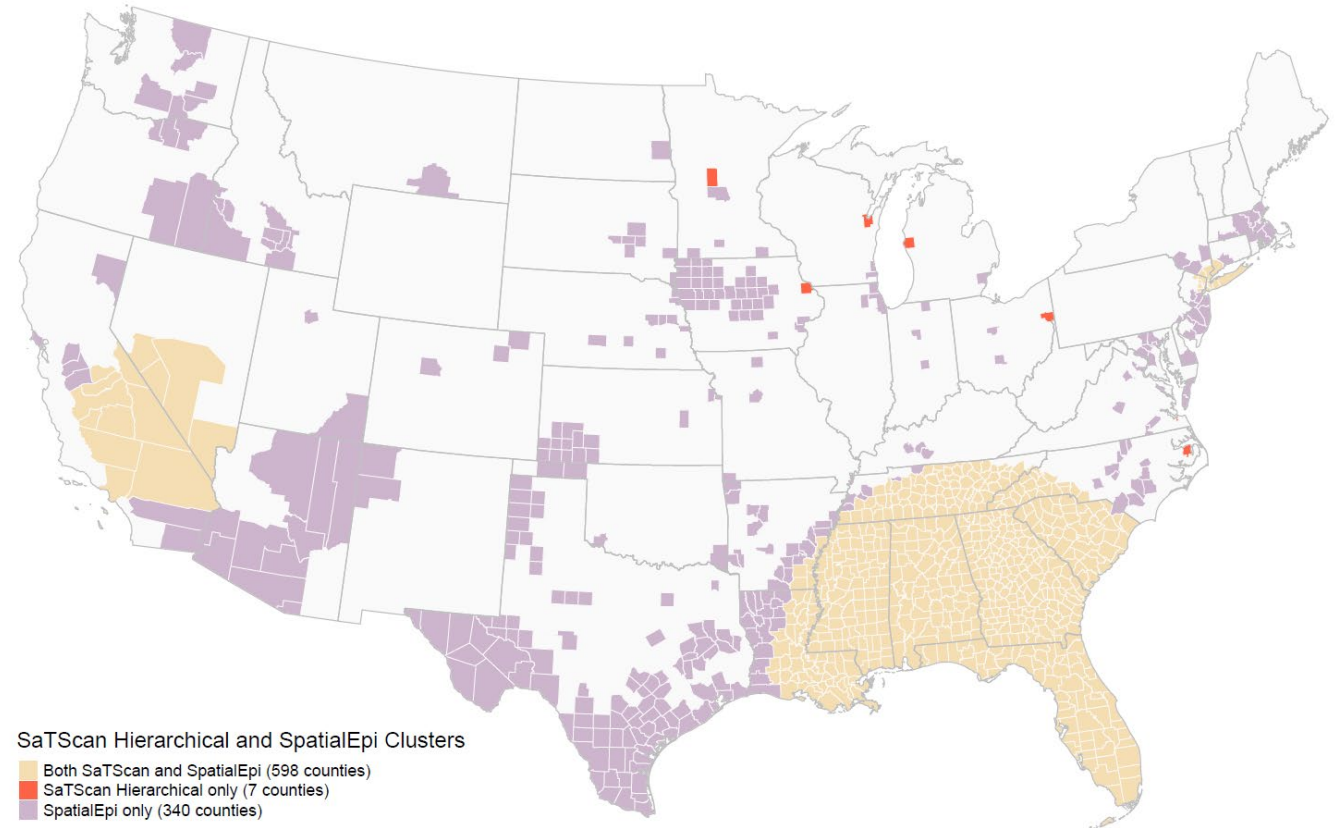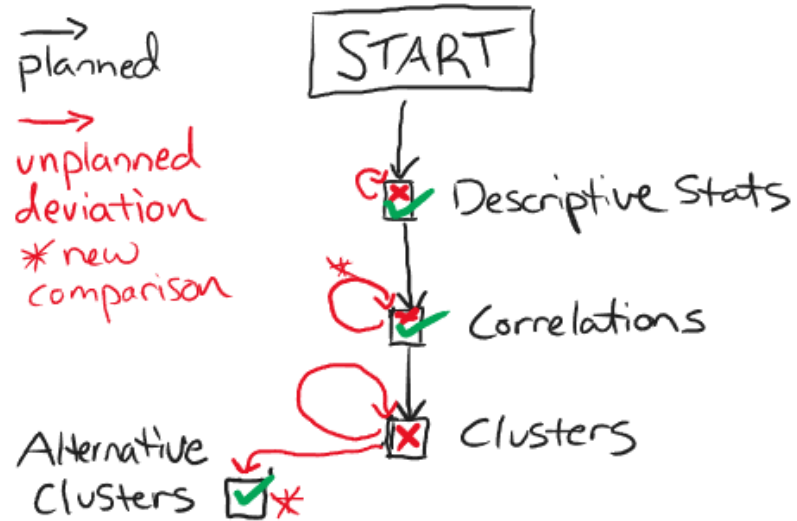
2. Attempt Reproduction
     ✚ code blocks
     ✚ check results
     ✚ unplanned deviations

```
178   American Community Survey (ACS) data for sociodemographic subcategories of people with
      disabilities can be accessed by using the `tidycensus` package to query the Census API. This
      requires an API key which can be acquired at
      [api.census.gov/data/key_signup.html](https://api.census.gov/data/key_signup.html).
179
180 ▾  ```{r API-Load-ACS, eval=FALSE}                                                    ⚙ ⊻ ▸
181   # If you wish to use a census API key, run the census_api_key() function in the console
182
183   # Query disability demographic data with geographic boundaries
184   acs <- get_acs(
185     geography = "county",
186     table = "S1810",
187     year = 2018,
188     output = "wide",
189     cache_table = TRUE,
190     geometry = TRUE,
191     keep_geo_vars = TRUE
192   )
193
```

*adding tidycensus code to query ACS data to the data sources section*
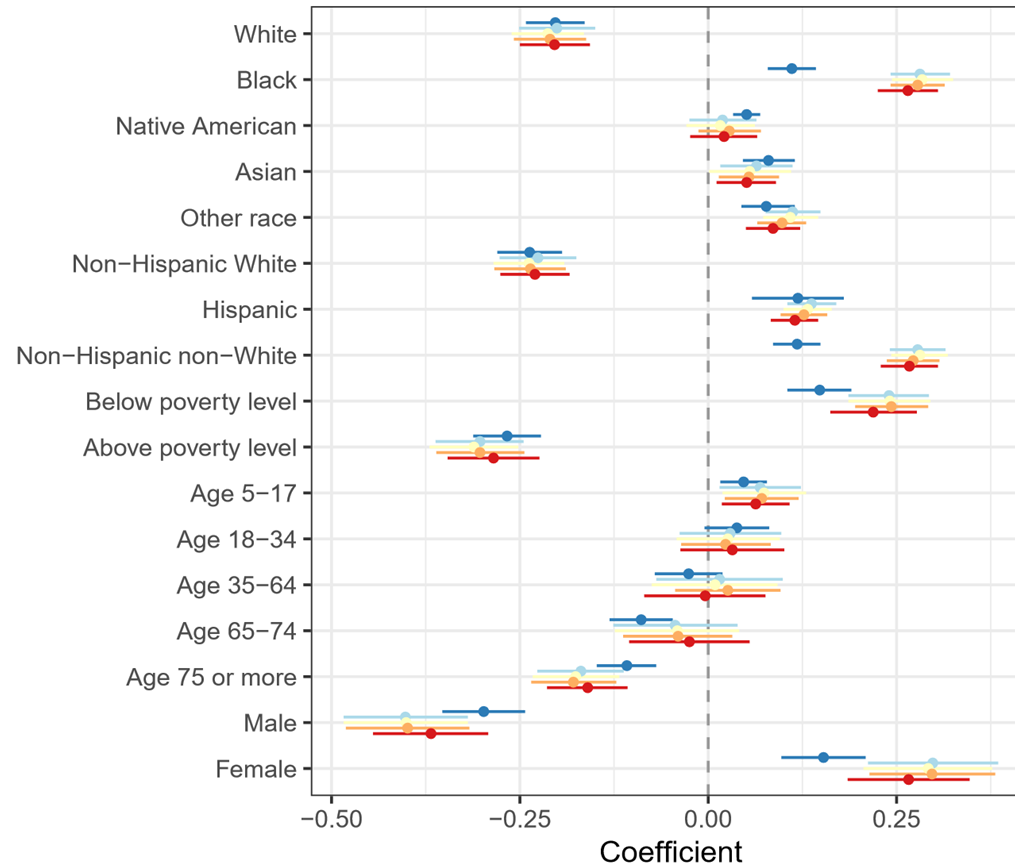
# Classroom Implementation III

Attempting a reproduction, addressing unexpected issues, and assessing outcomes

# Classroom Implementation III

Attempting a reproduction, addressing unexpected issues, and assessing outcomes

# Classroom Implementation IV

Create and share a reproducible, public report

1. Analysis Plan in Computational Notebook
   → PDF Preregistration

2. Attempt Reproduction
   + code blocks
   + check results
   + unplanned deviations

3. Analysis Report
   + output tables + figures
   + discuss unplanned deviations
   + conclusions
      → PDF Registered Report
      → HTML Website

# Outcomes and Resources

Eleven studies and counting…

| Original Study | Level | Pre-analysis | Compendium | Report | Publication |
|---|---|---|---|---|---|
| Malcomb et al 2014 | Undergraduate | GitHub | GitHub; OSF | — | —- |
| Kang et al 2020 | Undergraduate | GitHub | GitHub; OSF | OSF | *International Journal of Health Geographics* - In Review |
| Mollalo et al 2020 | Graduate | GitHub | GitHub; OSF | GitHub | *Geographical Analysis* |
| Saffary et al 2020 | Graduate | GitHub | GitHub; OSF | GitHub | *Geographical Analysis* |
| Vijayan et al 2021 | Graduate | GitHub | GitHub; OSF | GitHub | *Geographical Analysis* |
| Chakrabory 2021 | Undergraduate | OSF | GitHub; OSF | OSF | |
| DiMaggio et al 2020 | Graduate | OSF | GitHub; OSF | OSF | *Annals of Epidemiology* |
| Speilman et al 2020 | Undergraduate | OSF | GitHub; OSF | OSF | |
| Maldonado | Undergraduate | Github - Private | GitHub - Private | GitHub - Private | *Journal of Immigrant and Minority Health* - Submit Fall 2024 |
| Kodros | Graduate | Github - Private | GitHub - Private | GitHub - Private | |
| Brodie | Graduate | Github - Private | GitHub - Private; OSF - Private | Github - Private | *Nature* - Submit Fall 2024 |

# Outcomes and Resources



*Workshop Website*

[hegsrr.github.io/Workshop-SDSS-2024/](hegsrr.github.io/Workshop-SDSS-2024/)

- 5 Peer-reviewed Publications

- 11 Reproduction and Replication Studies

- 2 Surveys of Researcher Practices with interactive data visualizations

- **Research Compendium Template**

- Manual In Development

- Course Syllabi

- 9 RAs Mentored

- ~75 Students Engaged in R&R Studies

# Key Ideas for Teaching Reproducibility and Replicability

What to take away from this workshop

## (1) Teach more than data and code sharing.
Link R&R to the epistemology of science and open science

## (2) Attempt reproductions and replications with your students
Pedagogically rich way to teach GIScience, spatial data science, and topical knowledge

## (3) Reproduction attempts create intrinsic and extrinsic rewards
Reproduction attempts are rewarding for students, improve learning outcomes, produce publications

## (4) Use (and improve) our open educational materials
Templates, past reproductions, teaching materials

*Workshop Website*

*hegsrr.github.io/ Workshop-SDSS- 2024/*

# An Open Invitation

Please reach out to us.
We want to collaborate on R&R research.
We want to help those who are interested adopt this approach

2024 SDSS Workshop on Teaching R&R