

# Challenges to collaboration and reproducibility

Nathanael Rosenheim, PhD

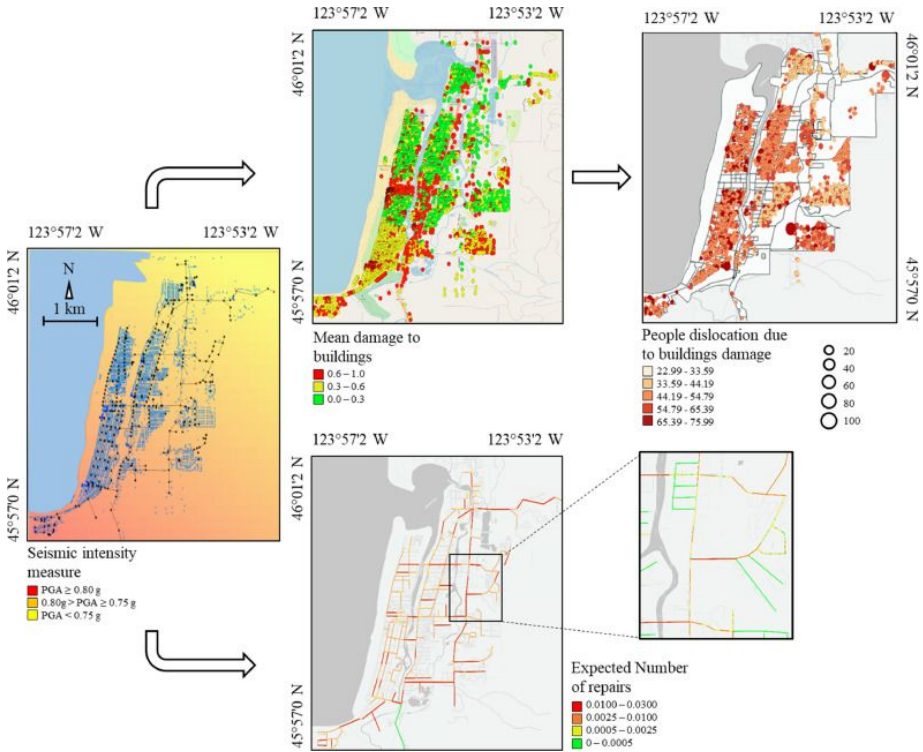
URSC 689

Feb 1, 2021

# Challenges Collaborating Across Universities

Community resilience model that combines work completed by Urban Planning researchers at TAMU and Civil Engineers at the University of Illinois -Urbana-Champaign.

All data shared via email. Models were run independently without shared code or version control.



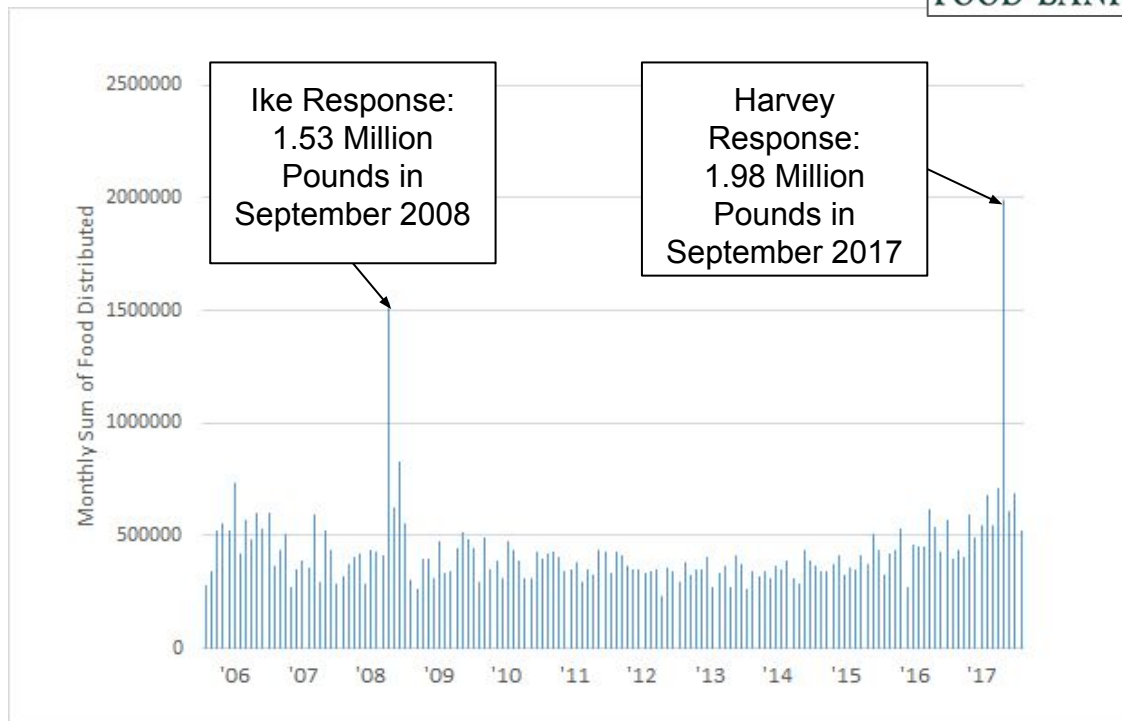


# Challenges Collaborating with Communities

Models of post disaster food aid distributions.

Community partner had a massive SQL database, but only way to access information was with a point-and-click user interface.

Took graduate student 40 hours to download data.



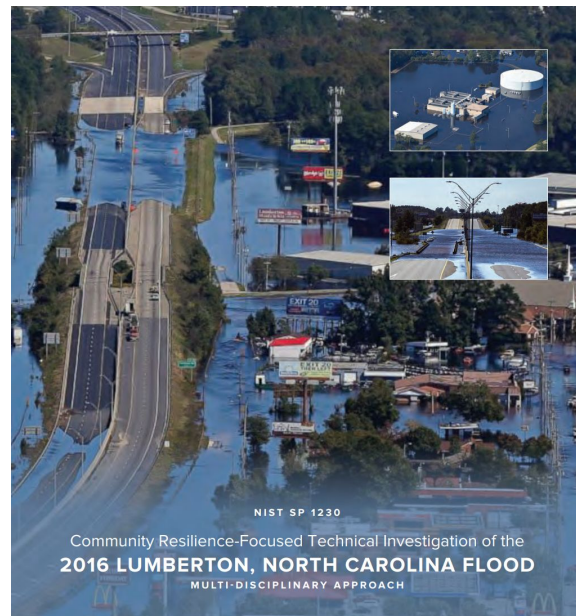
Source Data: Southeast Texas Food Bank Primarius Reports

# Challenges Collaborating with Federal Agencies

Post disaster field studies that combine engineering damage assessments with social science household and business level surveys.

Multi-year collaboration between 11 universities through the National Institute for Standards and Technology (NIST).

Diverse range of data collection, required Institutional Review Board (IRB) approval. Requires storage of confidential data with limited means of sharing and citing data.



#### EDITORS

John W. van de Lindt  
Walter Gillis Peacock  
Judith Mitroni-Reiser

This publication is available free of charge from:  
<https://doi.org/10.6028/NIST.SP.1230>

**NIST**  
National Institute of  
Standards and Technology  
U.S. Department of Commerce

# Challenges Collaborating with yourself

The most important collaborator is yourself...

Make it easy for your future self to find files, programs, data, and to continue the work you started.

# What is data?

Information obtained by scientific work and used for analysis.<sup>1</sup>

- Tabular Data
  - Survey responses
  - Administrative Data
- Metadata - Codebooks
- Relational Databases

	storeid	Q2_1
1	12	2. Manager
2	16	5. Other
3	41	5. Other
4	71	.
5	104	5. Other
6	123	2. Manager
7	125	5. Other
8	153	2. Manager
9	154	.
10	165	1. Owner
11	186	5. Other
12	202	5. Other
13	239	.
14	279	2. Manager
15	319	2. Manager
16	323	3. Owner and Manager
17	342	2. Manager
18	370	.
19	386	2. Manager
20	406	2. Manager
21	448	1. Owner
22	460	.
23	474	2. Manager

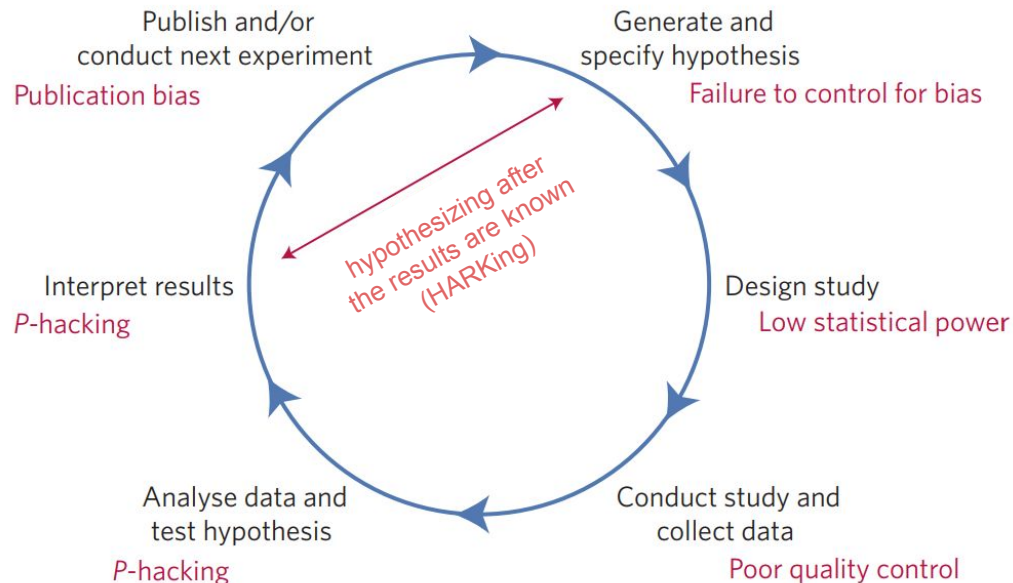
storeid	STOREID
<div><div>type: numeric (int)</div><div><div>range: [12,3605]</div><div>unique values: 135</div></div><div><div>mean: 1716.04</div><div>std. dev: 1061.55</div></div><div><div>percentiles:</div><div><div>10%279</div><div>25%832</div><div>50%1694</div><div>75%2564</div><div>90%3207</div></div></div></div>	
<div>storeid:<div><div>1. [SETX Survey Text] Store ID</div><div>2. Primary Key - unique ID randomly assigned when the sample frame was set.</div><div>3. Use STOREID to merge Coverage and Response Datasets.</div><div>4. [Citation] Rosenheim, N. et al 2018. Southeast Texas Food Retail Survey.</div><div>5. [Name of Saved Data File] RAPID17_1gv1_SNAP_SETX_RetailSurvey_2019-02-01/RAPID17_1gv1_SNAP_SETX_Reta &gt; ilSurvey_2019-02-01.dta</div><div>6. [Program to replicate Data File] RAPID17_1gv1_SNAP_SETX_RetailSurvey_2019-02-01.do</div><div>7. [Date data file was created] 1 Feb 2019 16:48:00</div></div></div>	
Q2_1	Question: 1
<div><div>type: numeric (byte)</div><div>label: Q2_1lbl_r</div><div><div>range: [1,5]</div><div>unique values: 5</div></div><div><div>units: 1</div><div>missing .: 31/135</div></div><div><div>tabulation:</div><div><div>Freq.</div><div>Numeric</div><div>Label</div></div><div><div>11</div><div>1</div><div>1. Owner</div></div><div><div>48</div><div>2</div><div>2. Manager</div></div><div><div>4</div><div>3</div><div>3. Owner and Manager</div></div><div><div>5</div><div>4</div><div>4. Assistant Manager</div></div><div><div>36</div><div>5</div><div>5. Other</div></div><div><div>31</div><div>.</div><div>Missing</div></div></div></div>	
<div>Q2_1:<div><div>1. [SETX Survey Text] 1. What is your role with this business? - Selected Choice</div><div>2. [Citation] Rosenheim, N. et al 2018. Southeast Texas Food Retail Survey.</div><div>3. [Name of Saved Data File] RAPID17_1gv1_SNAP_SETX_RetailSurvey_2019-02-01/RAPID17_1gv1_SNAP_SETX_Reta &gt; ilSurvey_2019-02-01.dta</div><div>4. [Program to replicate Data File] RAPID17_1gv1_SNAP_SETX_RetailSurvey_2019-02-01.do</div><div>5. [Date data file was created] 1 Feb 2019 16:48:00</div></div></div>	

Reference: 1. "data, n.". OED Online. December 2018. Oxford University Press.  
<http://www.oed.com.lib-ezproxy.tamu.edu:2048/view/Entry/296948?rskey=c3az3E&result=1> (accessed February 08, 2019).

# What is science?

“Science is an approach to knowledge... that strives to better approximate the state of nature by reducing errors in inferences.”<sup>1</sup>

“Conceptualize science is a toolbox of... tools designed to minimize mistakes [or bias].”<sup>1</sup>



Scientific method with potential threats of bias.<sup>2</sup>

Reference: 1. Lilienfeld, S. O., Sauvign , K. C., Lynn, S. J., Cautin, R. L., Latzman, R. D., & Waldman, I. D. (2015). Fifty psychological and psychiatric terms to avoid: a list of inaccurate, misleading, misused, ambiguous, and logically confused words and phrases. *Frontiers in Psychology*, 6, 1100.

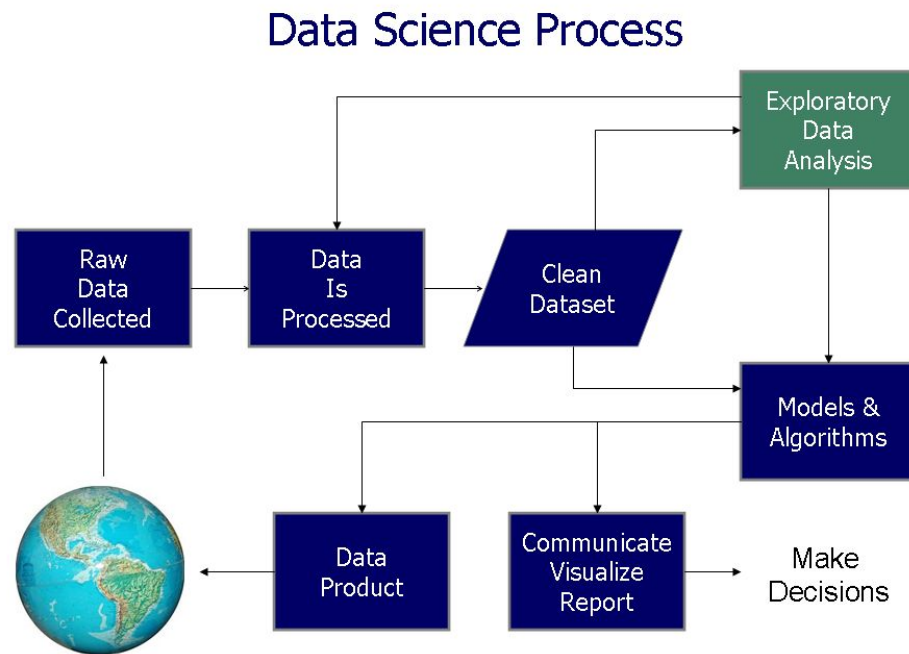
<https://doi.org/10.3389/fpsyg.2015.01100>

2. Munaf , M. R., Nosek, B. A., Bishop, D. V., Button, K. S., Chambers, C. D., du Sert, N. P., ... & Ioannidis, J. P. (2017). A manifesto for reproducible science. *Nature Human Behaviour*, 1, 0021. <https://doi.org/10.1038/s41562-016-0021>

# What is Data Science?

Data science is a set of tools designed to minimize bias associated with the analysis of data. “The discipline of turning raw data into understanding.”<sup>1</sup>

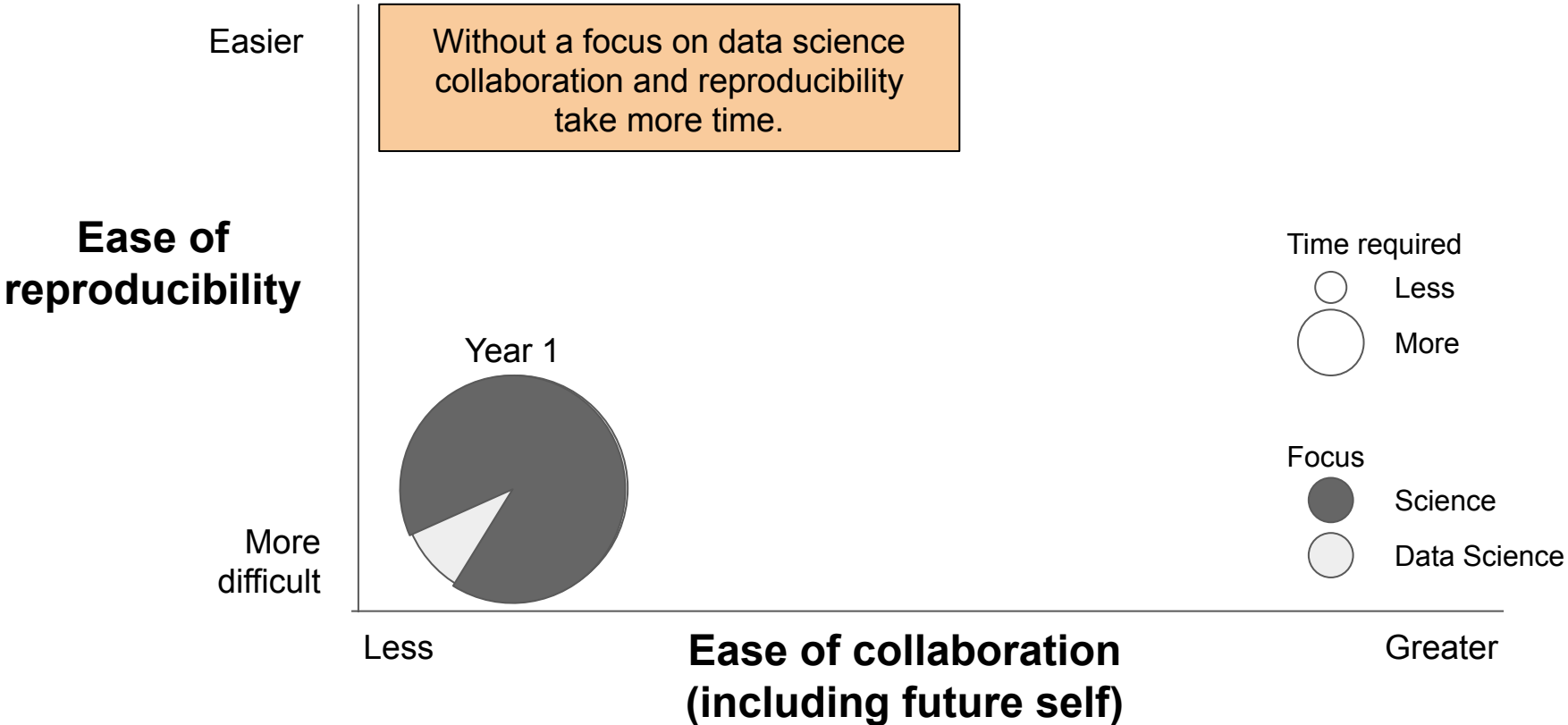
Example Tools/Concepts: Version Control, GitHub, Markdown [RMarkdown or Jupyter Notebook], Workflow, Repositories, Permanent Identifiers e.g. “handle” (hdl) or “digital object identifier” (doi)



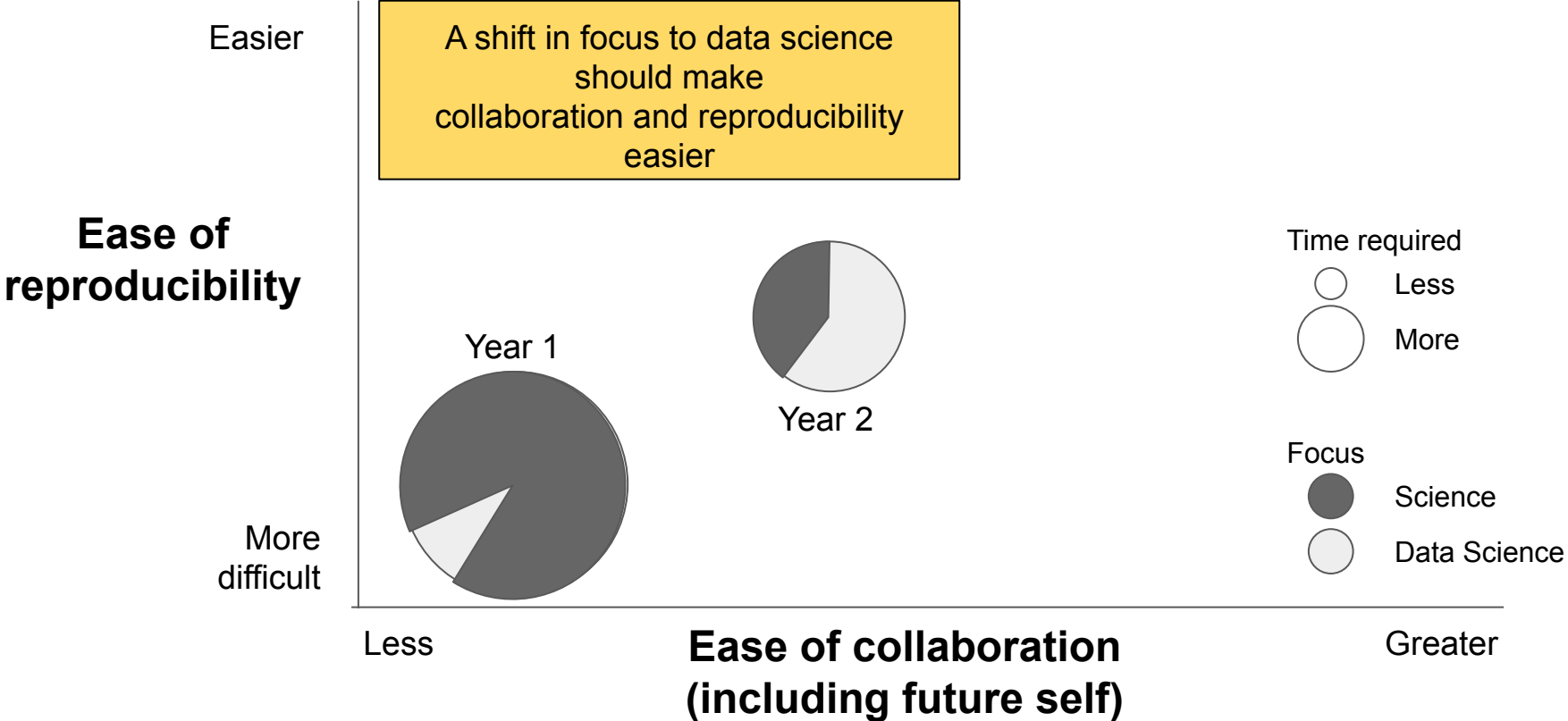
[2. Caldwell, J. \(2016\) A Data Science Solution to the Question "What is Data Science?" R-Bloggers](https://doi.org/10.1038/s41559-017-0160)



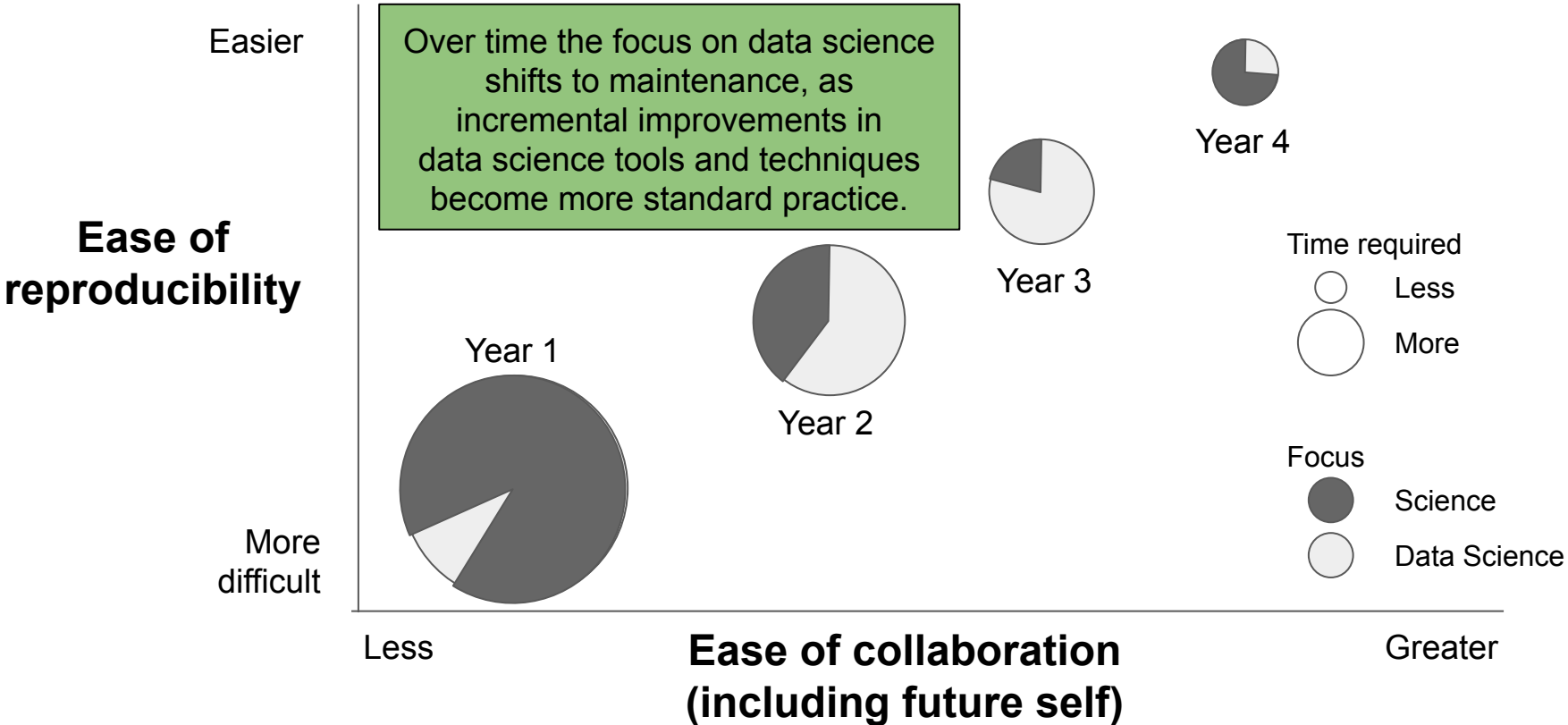
# Goal: Better science in less time



# Goal: Better science in less time



# Goal: Better science in less time



# Replication Standard - Individual or Social Contract?

Individual Responsibility	Social Contract
If asked a researcher should be able to provide the files to replicate published results.	Files to replicate published results are submitted at time of publication.
Emphasis on trust.	Emphasis on transparency.
Faith in the author.	Focus on openness.
Reinforcement of status.	Distributes power and access.

# Course objectives...

Help to overcome challenges to data sharing, documentation, publication, and analytics.

Help to develop a culture that supports a social contract for data replication.

To be a part of a community that bolsters data science and leads to more open, discoverable, reproducible research.

# Motivating Examples

Rosenheim, Nathanael; Day, Wayne; Seong, Kijin (2021) “Automated Neighborhood Characteristics for Community Resilience Planning.” DesignSafe-CI. <https://doi.org/10.17603/ds2-hj0p-bp40>.

Roy, Malini; Rosenheim, Nathanael (2021) “Longitudinal Social Vulnerability Data Exploration for Harris County Census Tracts.” DesignSafe-CI. <https://doi.org/10.17603/ds2-hn6r-dh03>.

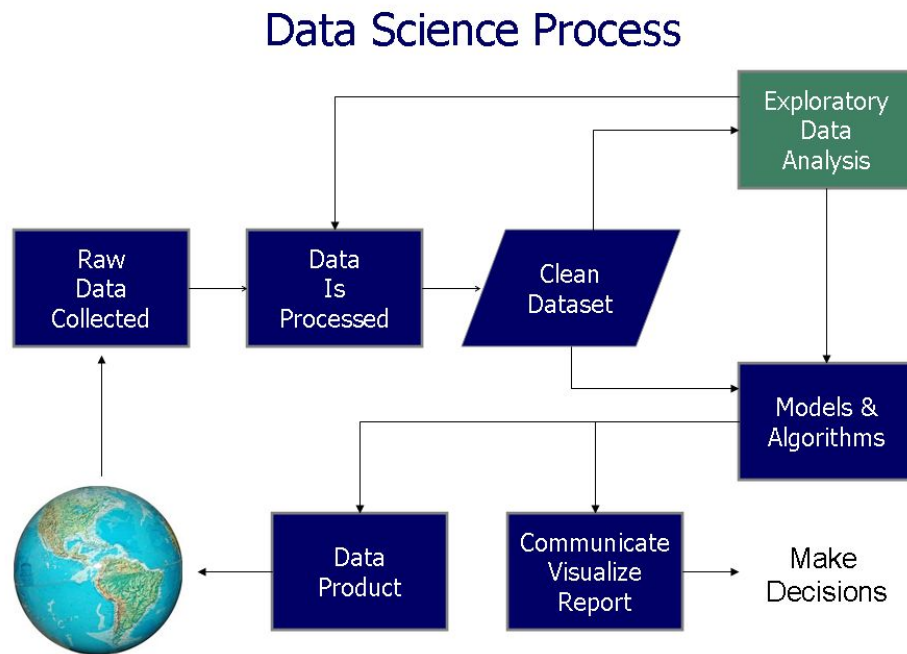
Gu, Donghwan, and Rosenheim, Nathanael. Demographic Analysis Workflow using Census API in Jupyter Notebook: 1990-2000 Population Size and Change. Ann Arbor, MI: Inter-university Consortium for Political and Social Research [distributor], 2020-07-30.  
<https://doi.org/10.3886/E120381V1>

Goodman, Cooper, Rosenheim, Nathanael, Day, Wayne, Gu, Donghwan, and Korukonda, Jayasaree. Population Distribution Workflow using Census API in Jupyter Notebook: Dynamic Map of Census Tracts in Boone County, KY, 2000. Ann Arbor, MI: Inter-university Consortium for Political and Social Research [distributor], 2020-07-31. <https://doi.org/10.3886/E120382V1>

# Motivating Example

Rosenheim, Nathanael; Day, Wayne;  
Seong, Kijin (2021) "Automated  
Neighborhood Characteristics for  
Community Resilience Planning."  
DesignSafe-CI.

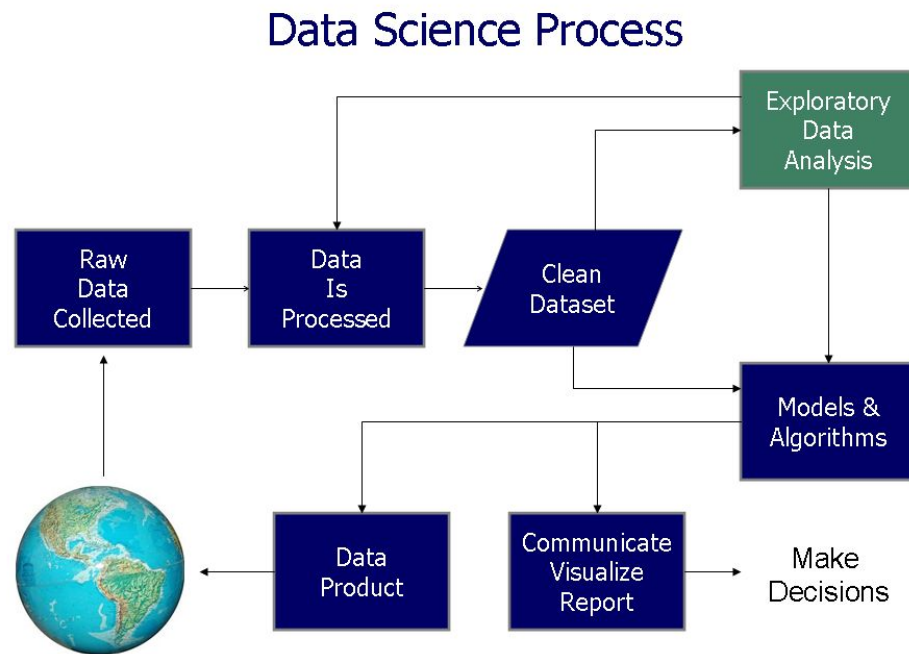
<https://doi.org/10.17603/ds2-hj0p-bp40>.



[2. Caldwell, J. \(2016\) A Data Science Solution to the Question "What is Data Science?" R-Bloggers](https://doi.org/10.1038/s41559-017-0160-2)

# Motivating Example

**Motivation:** Demonstrate the potential of Census API and Jupyter Notebooks for obtaining, cleaning, and exploring Census data.



Rosenheim, Nathanael; Day, Wayne; Seong, Kijin (2021) “Automated Neighborhood Characteristics for Community Resilience Planning.” DesignSafe-CI. <https://doi.org/10.17603/ds2-hj0p-bp40>.

[2. Caldwell, J. \(2016\) A Data Science Solution to the Question “What is Data Science?” R-Bloggers](https://doi.org/10.1038/s41559-017-0160-2)

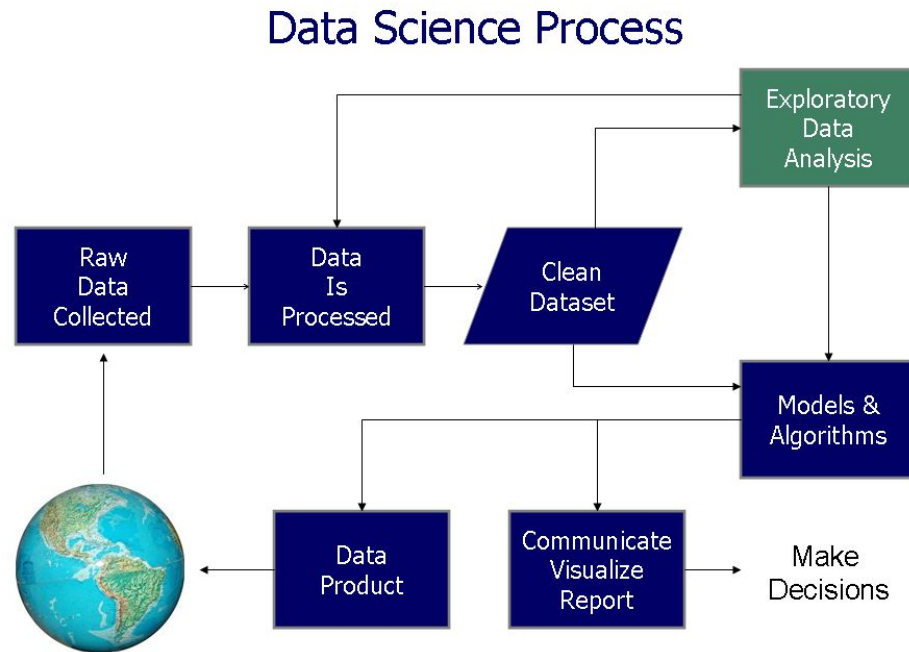


# Motivating Example

**Raw data collection:** 2010  
Decennial Census data using  
Census API and TIGER Shapefiles  
for Block Group Boundaries

Program automates data collection:

**IN-CORE\_CENSUSAPI\_DWNLD\_JOPLIN**  
**MSA\_2021-01-13.ipynb**



[2. Caldwell, J. \(2016\) A Data Science Solution to the Question "What is Data Science?" R-Bloggers](https://doi.org/10.1038/s41559-017-0160-2)

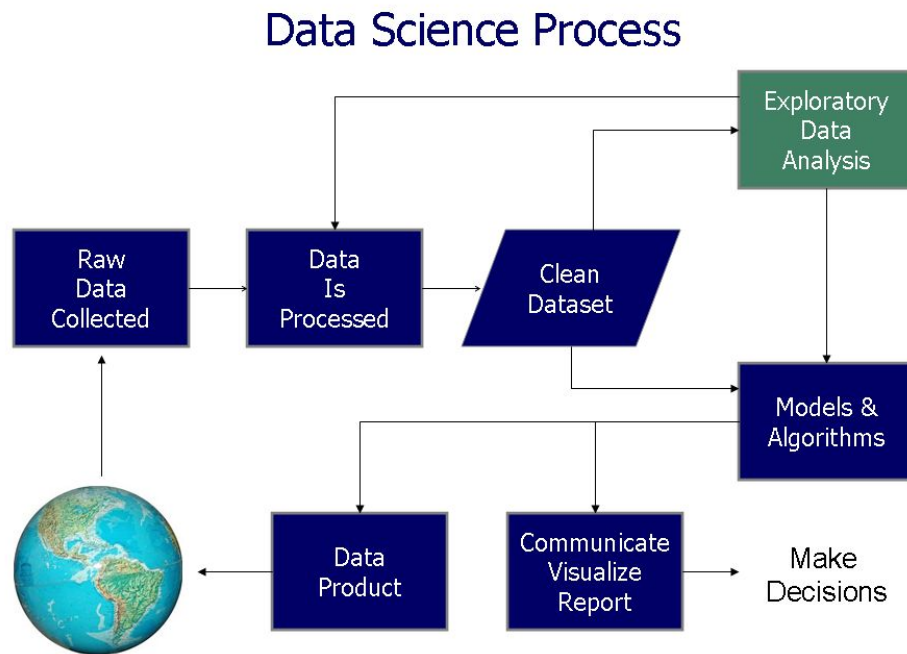
# Motivating Example

**Data is Processed:**

Jupyter Notebook

IN-CORE\_CENSUSAPI\_DWNLD\_JOPLIN

MSA\_2021-01-13.ipynb



[2. Caldwell, J. \(2016\) A Data Science Solution to the Question "What is Data Science?" R-Bloggers](https://doi.org/10.1038/s41559-017-0160-2)

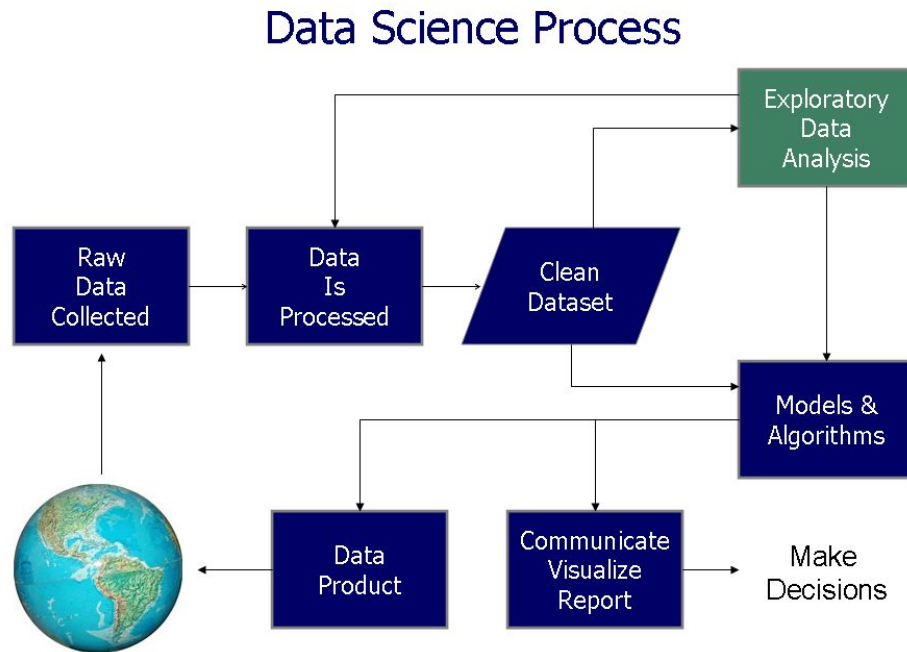
# Motivating Example

## Clean Dataset:

New variables are created and block group data is merged with shapefiles.

Program generates CSV, HTML, and Shp files

**IN-CORE\_CENSUSAPI\_DWNLD\_JOPLIN  
MSA\_2021-01-11.csv**



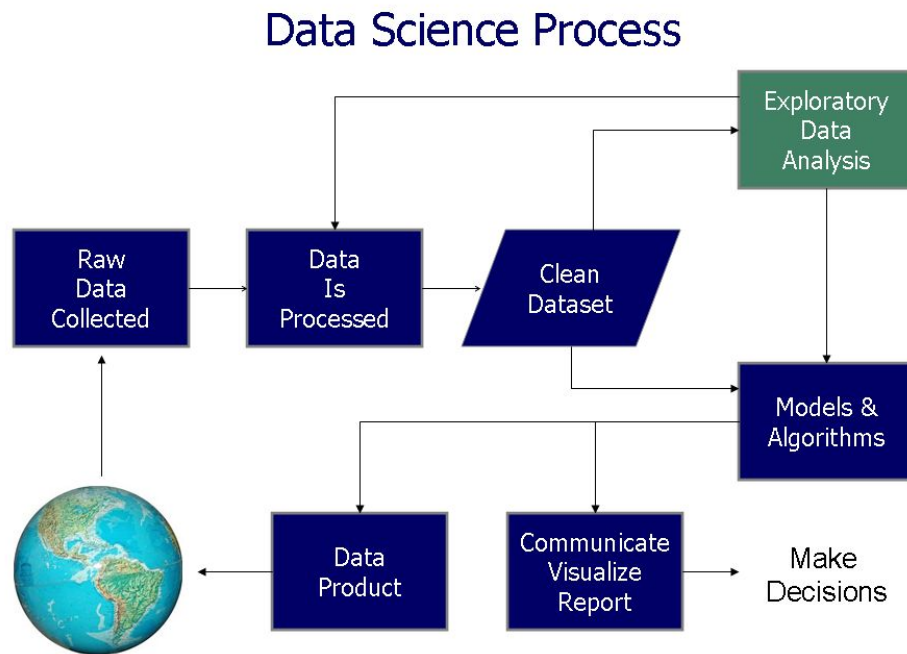
[2. Caldwell, J. \(2016\) A Data Science Solution to the Question "What is Data Science?" R-Bloggers](https://doi.org/10.1038/s41559-017-0160-2)

# Motivating Example

## Exploratory Data Analysis:

Jupyter Notebook includes  
descriptive statistics, tables, and  
exploratory maps

**IN-CORE\_CENSUSAPI\_DWNLD\_JOPLIN  
MSA\_2021-01-11\_map.html**



[2. Caldwell, J. \(2016\) A Data Science Solution to the Question "What is Data Science?" R-Bloggers](https://doi.org/10.1038/s41559-017-0160-2)

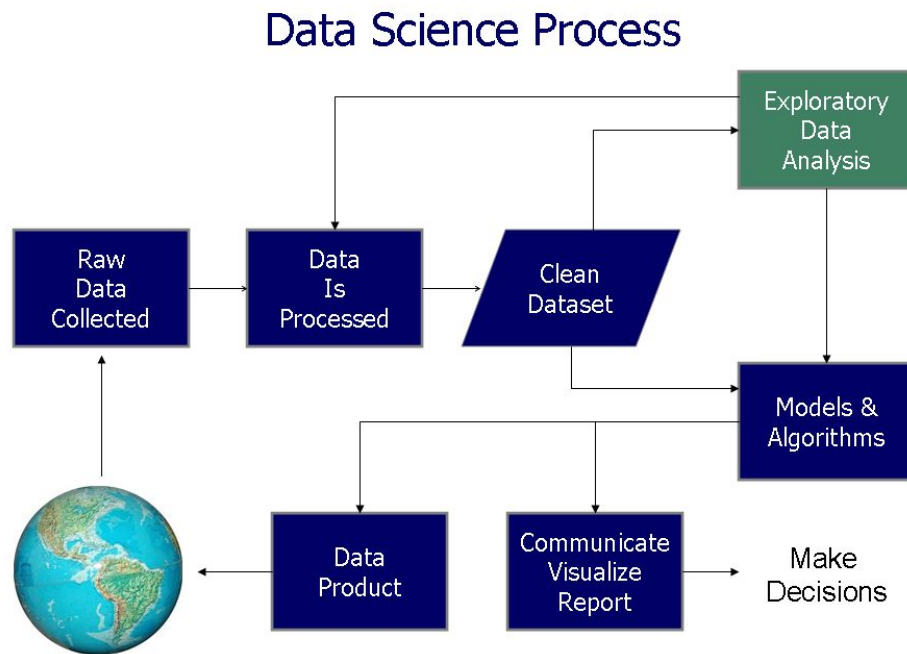
# Motivating Example

## Models and Algorithms:

Workflow does not include models

Algorithm saved as python program file:

**IN-CORE\_CENSUSAPI\_BGMAP\_2021-01-19T1403.py**



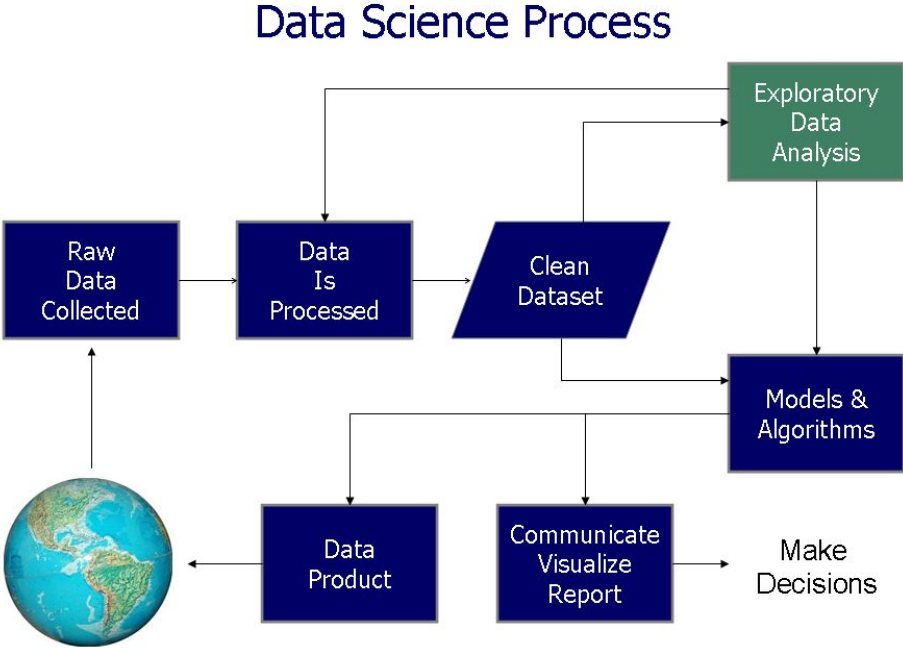
[2. Caldwell, J. \(2016\) A Data Science Solution to the Question "What is Data Science?" R-Bloggers](https://doi.org/10.1038/s41559-017-0160-2)

# Motivating Example

Communicate, Visualize, Report

`_ReadMeFirst_CensusAPI_BlockGroup_2`  
`021-01-19.docx`

Step by Step Instruction-How to  
Download Tables as a CSV  
`file_2021-01-07.pdf`



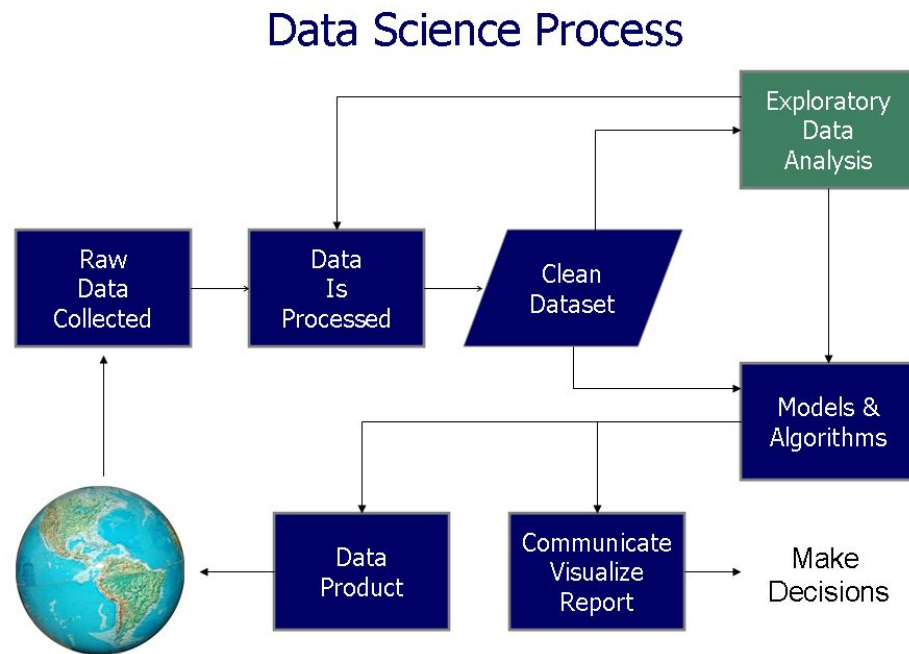
[2. Caldwell, J. \(2016\) A Data Science Solution to the Question "What is Data Science?" R-Bloggers](https://doi.org/10.1038/s41559-017-0160-2)

# Motivating Example

## Data Product:

Rosenheim, Nathanael; Day, Wayne;  
Seong, Kijin (2021) “Automated  
Neighborhood Characteristics for  
Community Resilience Planning.”  
DesignSafe-CI.

<https://doi.org/10.17603/ds2-hj0p-bp40>.



[2. Caldwell, J. \(2016\) A Data Science Solution to the Question "What is Data Science?" R-Bloggers](https://doi.org/10.1038/s41559-017-0160-2)

# What the project looks like in DesignSafe

The screenshot displays the DesignSafe-CL web interface. At the top, a teal header bar contains the NSF and NHERI logos. Below this, the 'DESIGNSAFE' logo is prominently displayed. A navigation bar includes links for 'Workspace', 'Learning Center', 'NHERI Facilities', 'NHERI Community', 'News', and 'Help'. A search bar is located on the right side of the navigation bar. The main content area is titled 'DATA DEPOT' and features a search bar and a list of actions: 'Rename', 'Move', 'Copy', 'Preview', 'Preview Images', 'Download', and 'Move to Trash'. The project details for 'PRJ-2978 | Automated Neighborhood Characteristics for Community Resilience Planning' are shown, including the authors 'Rosenheim, Nathanael; Day, Wayne; Seong, Kijin', the project type 'Other', and the award 'Center for Risk-Based Community Resilience Planning - NIST-70NANB20H008'. A red box highlights the text 'Project team included two PhD students - project took about 2 months from start to finish.' with two red arrows pointing to the authors' names in the project details.

designsafe-ci.org/data/browser/public/designsafe.storage.published/PRJ-2978

NSF | NHERI

DESIGNSAFE

Log in Register

Workspace Learning Center NHERI Facilities NHERI Community News Help

Search DesignSafe

DATA DEPOT

Find in Data Depot

Rename Move Copy Preview Preview Images Download Move to Trash

PRJ-2978 | Automated Neighborhood Characteristics for Community Resilience Planning

Download Dataset

Author Rosenheim, Nathanael; Day, Wayne; Seong, Kijin

Project Type Other

Awards Center for Risk-Based Community Resilience Planning - NIST-70NANB20H008

Published

Published (NEES)

Project team included two PhD students - project took about 2 months from start to finish.

Rosenheim, Nathanael; Day, Wayne; Seong, Kijin (2021) "Automated Neighborhood Characteristics for Community Resilience Planning." DesignSafe-CL. <https://doi.org/10.17603/ds2-hj0p-bp40>.



# Project files developed and created in DesignSafe

process to obtain, clean, and explore census block group data. The data outputs include CSV files required to run a post-disaster population displacement model currently in use by IN-CORE. This project applies the code to six community testbeds (Seaside, OR; Joplin, MO; Galveston, TX; Lumberton, NC; Memphis MSA; and Mobile, AL) to illustrate the generalizability of the code. The code utilizes Census API and may be modified to identify some of the socio-demographic and socio-economic characteristics of a neighborhood's social vulnerability.

PRJ-2978

<input checked="" type="checkbox"/> Name	Size	Last modified
<input type="checkbox"/>  <a href="#">_ReadMeFirst_CensusAPI_BlockGroup_2021-01-19.docx</a> <b>README</b>	79.9 kB	1/28/21 1:37 PM
<input type="checkbox"/>  <a href="#">IN-CORE_ACSdata_BGMAP_2021-01-19.ipynb</a> <b>Jupyter Notebook</b>	1.7 MB	1/28/21 1:37 PM
<input type="checkbox"/>  <a href="#">IN-CORE_BGMAP_2021-01-19</a> <b>Data</b>	--	1/28/21 1:22 PM
<input type="checkbox"/>  <a href="#">IN-CORE_BGMAP_2021-01-19.ipynb</a> <b>Jupyter Notebook</b>	6.5 MB	1/28/21 1:37 PM

# Project Developed Within Jupyter-DesignSafe

jupyter

LogoutControl Panel

FilesRunningClusters

Select items to perform actions on them.

UploadNew↺

0

/ MyProjects / PRJ-2978

Name↓Last ModifiedFile size

..

seconds ago

IN-CORE\_BGMAP\_2021-01-19

a day ago

IN-CORE\_CENSUSAPI\_DWNLD\_JOPLINMSA\_2021-01-11

a day ago

WorkKJS

WorkNPR

WorkWCD

IN-CORE\_ACSdata\_BGMAP\_2021-01-19.ipynb

10 days ago6.81 MB

IN-CORE\_BGMAP\_2021-01-19.ipynb

15 days ago245 kB

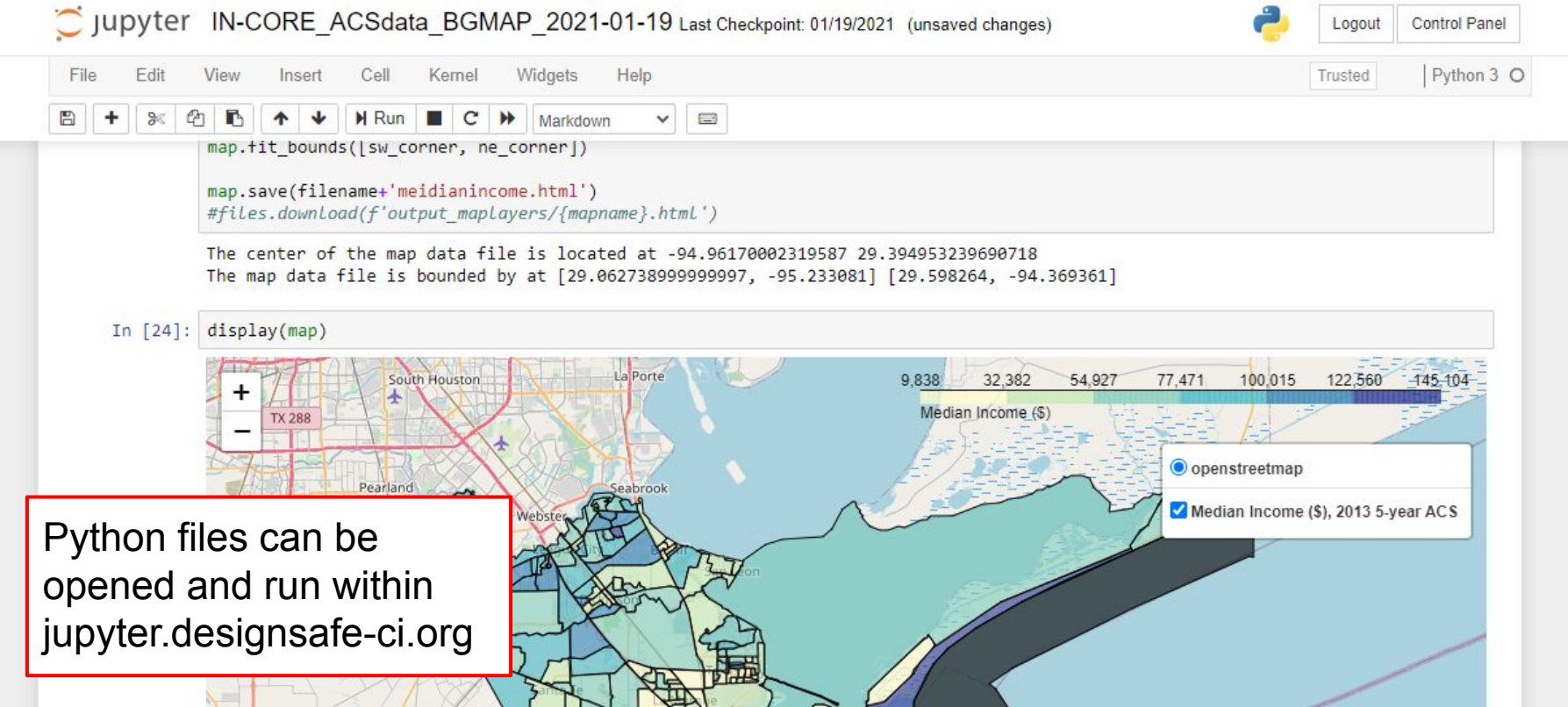
IN-CORE\_CENSUSAPI\_DWNLD\_JOPLINMSA\_2021-01-13.ipynb

10 days ago81.8 kB

\_ReadMeFirst\_CensusAPI\_BlockGroup\_2021-01-19.docx

Notice that internal project folders do not have to be published.

# Jupyter Notebook uses Python - Free and Open Source



# Thank you!

Nathanael Rosenheim

[nroseheim@arch.tamu.edu](mailto:nroseheim@arch.tamu.edu)

# Thoughts after the workshop

Data science may not be interested in reproducible, open, collaborative research.

What tools help make collaboration and replication easier - “Better science in less time”

Ideal dataset - provides individual level characteristics that identify a person or business. Provides information that helps distinguish and predict outcome. For example, which applicants to Texas A&M will decide to enroll? Need to have enough detail about each applicant to predict the outcome.

Ideal dataset - is linkable to other datasets that provide even more details and information.

# Thoughts after the workshop

Some discussion about the drawbacks of synthetic, derived or sanitized data.

Discussion about how “hiding” data leads to bad data. Need to reward making data available and more usefull.

Thoughts about decision makers interpretation of data analytics:

- Each decision maker (example a football coach) will have different ways they want to see the interpretation of data. For example, some like dashboards, others want just a few numbers.
- Emotions in the room make a difference
- The person analyzing the data needs to be very familiar with the data - actually needs to see the data collected.
- Need to “read the audience” as the story is being told.

# Thoughts after the workshop

Common theme across presentations “Real action comes from individual data”. Seemed to be true for the perspectives of public health, transportation, athletics, admissions.

When thinking about data collection start by planning to share the data.

An option for data sharing: Texas Data Repository <https://data.tdl.org/>