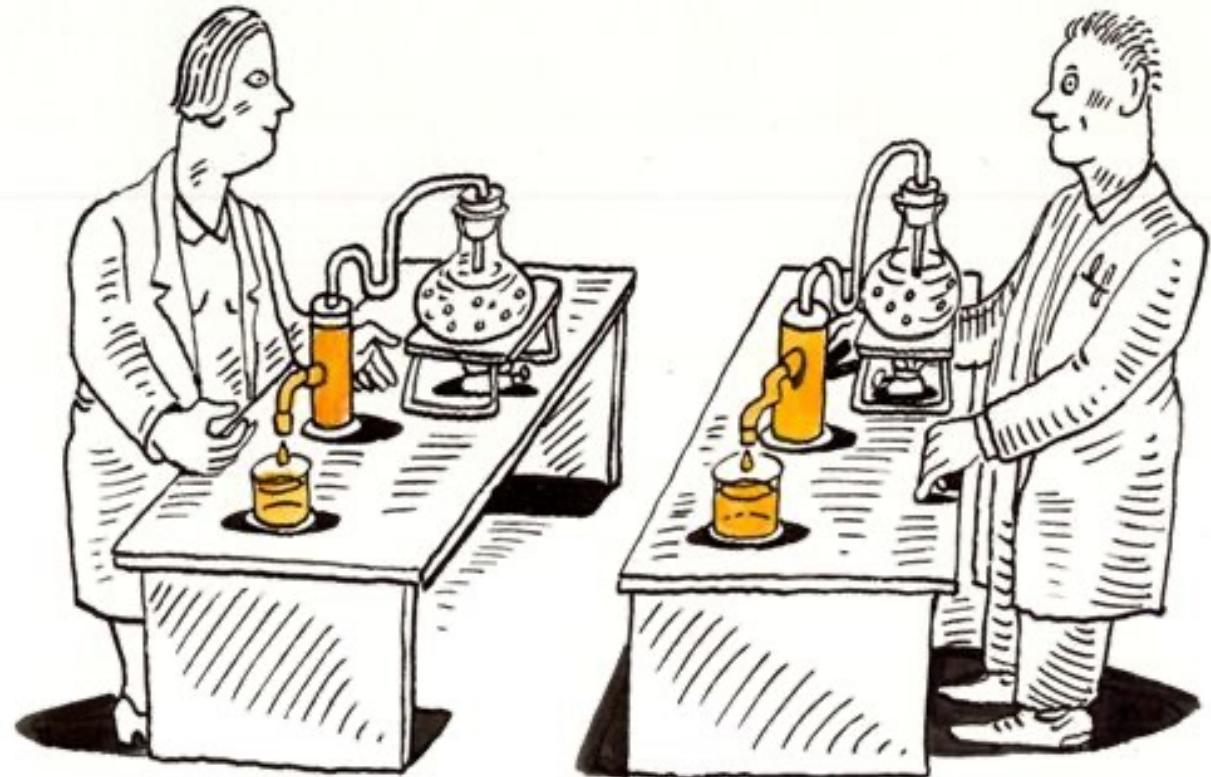


# Introduction to Reproducible Science: *Day 1*

Nelson Roque, PhD



Scan for  
Slides & Code  
available on GitHub



# Background

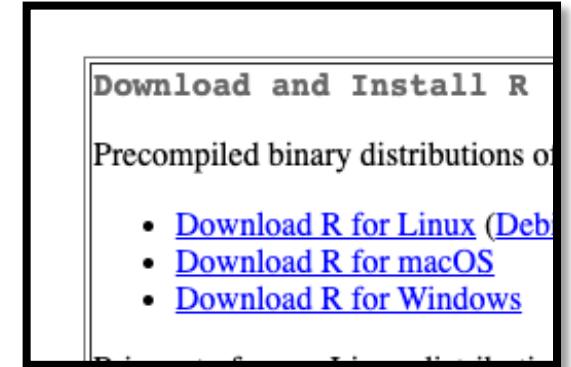
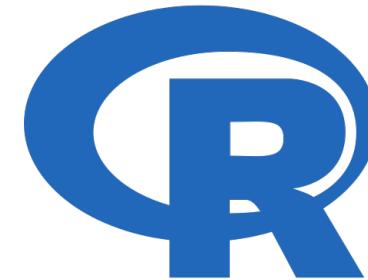
- Over the last decade I've engaged in professional development and research opportunities to become proficient across topics including data wrangling and modeling of text, image, video, and eye-tracking data, as well as more recently, sensor data.
- I look forward to training the next generation of scientists on code-based methods to enhance reproducibility and usability of their research results.

# About this Workshop

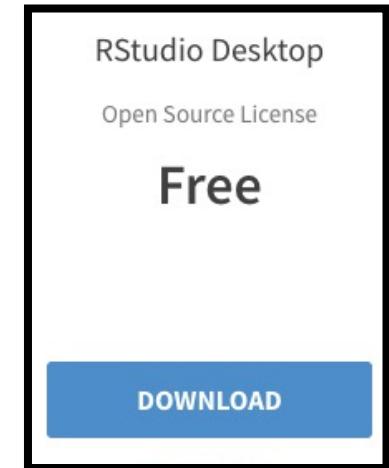
- 3 days: July 6, 8, 11, 11am to 2pm
- In-person, room 301Q, but also recorded
- All slides, code, and recordings will be made available at the website listed below, as well as Github
  - Code & Slides
    - [https://github.com/nelsonroque/contextlab\\_reproduciblescience\\_workshops](https://github.com/nelsonroque/contextlab_reproduciblescience_workshops)
  - Textbook
    - [https://nelsonroque.github.io/contextlab\\_introdatascience\\_webcourse/](https://nelsonroque.github.io/contextlab_introdatascience_webcourse/)

# Software Required for Workshop

- R - <https://cran.r-project.org/>



- RStudio -  
<https://www.rstudio.com/products/rstudio/download/>



# Workshop Learning Objectives

1. Describe various principles, tools and techniques supportive of open and reproducible science.
  1. List and describe the FAIR Principles (<https://www.go-fair.org/fair-principles>)
2. Develop a code-only pipeline to allow reproducibility of data prep and analyses.
  1. Learn about data wrangling principles (e.g., tidy data)
  2. Import, pre-process & visualize data using R
3. Develop a long-term learning plan for practicing reproducible science tools and techniques.
  1. Share recommended readings, activities for continued learning

# Agenda: Day 1

- What is Reproducible Science?
- Reproducible & FAIR Data Workflows
- Tools Supporting Reproducible Science
  - Overview of available tools
    - **Skill 1: Using Endnote for Reference Management**
    - **Skill 2: Using Git (and Github) for code management and collaboration**
  - Orientation to R, RStudio, RMarkdown
    - **Skill 3: R syntax primer**
- Data Science: Latest trends
- Long-term Learning Recommendations

# Agenda: Day 2

- Data wrangling and visualization of Big Data
  - **Skill 1: Data wrangling the Google Mobility dataset**
- Reproducible survey research
  - Qualtrics survey design tips
  - **Skill 2: Data wrangling Qualtrics data**
- Working with JSON data
  - **Skill 3: cleaning and visualizing keystroke JSON data**

# Agenda: Day 3

- Text mining
  - Skill 1: word and bigram frequency analysis
  - Skill 2: generating wordclouds
  - Skill 3: sentiment analysis
- Interacting with APIs and JSON data
  - Skill 4: querying API for results and data aggregation
- Closing Discussion & Q/A

# **What is Reproducible Science?**

# Icebreaker: Introductions

1. Name, lab affiliation
2. What does reproducible science mean to you?
3. Which tools have you previously used?
4. Which tools are you most excited to learn?



# Background

- A reproducibility crisis (Ioannidis, 2005; Open Science Collaboration, 2015) has emerged as a threat to the scientific enterprise.
  - Ioannidis, John P A. 2005. “Why Most Published Research Findings Are False.” *PLoS Medicine* 2 (8): e124. doi:10.1371/journal.pmed.0020124.
  - Munafò, M. R., Nosek, B. A., Bishop, D. V., Button, K. S., Chambers, C. D., Percie du Sert, N., ... & Ioannidis, J. (2017). A manifesto for reproducible science. *Nature human behaviour*, 1(1), 1-9.
  - Open Science Collaboration. 2015. “Estimating the Reproducibility of Psychological Science.” *Science* 349 (6251): aac4716–aac4716. doi:10.1126/science.aac4716.

# What is reproducible science?

- Making entire scientific process transparent (as allowable; often required by law, grants).
- Sharing experiment, raw data, and analysis code
- Detailed methods sections
- Sharing stimulus sets
- Being able to walk through a data analysis start (load raw data) to finish (manuscript analyses) in code
- Following FAIR Principles
- Pre-registering hypotheses and analysis plans (e.g., on OSF)

# Benefits of reproducible research

- increased likelihood that the research will be correct
- reproducibility makes it easier to check the research
- it is easier to reproduce the research independently
- easier to extend the research
- reusable code and instruction resulting in increased efficiencies

<https://www.displayr.com/what-is-reproducible-research>

# Reproducible versus replicable

- **Reproducible:** when the exact results can be reproduced if given access to the original data, software, or code
- **Replicable:** research results can be reproduced by independent researchers using different methods.

# Improving Usability of Reproducible Research

- Hundreds of reporting guidelines exist to enhance usability of research

The screenshot shows the homepage of the EQUATOR Network. At the top, there is a logo for "equator network" with a green globe icon. To the right of the logo, the text "Enhancing the QUAlity and Transparency Of health Research" is displayed. Below the logo, a navigation bar includes links for Home, About us, Library, Toolkits, Courses & events, News, Blog, Librarian Network, and Contact. A green banner below the navigation bar states "Your one-stop-shop for writing and publishing high-impact health research" and lists links for finding reporting guidelines, improving writing, joining courses, running training, enhancing peer review, and implementing guidelines. On the left, there is a section titled "Library for health research reporting" with a sub-section for "Search for reporting guidelines". On the right, there is a section titled "Reporting guidelines for main study types" listing various guidelines such as CONSORT, STROBE, PRISMA, SPIRIT, etc., each with a link to its respective page. At the bottom, there are three tabs: "Toolkits" (red), "EQUATOR highlights" (blue), and "News" (yellow). The "EQUATOR highlights" tab contains the text "6/01/2022 - ICMJE Recommendations updated to include new". The "News" tab contains the text "EQUATOR Network Newsletter April 2022".

EQUATOR resources in  
German | Portuguese |  
Spanish

Your one-stop-shop for writing and publishing high-impact health research

find reporting guidelines | improve your writing | join our courses | run your own training course | enhance your peer review | implement guidelines

Library for health research reporting

The Library contains a comprehensive searchable database of reporting guidelines and also links to other resources relevant to research reporting.

Search for reporting guidelines

Not sure which reporting guideline to use?

Reporting guidelines under development

Visit the library for more resources

Reporting guidelines for main study types

Randomised trials	CONSORT	Extensions
Observational studies	STROBE	Extensions
Systematic reviews	PRISMA	Extensions
Study protocols	SPIRIT	PRISMA-P
Diagnostic/prognostic studies	STARD	TRIPOD
Case reports	CARE	Extensions
Clinical practice guidelines	AGREE	RIGHT
Qualitative research	SRQR	COREQ
Animal pre-clinical studies	ARRIVE	
Quality improvement studies	SQUIRE	Extensions
Economic evaluations	CHEERS	

See all 527 reporting guidelines

Toolkits

EQUATOR highlights

News

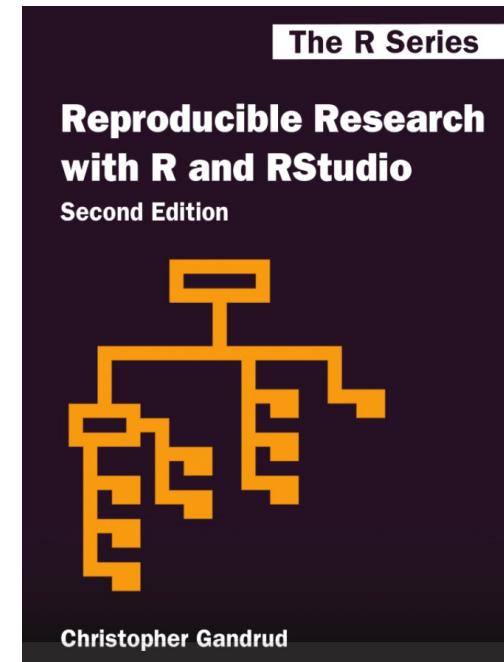
6/01/2022 - ICMJE Recommendations updated to include new

EQUATOR Network Newsletter April 2022

# **Reproducible & FAIR Data Workflows**

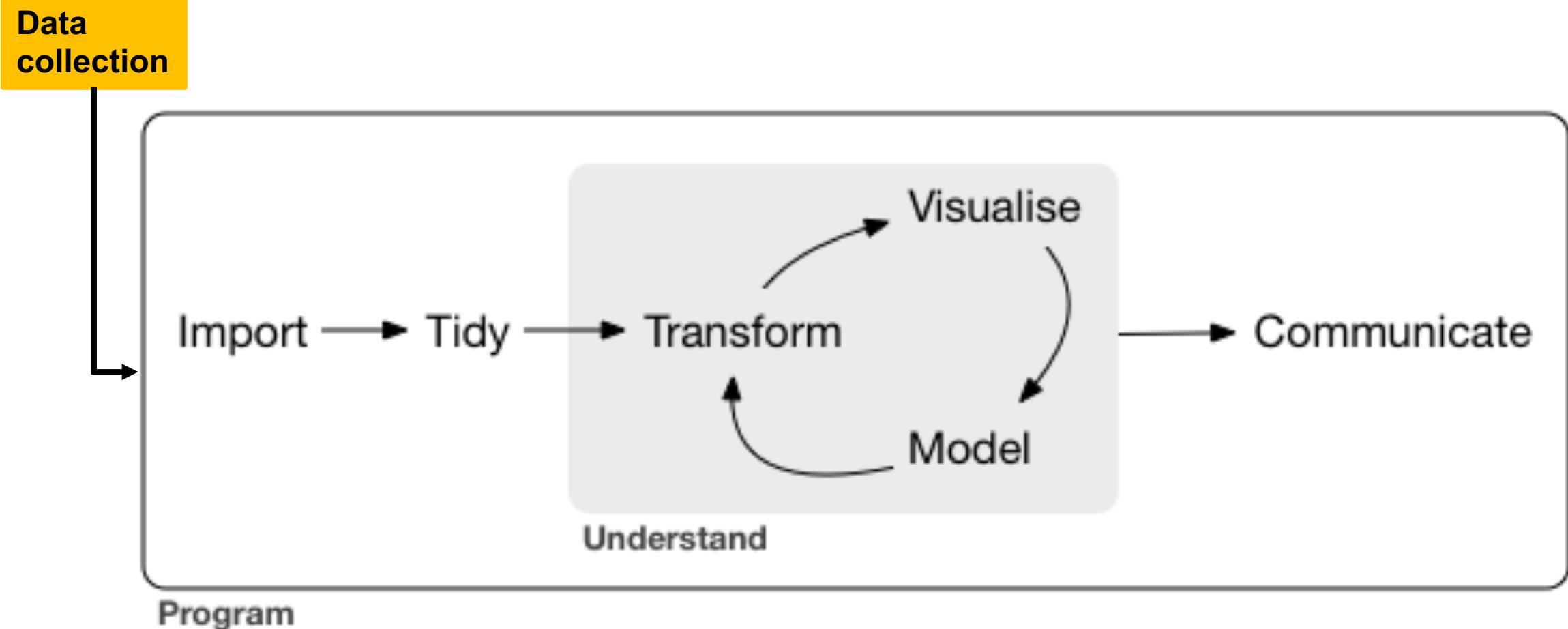
# Reproducible Workflows

- Documented, code-based approach to pre-processing and analytics
  - R (RMarkdown), Python (Jupyter notebooks)
- File versioning with clear changelogs
  - Git and Github



<https://monashdatafluency.github.io/r-rep-res/>

# Typical Data Science Workflow



Source: [R for Data Science](#)

# What is tidy data?

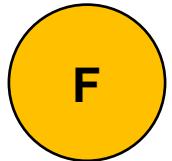
1. Every column is a variable.
2. Every row is an observation.
3. Every cell is a single value.

participant	condition	avg_response_time	perc_accuracy
9991	control	506	90
9992	active	516	96
9993	control	526	99

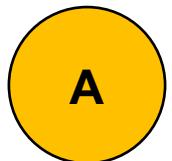
# Uncleaned Data Problems

- In a perfect data-generating system, data are already tidy
- In practice, we must first determine
  - data structure issues, including:
    - Column headers are values, not variable names.
    - Multiple variables are stored in one column.
    - Variables are stored in both rows and columns.
    - Multiple types of observational units are stored in the same table.
    - A single observational unit is stored in multiple tables.
  - Data quality issues
    - Values out of range
    - Improperly/inconsistently coded response options
    - Inconsistent records counts

# FAIR Principles



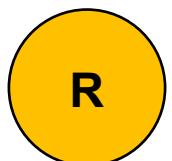
**Findable:** "The first step in (re)using data is to find them. Metadata and data should be easy to find for both humans and computers. Machine-readable metadata are essential for automatic discovery of datasets and services"



**Accessible:** "Once the user finds the required data, she/he/they need to know how they can be accessed, possibly including authentication and authorisation."



**Interoperable:** "The data usually need to be integrated with other data. In addition, the data need to interoperate with applications or workflows for analysis, storage, and processing."



**Reusable:** "The ultimate goal of FAIR is to optimise the reuse of data. To achieve this, metadata and data should be well-described so that they can be replicated and/or combined in different settings."

# **Tools Supporting Reproducible Science**

# Tools Supporting Reproducible Science

- R
- Latex, Markdown
- Python, Anaconda
- Github, Github Pages: <https://pages.github.com>
- Docker, VMs
- Infrastructure as Code (IaC; e.g., AWS CDK)
- IDEs (e.g., Visual Studio Code)
- Open Science Framework
- Documentation software: Docusaurus, <https://docusaurus.io>
- Code-based experiment creation software: Opensesame, Psychopy, jsPsych



# Markdown

Lightweight markup used to add formatting elements to plaintext documents.

Created by [John Gruber](#) in 2004, Markdown is now one of the world's most popular markup languages.

Using Markdown is different than using a [WYSIWYG](#) editor.

- In Word, you click buttons in GUI to apply formatting
- In Markdown, you add syntax to text to apply formatting

Element	Markdown Syntax
Heading	# H1 ## H2 ### H3
Bold	<b>bold text</b>
Italic	<i>italicized text</i>
Blockquote	> blockquote
Ordered List	1. First item 2. Second item 3. Third item
Unordered List	- First item - Second item - Third item
Code	`code`
Horizontal Rule	---
Link	[title](https://www.example.com)
Image	![alt text](image.jpg)

<https://www.markdownguide.org/getting-started/>

# What is Git?

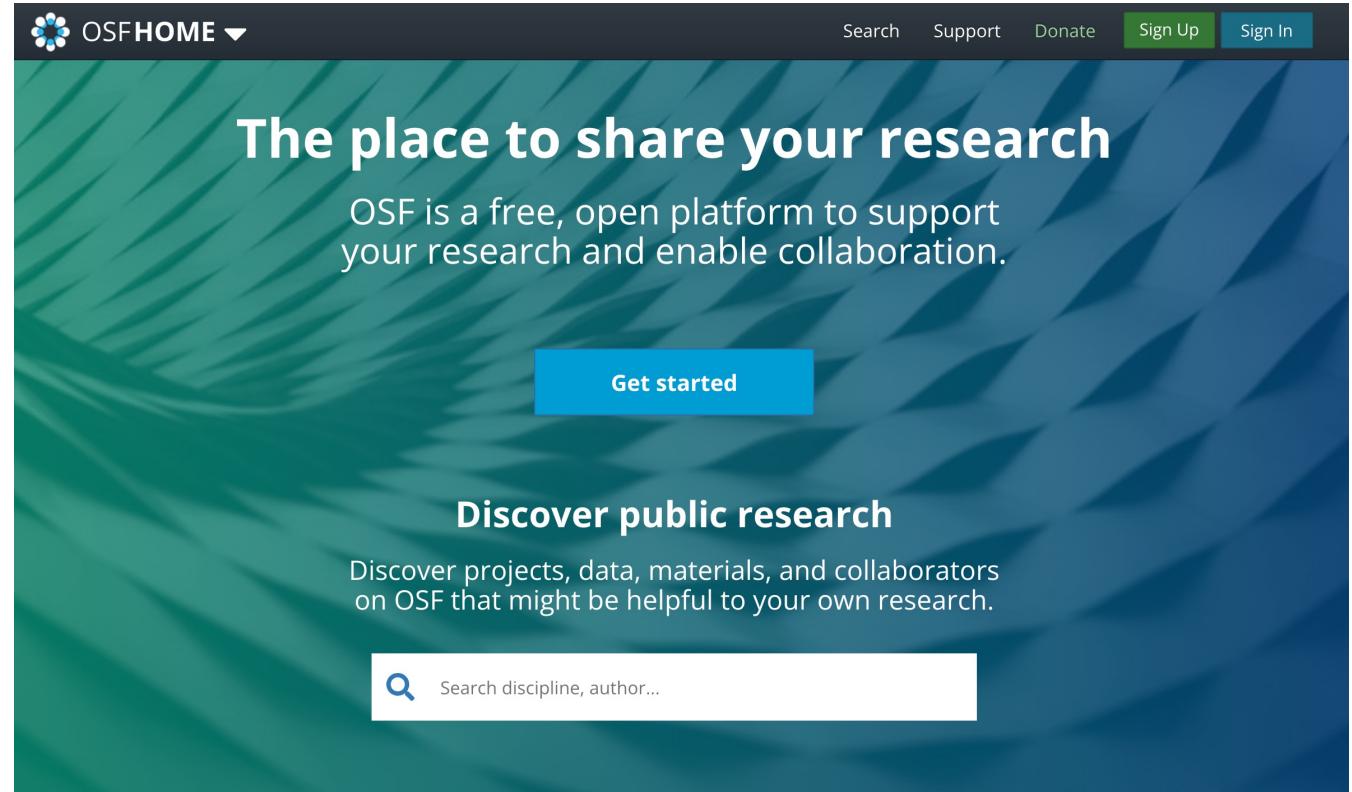
- A version control system to manage code; edits, branching, and collaboration.
- Git providers
  - Open source Git: <https://git-scm.com/>
  - Github - <https://desktop.github.com>
    - Get Github Cheatsheets: <https://training.github.com/downloads/github-git-cheat-sheet/>
  - Gitlab

# Git These Concepts

- **Repository:** a logical unit of storage for maintaining code/other assets. Can contain many folders, files.
- **.gitignore:** a file (with dot prefix) to specify what files should NOT be pushed to Github
- **.gitkeep:** a file (with dot prefix) to specify that a folder should be kept if blank as part of the project structure. This is helpful when creating templates.
- **commit:** create a record/snapshot of all files (or specific files) at a moment in time. You may tag collaborators, add a title and notes.
- **push:** 'storing' the result of prior commits
- **pull:** 'extracting' the result of prior commits
- **diff:** difference (line-specific) between two commits. Helpful when trying to remember when and where code edits were made.

# OSF

- Preregister study design, hypotheses, analysis plans and scripts
- Post-hoc curate documentation for your study across various components (IRB, surveys, stimuli, datasets, codebooks)



# PsyArXiv

- Publish pre-prints of your paper
- Supports immediate sharing of research
- <https://psyarxiv.com>



# Let's jump into Github

Download here:

<https://desktop.github.com>



# Endnote

- **What is Endnote?**

- a citation manager (like Zotero, Mendeley)
- Guide from the UCF library with installation links:  
<https://guides.ucf.edu/citations-endnote>

- **Why Endnote?**

- no need for manual tracking (deleting, updating) of references
- change citation format at click of a button
- capture PDFs for any citations at click of a button

- **Downloads:**

- [Mac](<http://ezproxy.library.ucf.edu/loggedin/EndNote20SiteInstaller.zip>)
- [Windows](<http://ezproxy.library.ucf.edu/loggedin/EndNote20.zip>)

# **Let's jump into Endnote**

Download here:

<https://guides.ucf.edu/citations-endnote>



**EN**

# **Orientation to R, RStudio, RMarkdown**

# What is R? RStudio?

- **What is R?**
  - Statistical programming language (released ~1993). Similar to `S` (released ~1976)
  - Download R: <https://cran.r-project.org>
- **What is RStudio?**
  - Integrated development environment (IDE) for R; view data objects, plots, and file directory all in one place.
  - Download Rstudio: <https://www.rstudio.com/products/rstudio/download>

# What are R packages?

- **Libraries of code to accomplish specific functions**
  - e.g., data viz, machine learning, data wrangling, compute effect sizes
- **Think of CRAN as an equivalent to iOS App Store, or Google Play Store**
  - CRAN is where you download published R packages (directly from R, RStudio)
- **A developer can also make packages available directly in source code, or via tools like Github, for example:**

```
install.packages("devtools")
devtools::install_github("nelsonroque/ruf")
```

# What is RMarkdown?

- Reproducible notebook (akin to Jupyter Notebooks in Python) where you can write a narrative using Markdown syntax and embed code/plots throughout.
  - Based on Markdown syntax: <https://www.markdownguide.org/basic-syntax/>
- With R Markdown, you can generate HTML, PDF, MS Word documents, PowerPoint presentations, and books!
- Learn More: <http://rmarkdown.rstudio.com>

The screenshot shows the RStudio interface with a dark theme. A red box highlights the top-left area where a script named "pipeline.R" is open. The script contains the following code:

```
1 # load libs required for pipeline
2 library(tidyverse)
3 library(readr)
4
5 # run scripts in your pipeline
6 source("scripts/load_data.R")
7 source("scripts/preprocess_data.R")
8 source("scripts/aggregate_data.R")
9 source("scripts/data_viz.R")
10 source("scripts/simple_stats.R")
```

The R console at the bottom displays the standard R startup message, indicating it's version 4.1.2 running on a Mac OS X system. The Global Environment panel on the right shows that the environment is currently empty.

Go to file/function | Addins

pipeline.R x

Source on Save | Run | Source

```
1 # load libs required for pipeline
2 library(tidyverse)
3 library(readr)
4
5 # run scripts in your pipeline
6 source("scripts/load_data.R")
7 source("scripts/preprocess_data.R")
8 source("scripts/aggregate_data.R")
9 source("scripts/data_viz.R")
10 source("scripts/simple_stats.R")
```

10:25 (Top Level) R Script

Console Terminal Jobs

R 4.1.2 ~/Documents/Github/contextlab\_vr\_fingertapping/ Platform: x86\_64-apple-darwin17.0 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.  
You are welcome to redistribute it under certain conditions.  
Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors.  
Type 'contributors()' for more information and  
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or  
'help.start()' for an HTML browser interface to help.  
Type 'q()' to quit R.

> plot(x=rnorm(100))
> `|`

Environment History Files Connections Packages Git Tutorial

New Folder New Blank File Delete Rename More

Home Documents Github contextlab\_vr\_fingertapping scripts

Name	Size	Modified
..		
aggregate_data.R	1.7 KB	Jun 1, 2022, 7:23 PM
data_viz.R	562 B	Jun 1, 2022, 7:29 PM
load_data.R	1.2 KB	Jun 1, 2022, 7:09 PM
preprocess_data.R	696 B	Jun 1, 2022, 7:34 PM
simple_stats.R	124 B	Jun 1, 2022, 7:18 PM

# Files

Plots Help Viewer

Zoom Export Publish

norm(100)

Index

A screenshot of the RStudio interface. The top menu bar includes 'File', 'Edit', 'Source', 'Run', 'View', 'Tools', 'Help', and 'Addins'. The title bar shows 'contextlab\_vr\_fingertapping'. The left panel contains a script named 'pipeline.R' with the following code:

```
1 # load libs required for pipeline
2 library(tidyverse)
3 library(readr)
4
5 # run scripts in your pipeline
6 source("scripts/load_data.R")
7 source("scripts/preprocess_data.R")
8 source("scripts/aggregate_data.R")
9 source("scripts/data_viz.R")
10 source("scripts/simple_stats.R")
```

The right panel shows the 'Environment' tab with the message 'Environment is empty'.

# Console: run one-off commands

The R console window is highlighted with a red border. It displays the following startup messages:

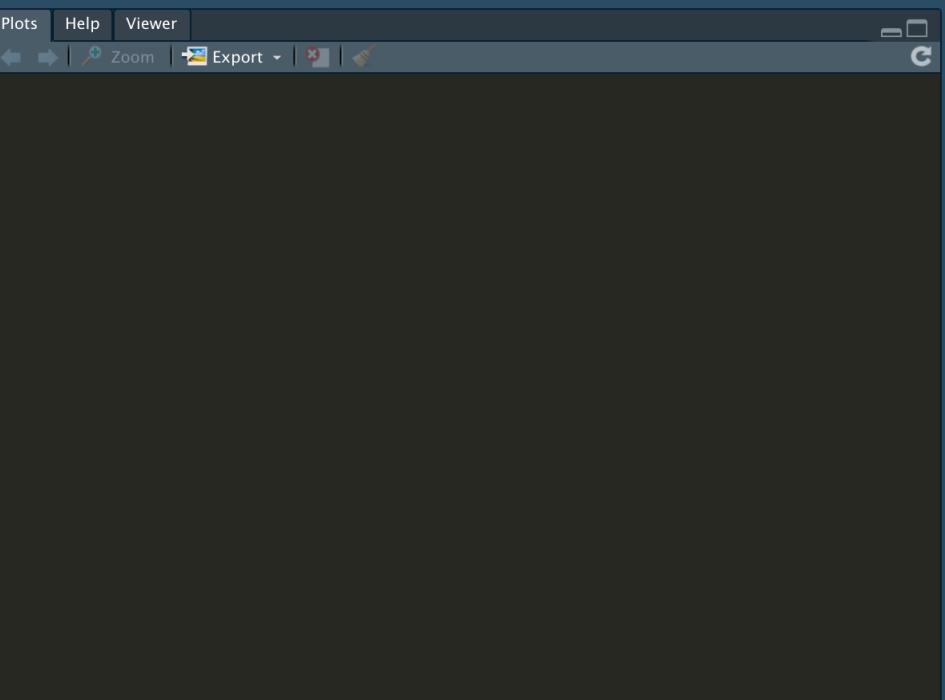
```
10:25 (Top Level) R Script
Console Terminal Jobs
R 4.1.2 · ~/Documents/Github/contextlab_vr_fingertapping/
Copyright (C) 2021 The R Foundation for Statistical Computing
Platform: x86_64-apple-darwin17.0 (64-bit)

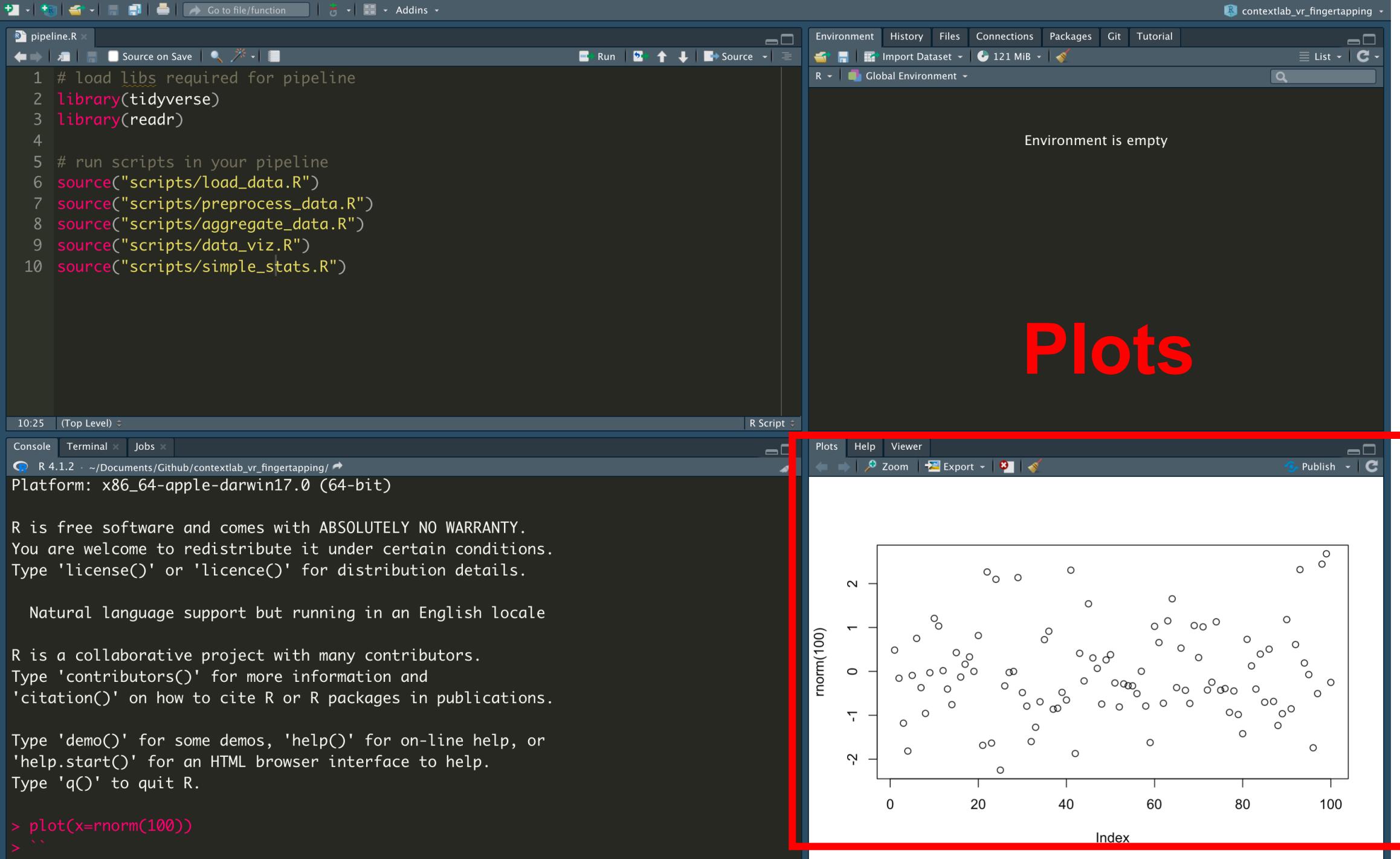
R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

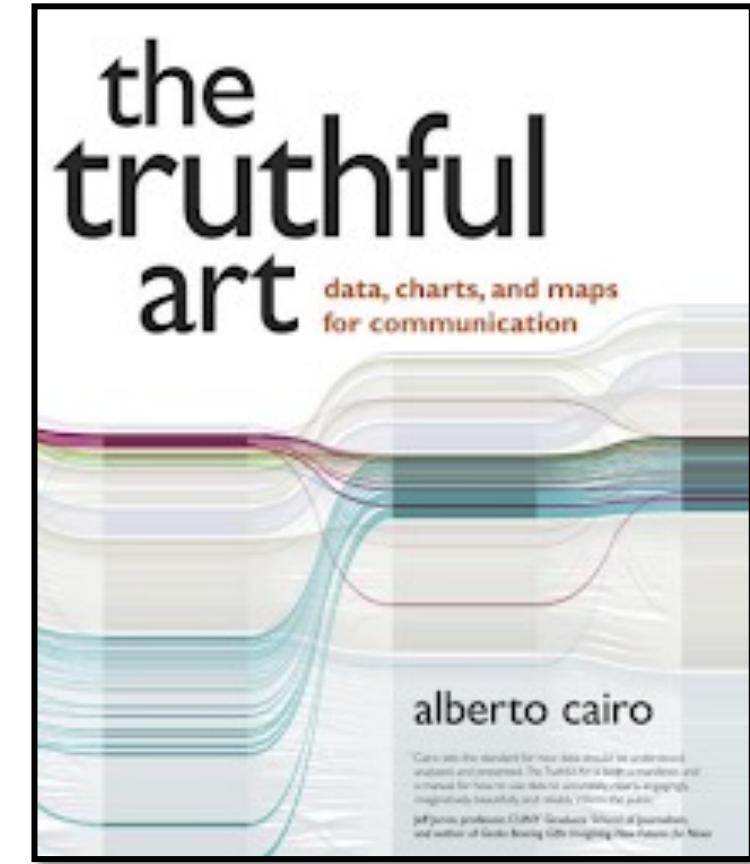
Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.
```





# PSA: Five Qualities of Great Visualizations

1. **Truthful:** transparency in decision making and reporting
2. **Functional:** understandable, uncluttered
3. **Beautiful:** plots can be art too
  1. <https://playingthearchive.com>
4. **Insightful:** reveal trends or relationships
5. **Enlightening:** integration of 1-4 executed successfully



The screenshot shows the RStudio interface with several panes:

- Code Editor (Top Left):** Displays the contents of `pipeline.R`. The code includes loading libraries (`tidyverse`, `readr`), running scripts in a pipeline, and source files for data loading, preprocessing, aggregation, visualization, and simple statistics.
- Environment (Top Right):** Shows the Global Environment pane with the message "Environment is empty".
- Help (Bottom Right):** A large red box highlights the Help pane for the `cut` function. The title is "Convert Numeric to Factor".
  - Description:** `cut` divides the range of `x` into intervals and codes the values in `x` according to which interval they fall. The leftmost interval corresponds to level one, the next leftmost to level two and so on.
  - Usage:** `cut(x, ...)`
  - Arguments:**
    - `x`: a numeric vector which is to be converted to a factor by cutting.
    - `breaks`: either a numeric vector of two or more unique cut points or a single number (greater than or equal to 2, giving the number of intervals) indicating the number of intervals into which the range should be divided.
- Console (Bottom Left):** Displays the R startup message, license information, and a welcome message about natural language support.

The screenshot shows the RStudio interface with the following components:

- Top Bar:** Contains icons for file operations (New, Open, Save, Print, etc.), a search bar, and tabs for "Go to file/function" and "Addins".
- Left Panel:** Shows the code editor with a script named "pipeline.R". The code includes imports for tidyverse and readr, and sources several other scripts for data loading and processing.
- Right Panel:** Shows the "Global Environment" pane with two objects:
  - df:** 100 obs. of 3 variables
  - fit:** Formal class lmerMod

# Viewer

The screenshot shows the RStudio console with the following output:

```
R 4.1.2 · ~/Documents/Github/contextlab_vr_fingertapping/ ↵
boundary (singular) fit: see help('isSingular')
Linear mixed model fit by REML ['lmerMod']
Formula: y ~ x + (1 | group)
Data: df
REML criterion at convergence: 291.8917
Random effects:
 Groups   Name        Std.Dev.
 group    (Intercept) 0.000
 Residual           1.024
Number of obs: 100, groups:  group, 17
Fixed Effects:
(Intercept)      x
 -0.06577    -0.12750
optimizer (nloptwrap) convergence code: 0 (OK) ; 0 optimizer warnings; 1 lme4 warnings
> fit = lme4::lmer(y ~ x + (1|group), data = df)
boundary (singular) fit: see help('isSingular')
> sjPlot::tab_model(fit)
> |
```

The screenshot shows the "Viewer" pane of the sjPlot package, which displays a summary table of the fitted model. The table is titled "y" and includes columns for Predictors, Estimates, CI, and p.

Predictors	Estimates	CI	p
(Intercept)	-0.07	-0.27 – 0.14	0.524
x	-0.13	-0.34 – 0.09	0.240

**Random Effects**

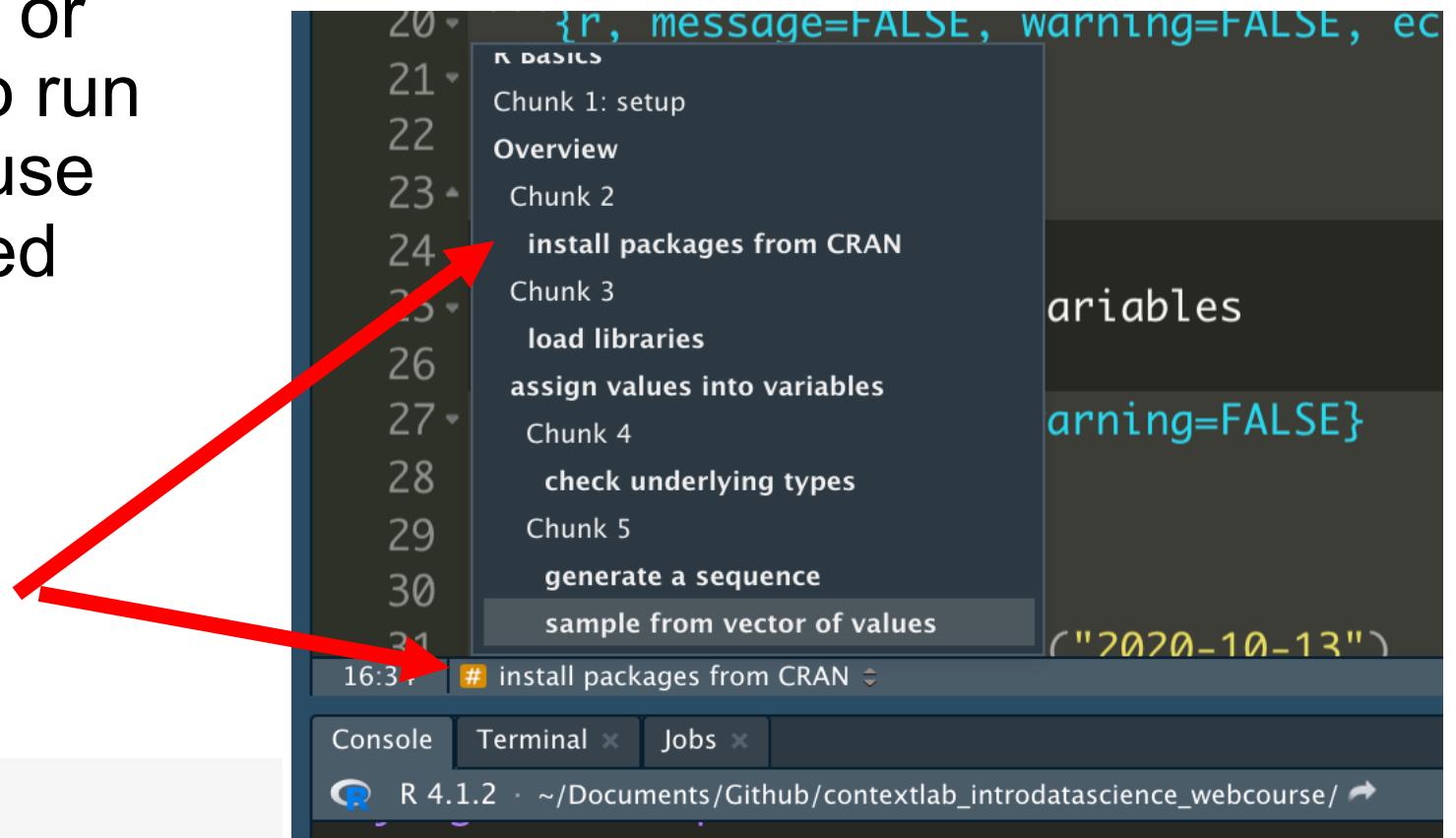
$\sigma^2$	1.05
$\tau_{00}$ group	0.00
N_group	17

Observations: 100  
Marginal R<sup>2</sup> / Conditional R<sup>2</sup>: 0.014 / NA

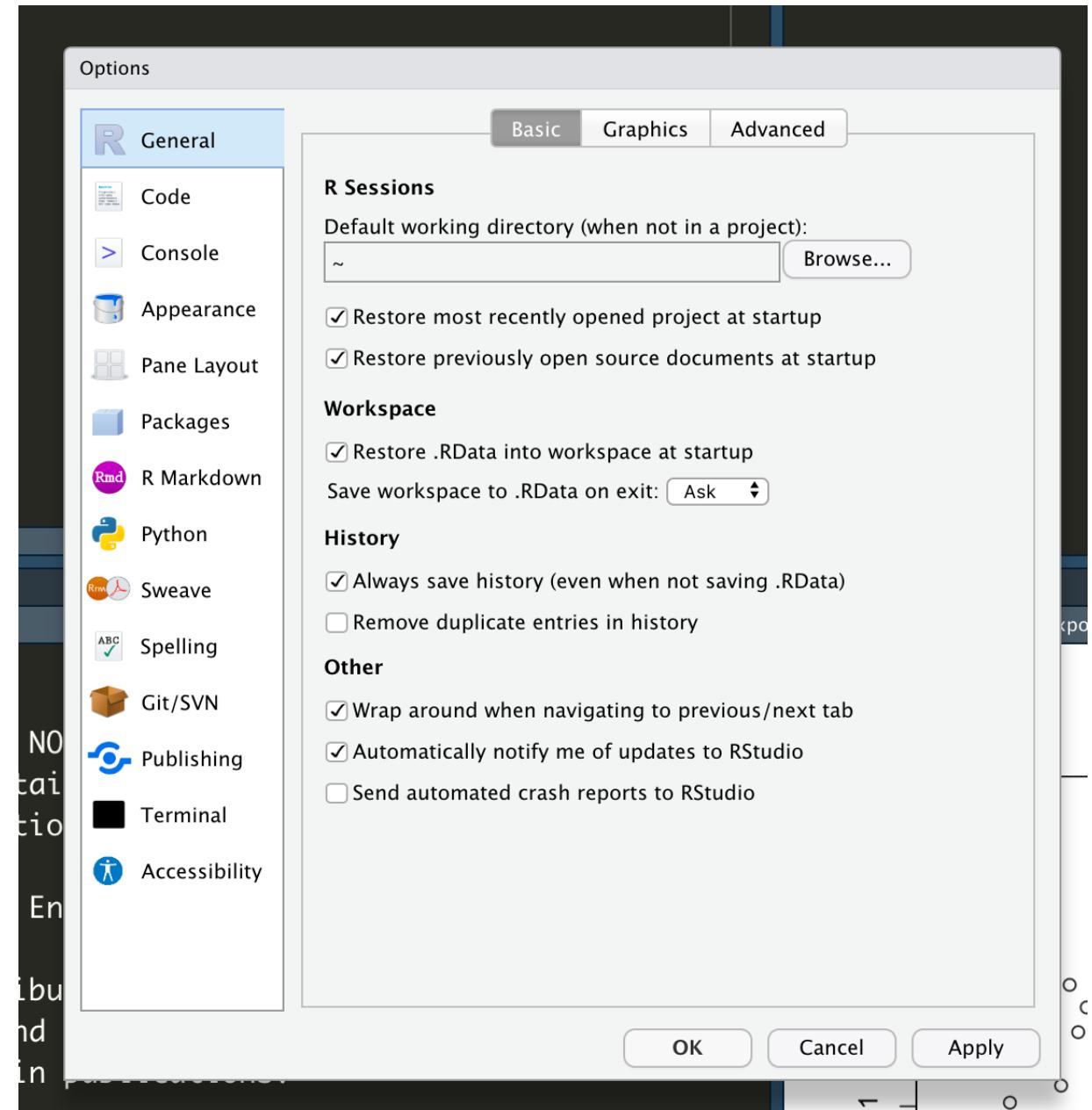
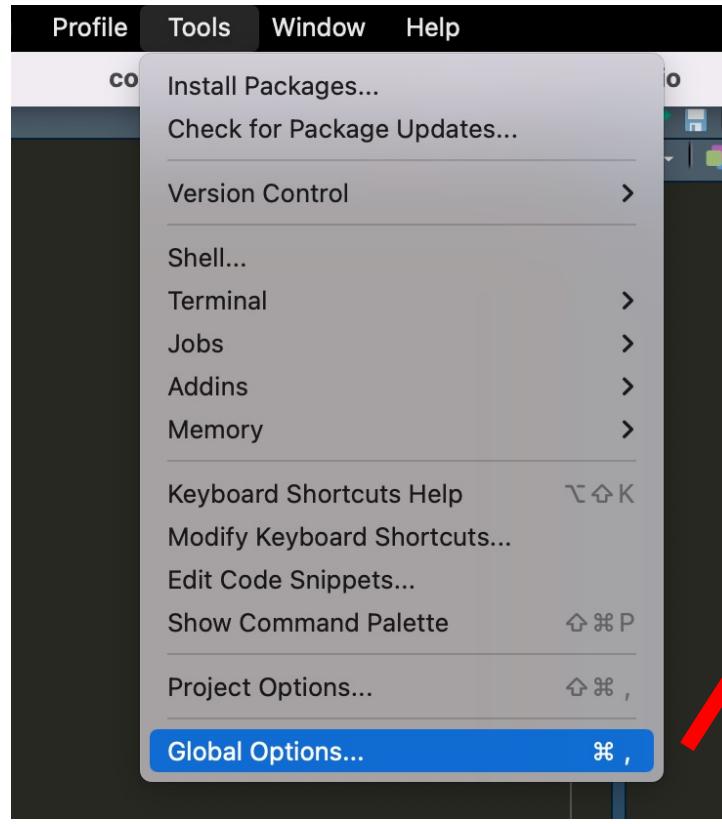
# Shortcuts

- Command+Enter (Mac) or Ctrl+Enter (Windows) to run the line where your mouse cursor is (or a highlighted selection of code)
- Add “----” inside of a comment to create a shortcut

```
15 ````{r}
16 # install packages from CRAN ----| 
17 install.packages(c("xlsx", "tidyverse", "devtools"))
18 ````
```



# Customizing RStudio



# **Intro to R Syntax**

# Basic data types & structures

- Character: text data
- Numeric: all real numbers with or without decimal values
- Integer: real values without decimal points
- Logical: TRUE, FALSE
- Raw: raw bytes of data
- List: named or numbered array of entries
- Data frame and tibble: rectangular representation of data
- Vector
- Matrix

<https://www.programiz.com/r/data-types>

# Pseudocode Exercise

You allowed participants on Qualtrics to enter state of residence as free text entry, resulting in e.g., Alabama, as: 'AL', 'Al', 'Alabama', 'ALABAMA'. Ideally you want them all to be two upper case letters.

**Goal:** in words only (no code), describe your approach to solving the data problem above.



# Let's jump into R

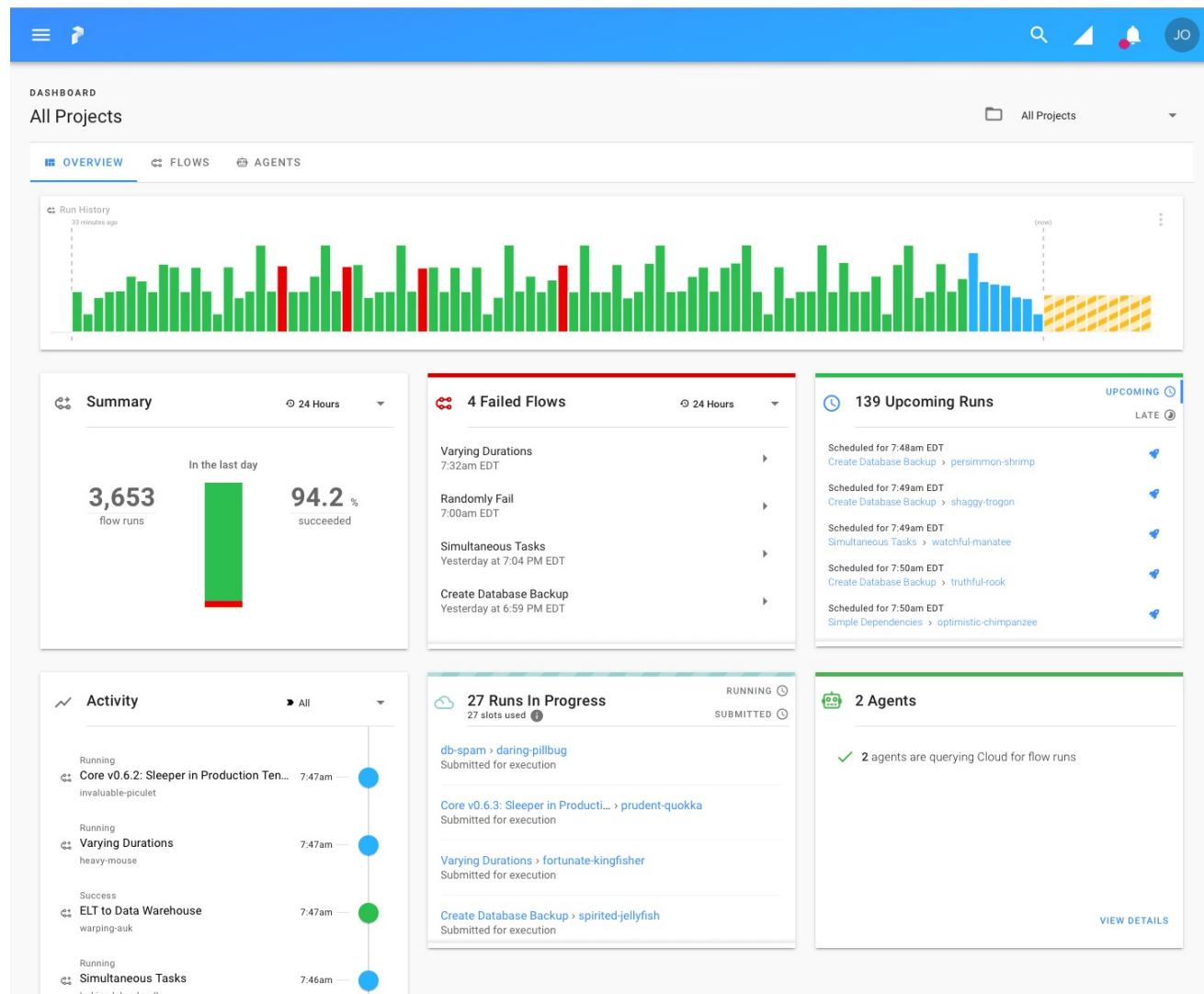
Scan the QR code to view the code we will work with



# Data Science: Latest Trends

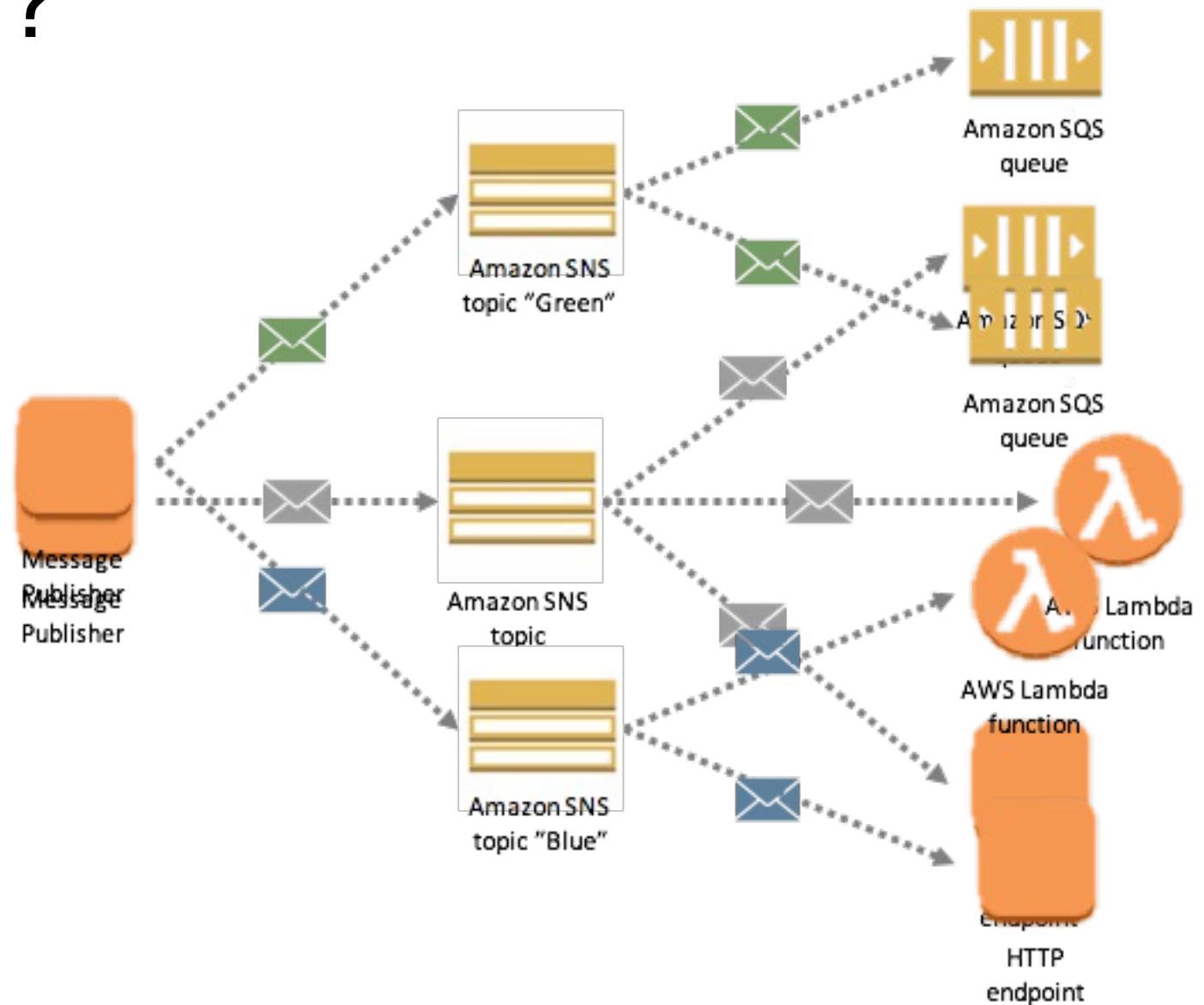
# ETL Pipelines

- Extract
  - Programmatically download data from origin (database, file server)
- Transform
  - Manipulate, score, reformat data as needed for computational needs
- Load
  - Re-save transform dataset(s) in new location on origin server (or other server)



# What is PubSub?

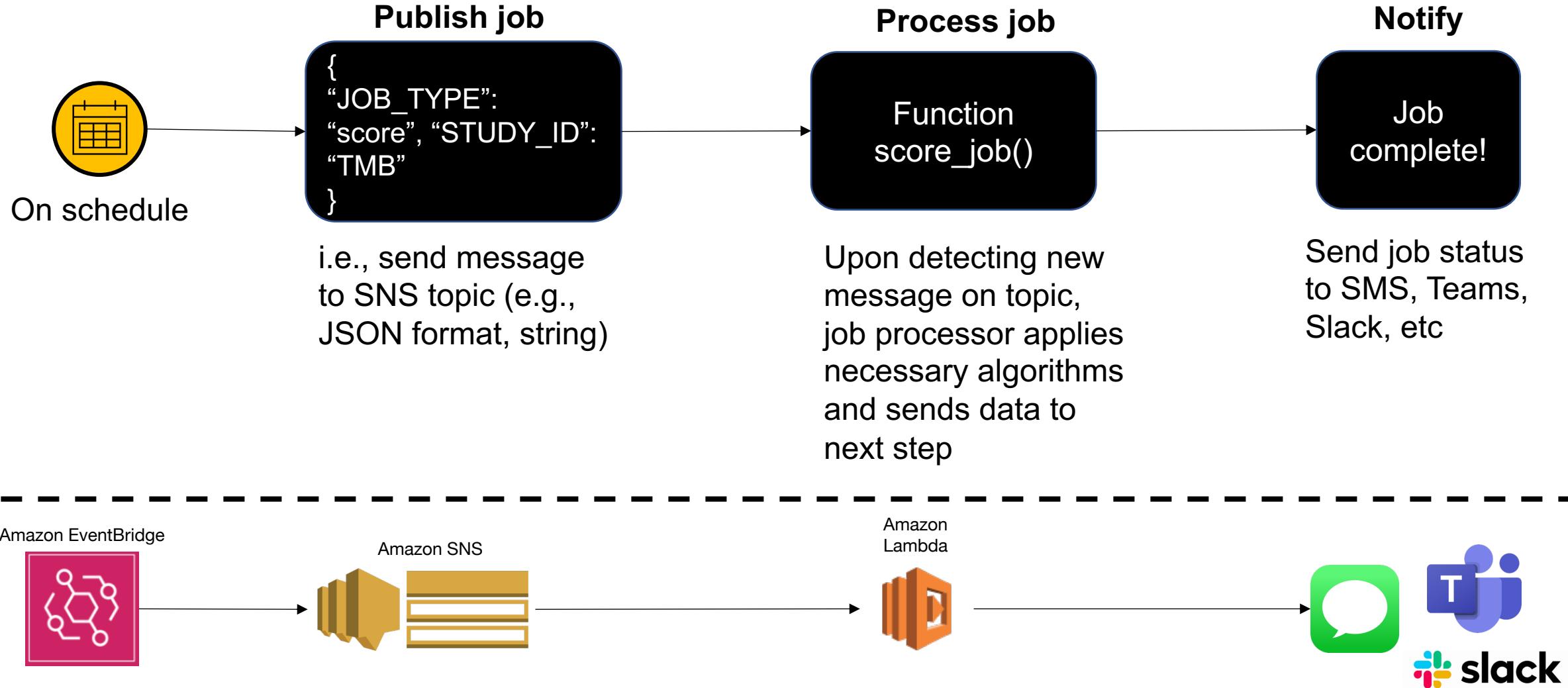
- Messaging pattern for loosely coupled communication between highly cohesive components
- **Why PubSub?**
  - Scalable, fault-tolerant, and highly available



# How I've Used PubSub Messaging

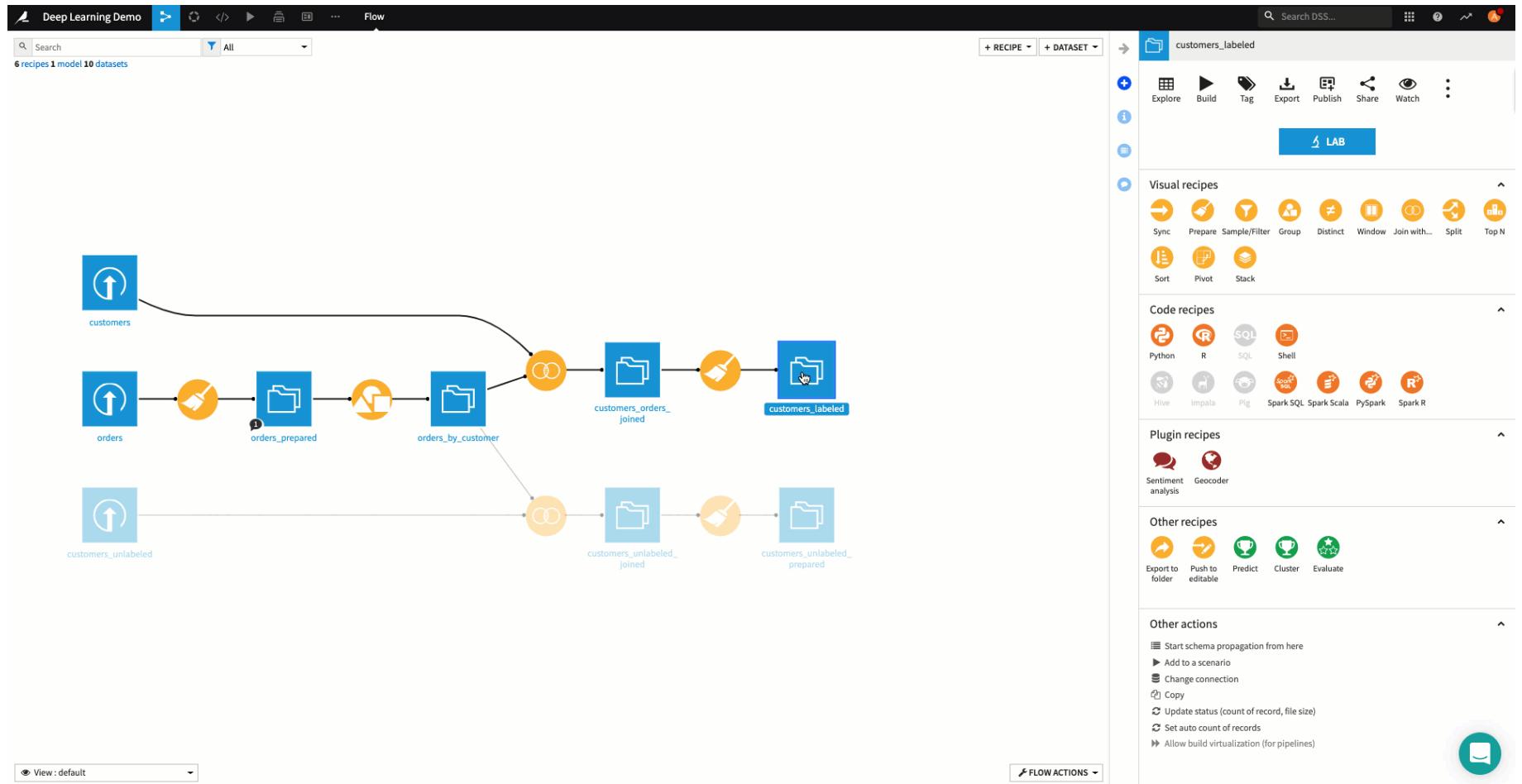
- Send a message to a topic for each device in a study, and query all data for last X days from device API
  - e.g., Fitbit, Air quality devices
- Send a message to a topic for each cognitive task, query respective data, score it, and publish to API/store in new database table

# Automating data jobs with PubSub



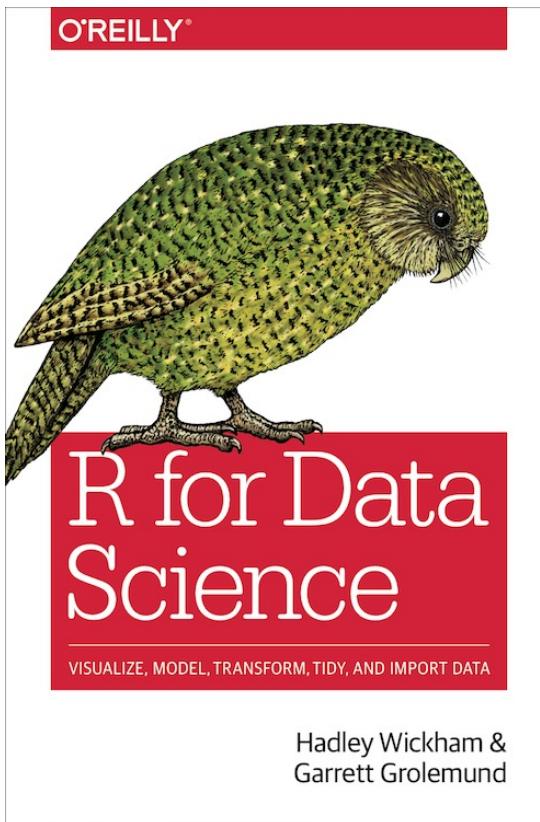
# The Future of Data Workflows

Seamless integration platforms, for example:

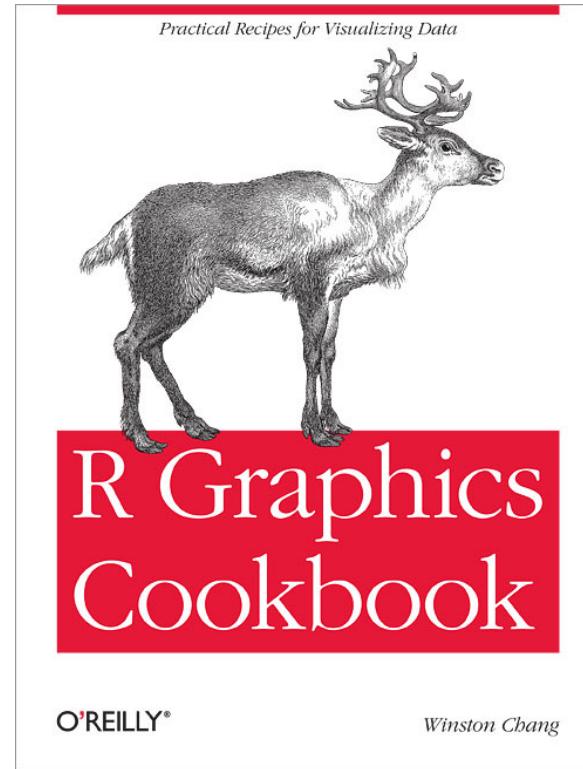


# **Long-term learning**

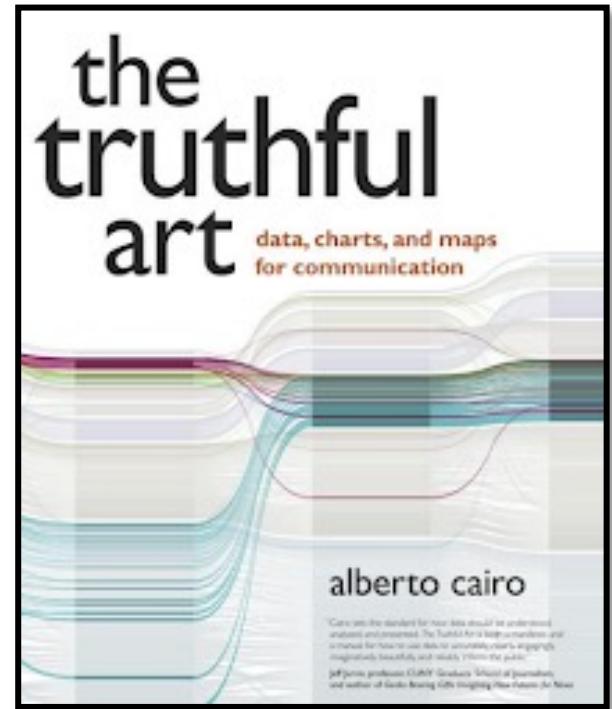
# Reading Recommendations



<https://r4ds.had.co.nz>

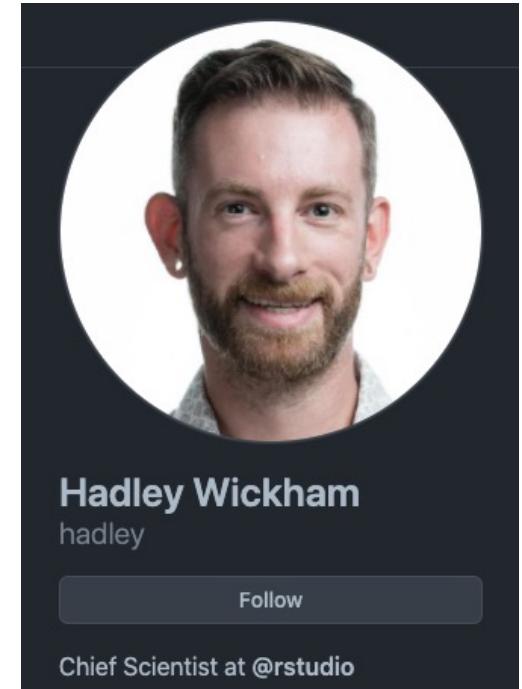


<https://r-graphics.org>



# Tutorial Recommendations

- Anything by Hadley Wickham
- <https://www.r-bloggers.com/>
- <https://www.r-statistics.com/>
- <https://blog.revolutionanalytics.com/>
- <https://r-charts.com/>
- <http://www.cookbook-r.com/Graphs/>
- <https://plotly.com/r/>



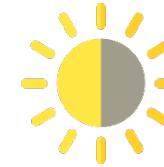
# Download Open Data

- Kaggle
  - <https://www.kaggle.com/datasets>
- Data.gov
  - <https://www.data.gov/>

# Download Cheatsheets

- Data Import
  - <https://github.com/rstudio/cheatsheets/raw/master/data-import.pdf>
- Data Wrangling Cheatsheet
  - <https://www.rstudio.com/wp-content/uploads/2015/02/data-wrangling-cheatsheet.pdf>
- Data Transformation with dplyr
  - <https://github.com/rstudio/cheatsheets/raw/master/data-visualization.pdf>
- String Manipulation
  - <https://github.com/rstudio/cheatsheets/raw/master/strings.pdf>
- Work with dates/times
  - <https://github.com/rstudio/cheatsheets/raw/master/lubridate.pdf>
- R Markdown
  - <https://github.com/rstudio/cheatsheets/raw/master/rmarkdown.pdf>
- More cheatsheets
  - <https://www.rstudio.com/resources/cheatsheets/>

# thank you



[nelson.roque@ucf.edu](mailto:nelson.roque@ucf.edu)



[nelsonroque.com](http://nelsonroque.com)



<https://github.com/nelsonroque>