

# Introduction to R

Day 1: Data Cleaning & Wrangling

Shirley Wang  
July 27, 2021

# Day 1 Goals

- Download & install R & RStudio
- Become familiar with the RStudio environment
- Operations
- Objects (what they are & how to create them)
- Data types
- Functions

# Day 1 Goals

- Download & install R & RStudio
- Become familiar with:
- Operations
- Objects (what they are)
- Data types
- Functions



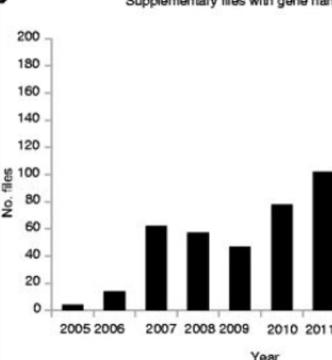
# Why use R?

Ethan Mollick  @emollick

A third of all genetics papers published in Nature over a decade (and 20% across all journals) had errors due to the fact that many gene have names like SEPT2 (the official name of Septin 2), which were automatically coded as dates by Microsoft Excel.

[genomebiology.biomedcentral.com/articles/10.11...](https://genomebiology.biomedcentral.com/articles/10.11...)

Supplementary files with gene name



8:23 PM · Mar 12, 2021 · Twitter for iPad

505 Retweets 181 Quote Tweets 1,450 Likes

Also an Excel error cost led to an estimated 1,500 deaths in the UK last year. The UK COVID contact tracing effort saved to an old Excel file format (.XLS) rather than the new one (.XLSX) & smaller row limits resulted in names being dropped. They were never contacted by tracers.

Does Contact Tracing Work?  
Quasi-Experimental Evidence from an  
Excel Error in England\*

Thiemo Fetzer<sup>†</sup> Thomas Graeber<sup>†</sup>  
November 24, 2020

Abstract

Contact tracing has been a central pillar of the public health response to the COVID-19 pandemic. Yet, contact tracing measures face substantive challenges in practice and well-identified evidence about their effectiveness remains scarce. This paper exploits quasi-random variation in COVID-19 contact tracing. Between September 25 and October 2, 2020, a total of 15,841 COVID-19 cases in England (around 15 to 20% of all cases) were not immediately referred to the contact tracing system due to a data processing error. Case information was truncated from an Excel spreadsheet after the row limit had been reached, which was discovered on October 3. There is substantial variation in the degree to which different parts of England areas were exposed – by chance – to delayed referrals of COVID-19 cases to the contact tracing system. We show that more affected areas subsequently experienced a drastic rise in new COVID-19 infections and deaths alongside an increase in the positivity rate and the number of test performed, as well as a decline in the performance of the contact tracing system. Conservative estimates suggest that the failure of timely contact tracing due to the data glitch is associated with more than 125,000 additional infections and over 1,500 additional COVID-19-related deaths. Our findings provide strong quasi-experimental evidence for the effectiveness of contact tracing.

Keywords: HEALTH, CORONAVIRUS  
JEL Classification: I31, Z18

Also an Excel error cost led to an estimated 1,500 deaths in the UK last year. The UK COVID contact tracing effort saved to an old Excel file format (.XLS) rather than the new one (.XLSX) & smaller row limits resulted in names being dropped. They were never contacted by tracers.

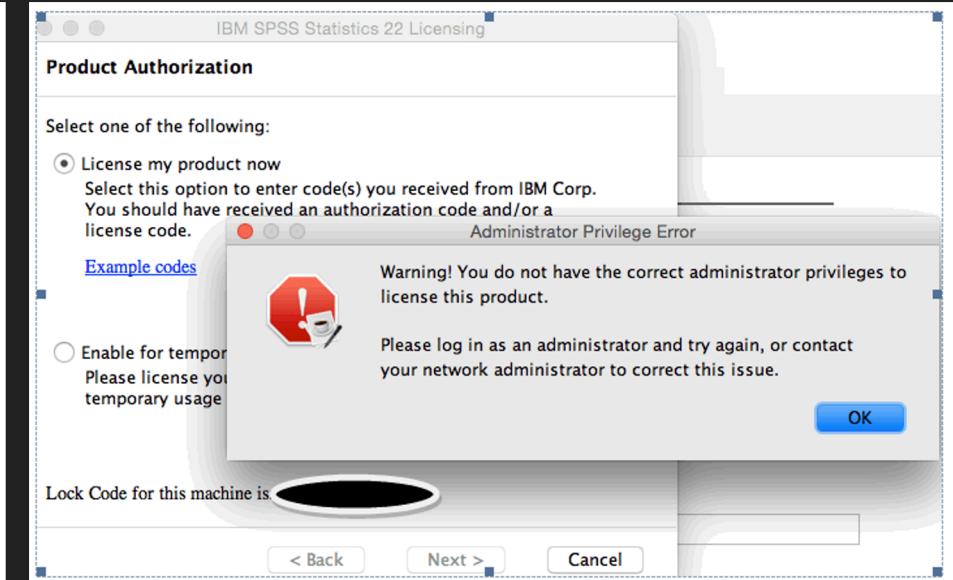
Does Contact Tracing Work?  
Quasi-Experimental Evidence from an  
Excel Error in England\*

Thiemo Fetzer<sup>†</sup> Thomas Graeber<sup>†</sup>  
November 24, 2020

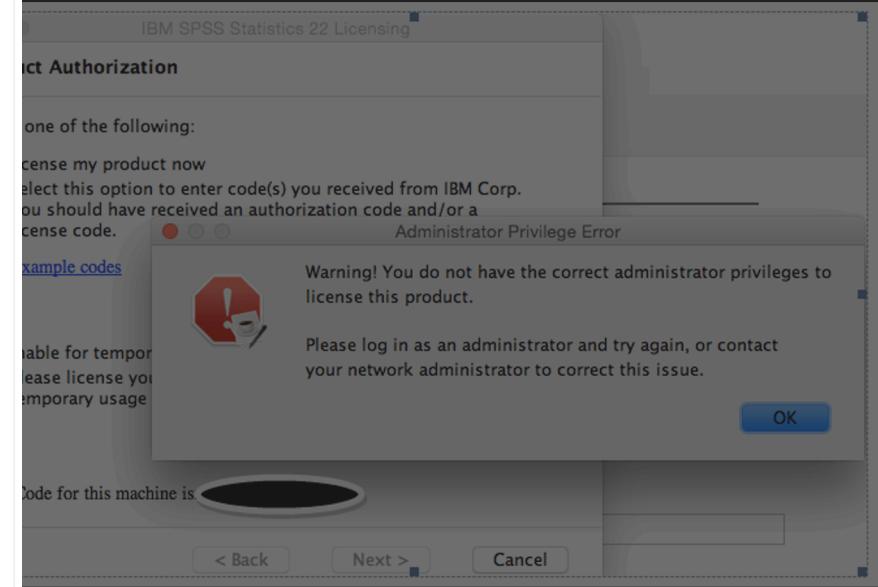
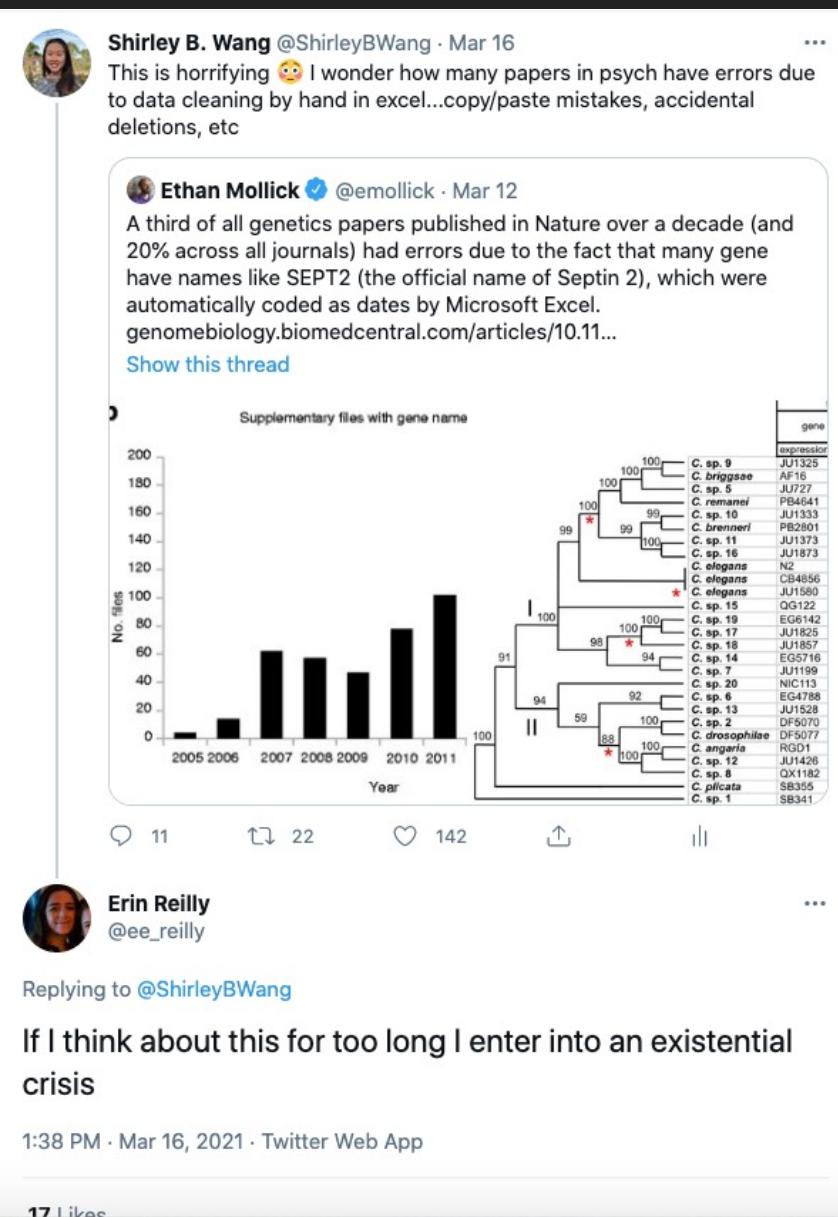
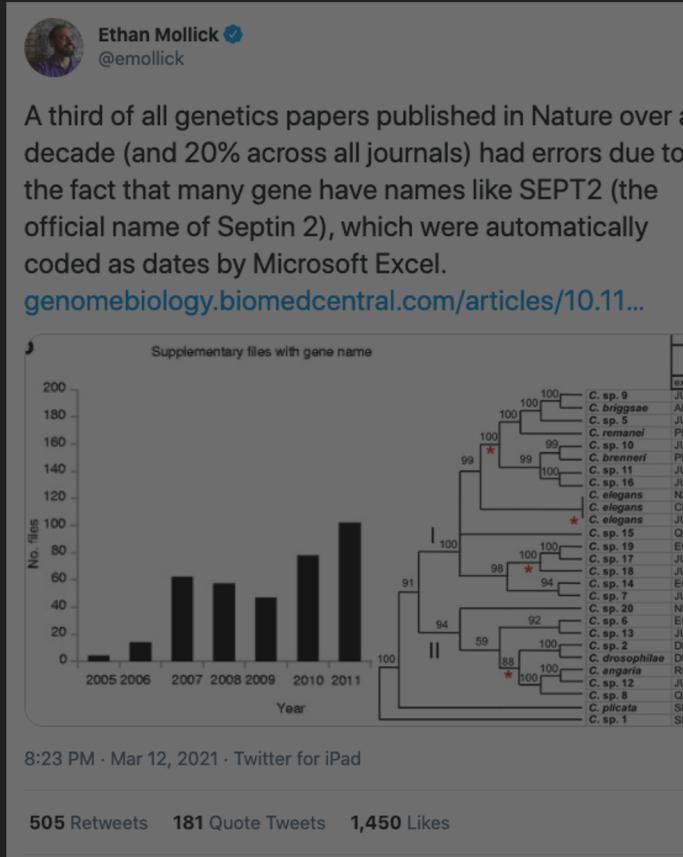
Abstract

Contact tracing has been a central pillar of the public health response to the COVID-19 pandemic. Yet, contact tracing measures face substantive challenges in practice and well-identified evidence about their effectiveness remains scarce. This paper exploits quasi-random variation in COVID-19 contact tracing. Between September 25 and October 2, 2020, a total of 15,841 COVID-19 cases in England (around 15 to 20% of all cases) were not immediately referred to the contact tracing system due to a data processing error. Case information was truncated from an Excel spreadsheet after the row limit had been reached, which was discovered on October 3. There is substantial variation in the degree to which different parts of England areas were exposed – by chance – to delayed referrals of COVID-19 cases to the contact tracing system. We show that more affected areas subsequently experienced a drastic rise in new COVID-19 infections and deaths alongside an increase in the positivity rate and the number of test performed, as well as a decline in the performance of the contact tracing system. Conservative estimates suggest that the failure of timely contact tracing due to the data glitch is associated with more than 125,000 additional infections and over 1,500 additional COVID-19-related deaths. Our findings provide strong quasi-experimental evidence for the effectiveness of contact tracing.

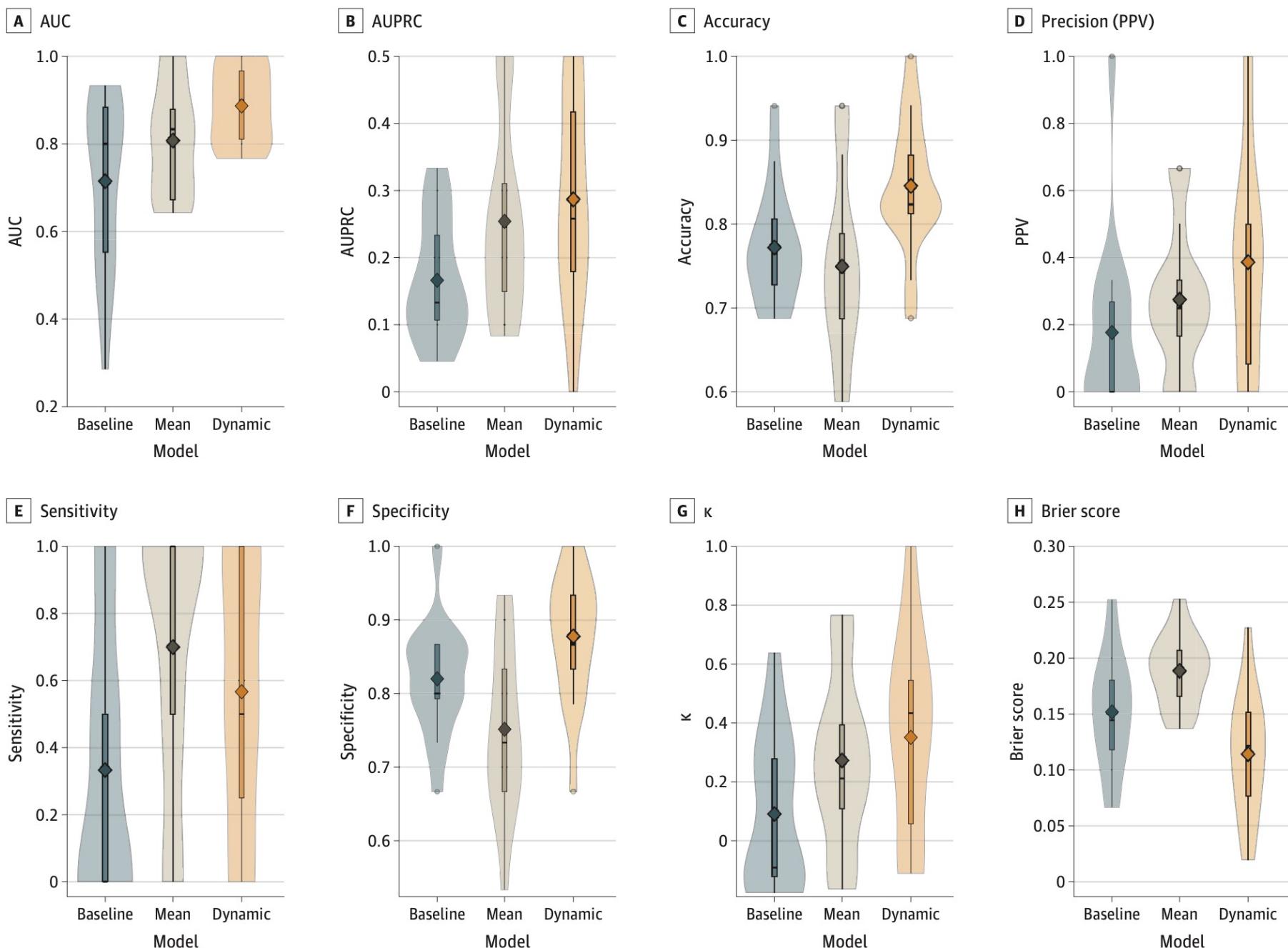
Keywords: HEALTH, CORONAVIRUS  
JEL Classification: I31, Z18



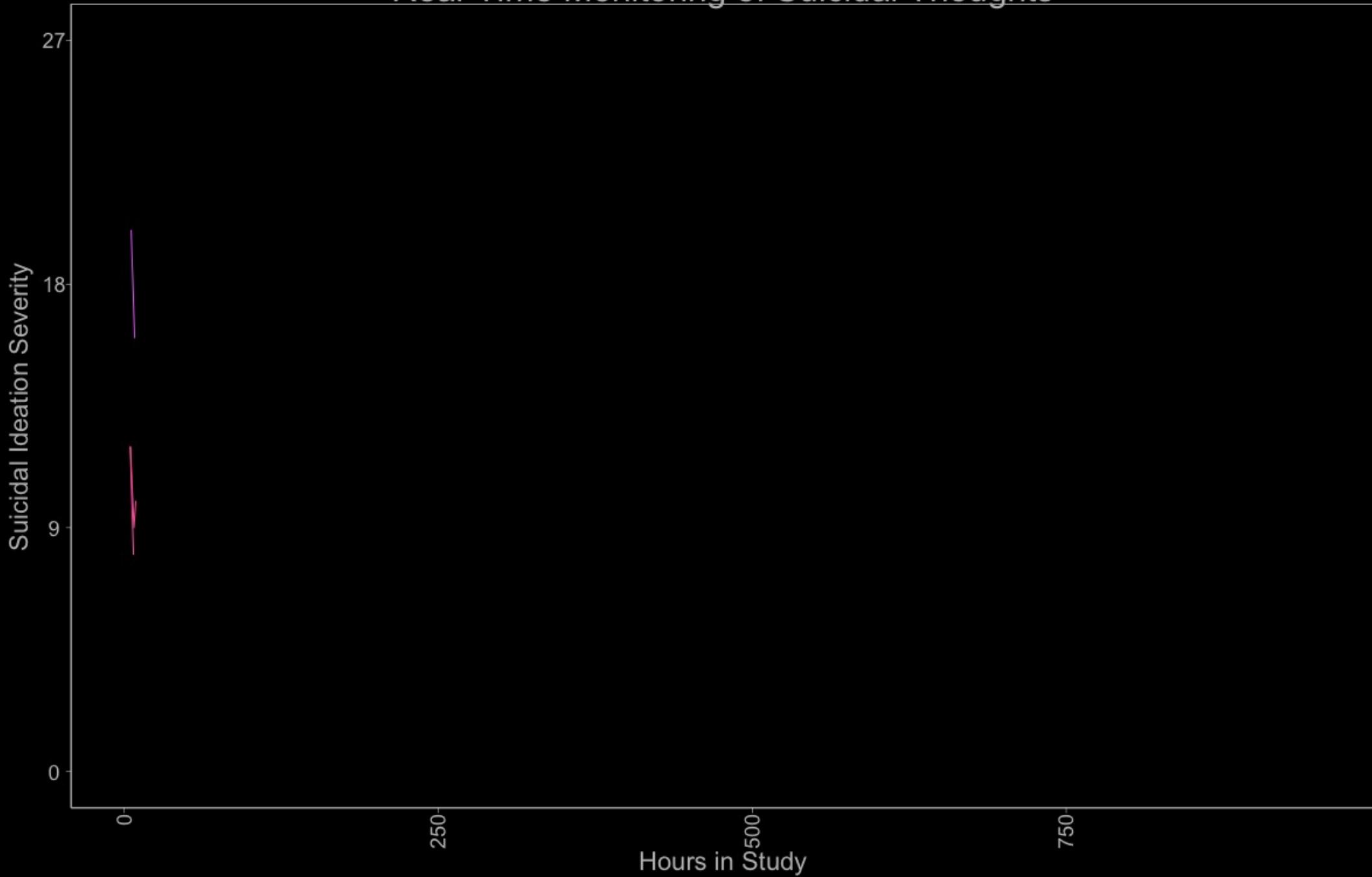
# Why use R?



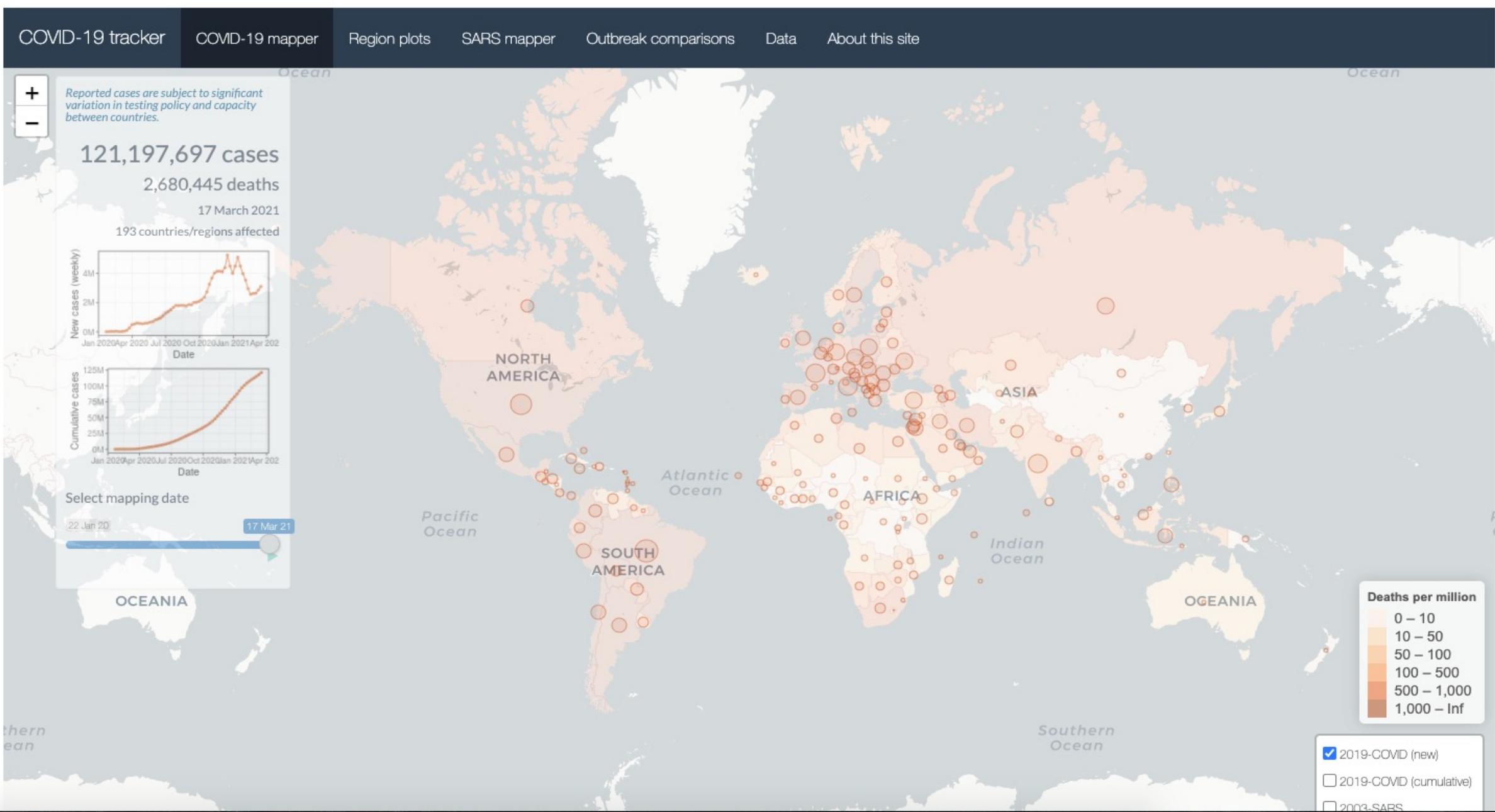
**Figure 1. Suicide Attempt Prediction Model Metrics**



# Real-Time Monitoring of Suicidal Thoughts



## covid19, epidemiology



[Introduction](#)[How It Works](#)[Code Chunks](#)[Inline Code](#)[Code Languages](#)[Parameters](#)[Tables](#)[Markdown Basics](#)[Output Formats](#)[Notebooks](#)[Slide Presentations](#)[Dashboards](#)[Websites](#)[Interactive Documents](#)[Cheatsheets](#)

# Slide Presentations

R Markdown renders to four presentation formats:

- `beamer_presentation` - PDF presentations with beamer
- `ioslides_presentation` - HTML presentations with ioslides
- `slidy_presentation` - HTML presentations with slidy
- `powerpoint_presentation` - PowerPoint presentation
- `revealjs::revealjs_presentation` - HTML presentations with reveal.js

Each format will intuitively divide your content into slides, with a new slide beginning at each first or second level header.

Insert a horizontal rule ( `***` ) into your document to create a manual slide break. Create incremental bullets with `>-`, as in the .Rmd file below, which is available [here](#) on RStudio Cloud.

The screenshot shows the RStudio interface with two panes. The left pane displays the R Markdown (.Rmd) source code, and the right pane shows the resulting HTML presentation.

**Left Pane (Source Code):**

```
1 ---  
2 title: "Viridis Presentation"  
3 output:  
4   revealjs::revealjs_presentation:  
5     theme: league  
6 ---  
7  
8 ```{r include = FALSE}  
9 knitr::opts_chunk$set(echo = FALSE)  
10 library(viridis)  
11 ```  
12  
13 The [viridis](https://github.com/sjmgarnier/viridis) package contains four color palettes, revealed in the plots that follow.  
14  
15 >- Viridis  
16 >- Magma  
17 >- Inferno  
18 >- Plasma  
19  
20 Each plot displays a contour map of the Maunga Whau volcano in Auckland, New Zealand.  
21  
22 ## Viridis colors  
23  
24 ```{r}
```

**Right Pane (HTML Output):**

The generated HTML page has a dark background with a large, semi-transparent watermark-like image of a volcano. The title "PLASMA COLORS" is displayed prominently at the top. Below the title, there is a caption: "Each plot displays a contour map of the Maunga Whau volcano in Auckland, New Zealand."



## Shirley Wang

Ph.D. candidate  
Harvard University



## About

I am a PhD candidate in clinical psychology with a secondary in computational science and engineering at Harvard University. I work with [Matt Nock](#) and am funded by the National Science Foundation Graduate Research Fellowship Program. My research examines why people engage in behaviors that are harmful to themselves, including eating disorder behaviors, nonsuicidal self-injury, and suicide. I use a wide range of methods to study these problems, including laboratory-based behavioral experiments, real-time monitoring, and large-scale longitudinal studies. I am particularly interested in using mathematical and computational modeling to formalize theories in psychopathology.

I value diversity, inclusion, and belonging, and am committed to promoting these values in research, academia, clinical practice, and beyond. I believe that increasing representation of people from marginalized and historically disadvantaged backgrounds is critical for conducting ethical, comprehensive, and innovative clinical science research.

### Interests

- eating disorders
- suicide
- computational modeling
- real-time monitoring

### Education

-  Ph.D. in Clinical Science, expected 2023  
Harvard University
-  A.M. in Clinical Science, 2019  
Harvard University

## Preface

- Why read this book
- Structure of the book
- Software information and convention...
- Acknowledgments

About the Author

## 1 Introduction

- 1.1 Motivation
- 1.2 Get started
- 1.3 Usage
- 1.4 Two rendering approaches
- 1.5 Some tips

## 2 Components

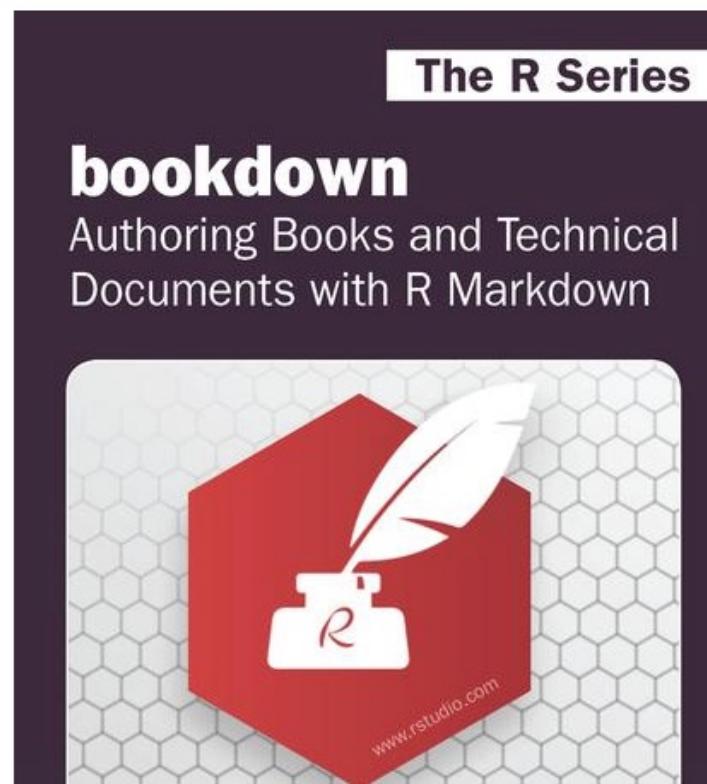
- 2.1 Markdown syntax
  - 2.1.1 Inline formatting
  - 2.1.2 Block-level elements
  - 2.1.3 Math expressions
- 2.2 Markdown extensions by bookdown
  - 2.2.1 Number and reference examples
  - 2.2.2 Theorems and proofs
  - 2.2.3 Special headers
  - 2.2.4 Text references

# bookdown: Authoring Books and Technical Documents with R Markdown

Yihui Xie

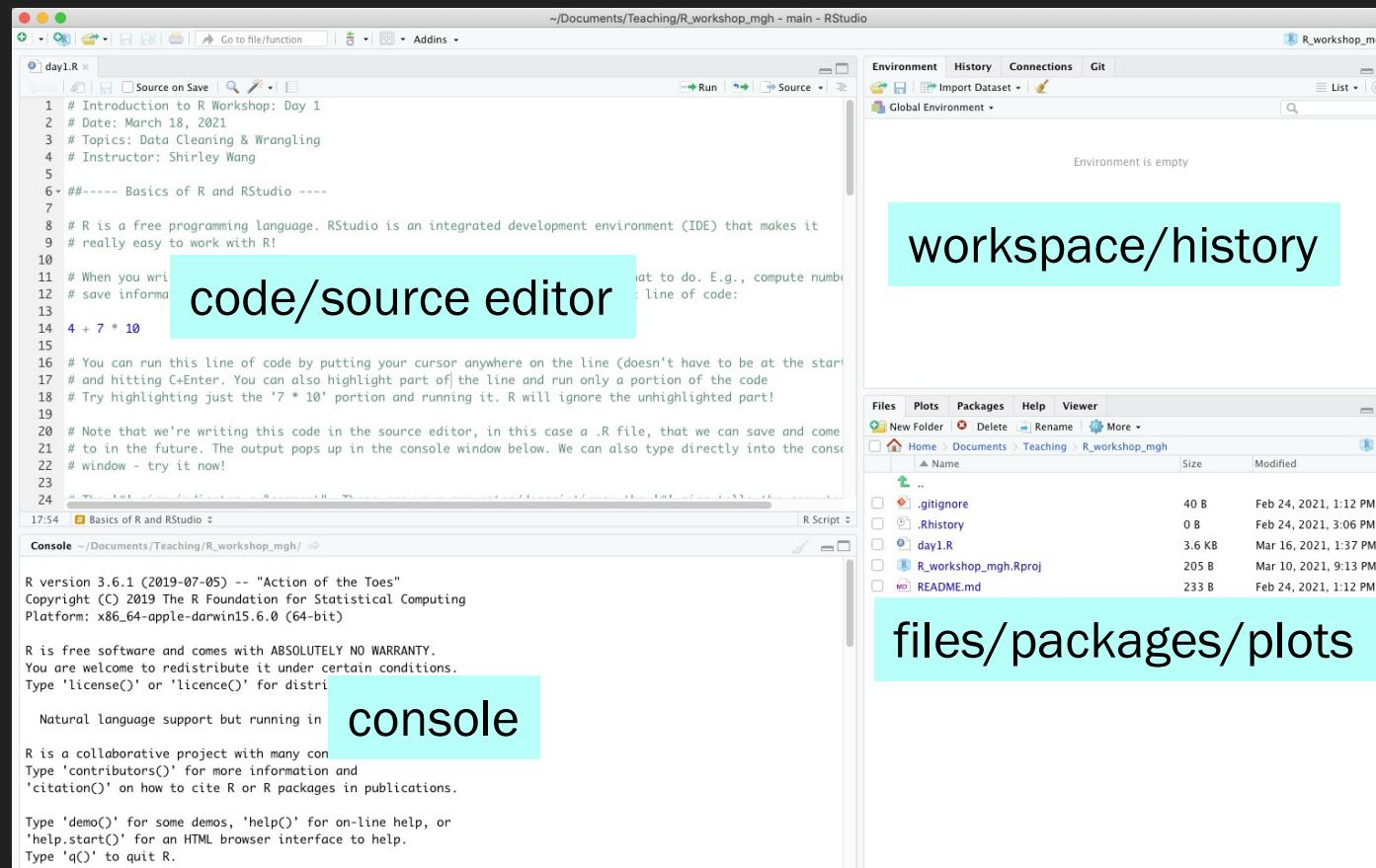
2021-03-15

## Preface



# R & RStudio Orientation

- R is a free programming language
- RStudio is an IDE that makes it easier write, run, and save R code!



# R Basics

- Case sensitive
- Run lines of code with cmd+Enter
- Get help with ? or `help()`
- Comment your code with #

# R Basics

- Case sensitive
- Run lines of code with cmd+Enter
- Get help with ? or
- Comment your code

When you trying to look at  
the code you wrote a month ago

IT'S SOME KIND OF ELVISH

I CAN'T READ IT

# Operations in R

## Arithmetic Operators

Operator	Description
+	addition
-	subtraction
*	multiplication
/	division
<sup>^ or </sup> **	exponentiation
x %% y	modulus (x mod y) 5%%2 is 1
x %/% y	integer division 5%/%2 is 2

## Logical Operators

Operator	Description
<	less than
<=	less than or equal to
>	greater than
>=	greater than or equal to
==	exactly equal to
!=	not equal to
!x	Not x
x   y	x OR y
x & y	x AND y
isTRUE(x)	test if X is TRUE

# Objects in R

- Objects store information (can be an entire dataset, a single value, a string of text, an image, etc)
- Assign values using ‘`<-`’
- Value on the right-hand side gets assigned to object name on the left-hand side (e.g., `myname <- 'shirley'`; `myage <- 26`)
- Rules for object names:
  - Must start with a letter
  - Can only contain letters, numbers, periods, underscores
  - Case-sensitive

# Data Types

- Numeric
- Integer
- Character
- Factor
- Logical
- Check data types in R with `class()`

# Data Structures

- Vector
- Matrix
- Array
- Data frame
- List

# Data Structures

- Vector
- Matrix
- Array
- Data frame
- List

~/Documents/Teaching/R\_workshop\_mgh - main - RStudio

	active	afraid	alert	angry	anxious	aroused	ashamed	astonished	at.ease	at.rest	attentive	blue	bored	calm
1	1	1	1	0	1	1	0	0	1	1	1	1	2	1
2	1	0	1	0	0	0	0	0	1	1	0	0	0	1
3	1	0	0	0	0	0	0	0	2	2	0	0	2	1
4	1	0	1	0	1	1	1	0	1	2	1	1	0	1
5	2	0	1	0	NA	2	0	3	3	1	1	0	1	3
6	2	0	1	0	NA	1	0	0	1	1	1	1	1	2
7	0	0	1	0	NA	0	0	0	1	1	1	3	2	2
8	0	0	0	0	NA	0	0	0	0	1	1	0	0	1
9	1	0	0	0	NA	1	0	0	2	0	0	1	0	0
10	0	0	2	0	NA	0	0	0	1	0	1	0	0	2
11	0	0	0	0	NA	0	0	0	2	0	0	0	2	2
12	1	0	1	0	NA	0	0	0	0	0	0	0	1	1
13	0	0	0	0	NA	0	0	0	0	0	0	0	1	0
14	2	0	1	0	NA	1	0	0	1	1	1	0	0	2
15	0	0	2	0	NA	0	0	0	3	3	3	0	0	3
16	0	0	0	0	NA	0	0	0	1	1	1	0	1	1
17	0	0	1	0	NA	1	0	0	1	1	1	0	1	1
18	3	0	2	0	NA	0	0	0	3	1	2	1	1	3
19	0	0	0	0	NA	0	0	0	0	0	0	0	0	0

Showing 1 to 20 of 3,896 entries, 92 total columns

# Data Frames

- Access columns with ‘\$’ (e.g., `df$age` accesses the column named ‘age’ in the data frame ‘df’)
- `View(df)` to view data frame in excel-like format
- `dim(df)` to get dimensions (# rows, # columns)
  - `nrow(df)`
  - `ncol(df)`
- `str(df)` to view structure of data frame
- `colnames(df)` to view all column (variable) names

# Functions

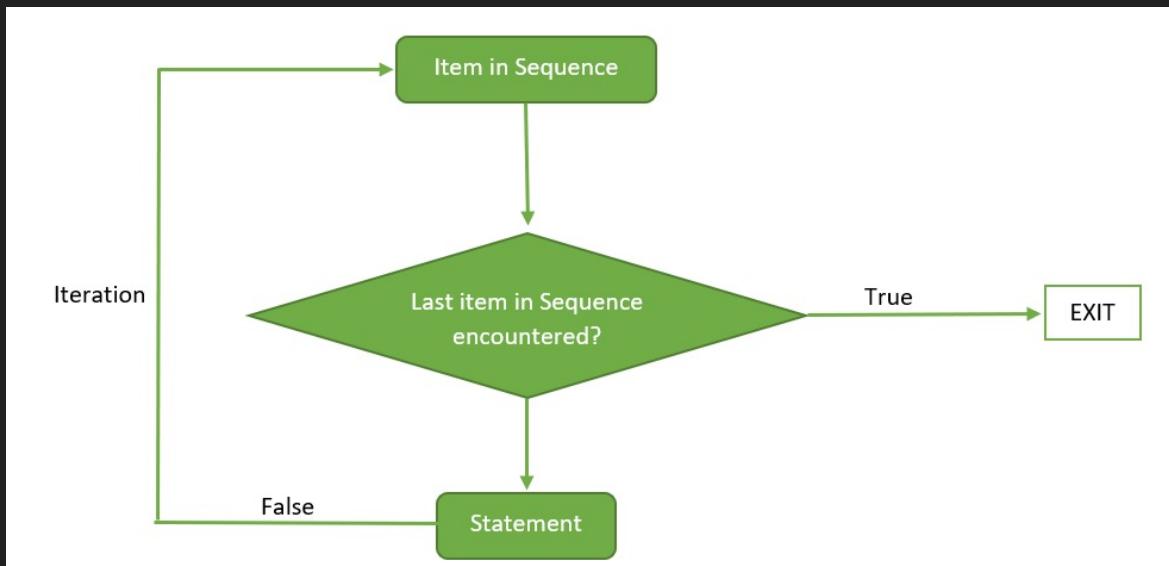
- Sets of instructions/statements to complete a specific task; takes input (e.g., numbers) and returns an output (e.g., mean of those numbers)
- R comes with many built-in functions (e.g., `mean()`, `min()`, `max()`, `print()`)
- You can also write your own functions to perform specific tasks!
- `function(arguments)`
- `myFunction <- function(input) {# do something; return output}`

# Packages

- R packages are bundles of functions, documentation, and data to perform tasks (e.g., machine learning, natural language processing, data visualization)
- Anyone can build an R package to share with other R users!
- Hosted on the central R repository CRAN: <https://cran.r-project.org/>
- Currently over 17,000 packages available
- Packages need to be installed on your computer, and loaded at the start of each new R session
- `install.packages('ggplot2')`
- `library(ggplot2)`

# Loops

- Loops are useful to repeat a block of code multiple times
- Rather than copy/pasting code several times, you can iterate through a loop! This can shorten hundreds of lines of code to just a few lines.



# Thank you!

## Helpful resources:

- <https://rstudio.com/resources/cheatsheets/>
- <https://www.r-bloggers.com/>
- <https://psyr.djnavarro.net/>
- <https://swirlstats.com/>
- <https://r4ds.had.co.nz/>
- <https://paulvanderlaken.com/2017/08/10/r-resources-cheatsheets-tutorials-books/>
- <https://www.learnr4free.com/en/index.html>
- [https://docs.google.com/document/d/1qtdiLbU32F\\_AVNRIF7d23wvPSgTCYqr5TFeBwbncPcw/edit](https://docs.google.com/document/d/1qtdiLbU32F_AVNRIF7d23wvPSgTCYqr5TFeBwbncPcw/edit)

