

# Data Wrangling with R and the Tidyverse

**Session 1**

Kieran Healy  
Statistical Horizons, April 2021

# Housekeeping

# Housekeeping

10am till 2pm US EST

# Housekeeping

10am till 2pm US EST

Lab session from 5pm to 6pm US EST

# Housekeeping

10am till 2pm US EST

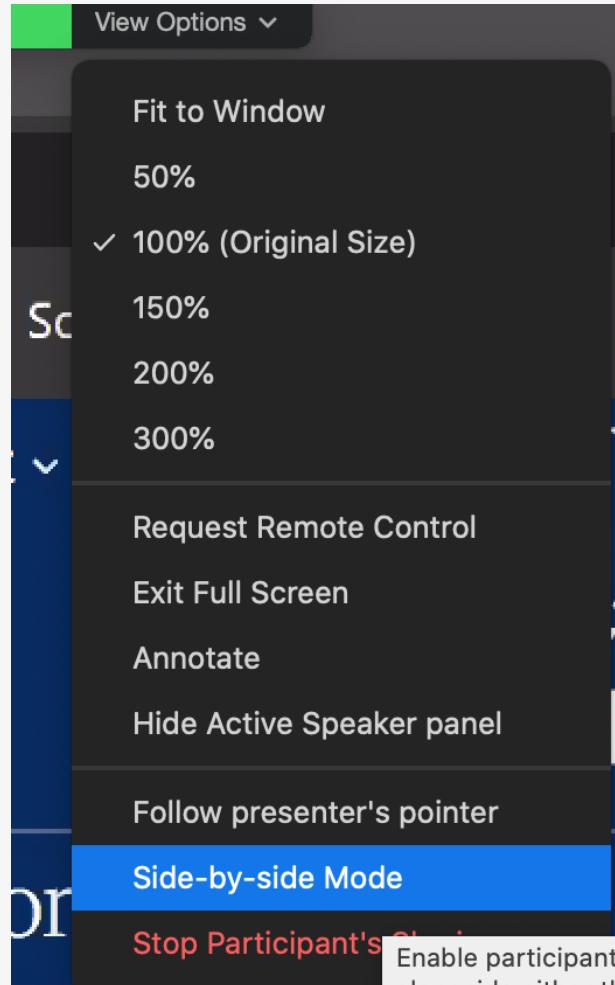
Lab session from 5pm to 6pm US EST

Use the Zoom chat to ask questions, or raise a hand with 

# In between class sessions



# For a better Zoom experience



If you're watching in full-screen view and I'm sharing my screen, then from Zoom's "View options" menu *turn off* "Side-by-Side" mode.

# This Morning's Goals

# This Morning's Goals

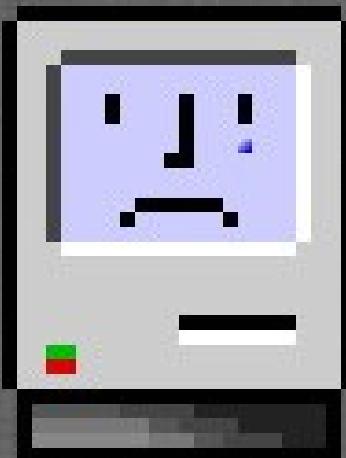
Some big-picture motivation & perspective

Getting familiar with RStudio and its relationship to R

Getting oriented to R and how it thinks

**DATA ANALYSIS**  
is mostly  
**DATA WRANGLING**

# Wrangling data is frustrating



# Can we make it **fun**?



# Can we make it **fun**?



No.

# Can we make it **fun**?



No.

⇒ Not *this* much fun, at any rate

# OK but can we eliminate frustration?



# OK but can we eliminate frustration?



Also no.

# OK but can we eliminate frustration?



Also no.

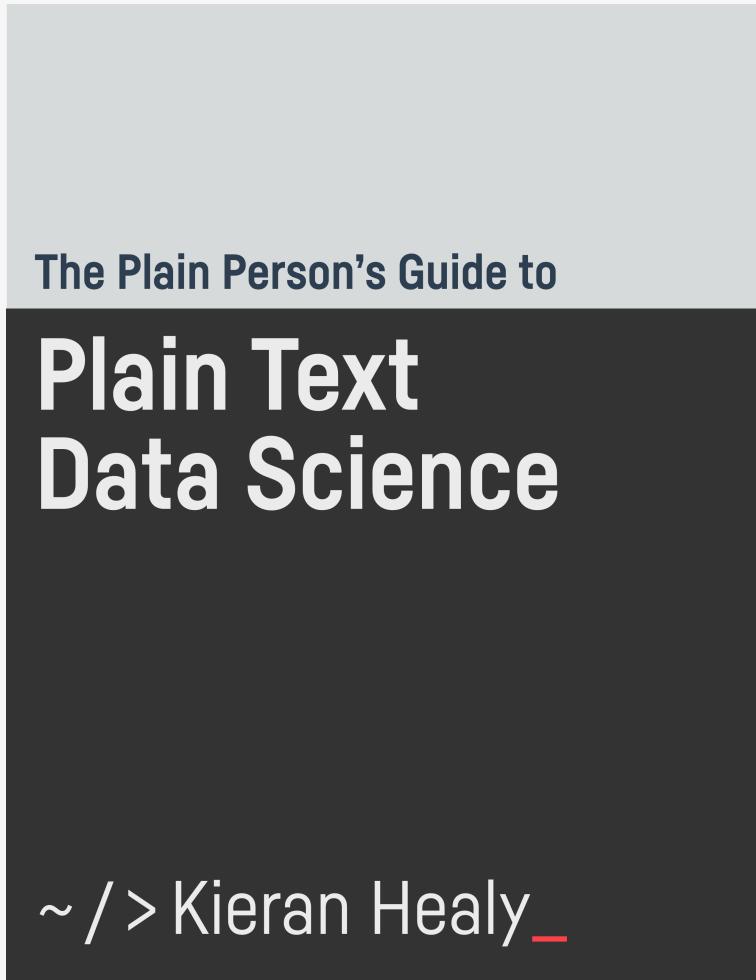
Sorry.

**HOWEVER, WE CAN  
MAKE IT *WORK***

# HOWEVER, WE CAN MAKE IT **WORK**

Also, it's weirdly satisfying once you get into it.

# We take a broadly *Plain Text* approach



# We take a broadly *Plain Text* approach

The Plain Person's Guide to

## Plain Text Data Science

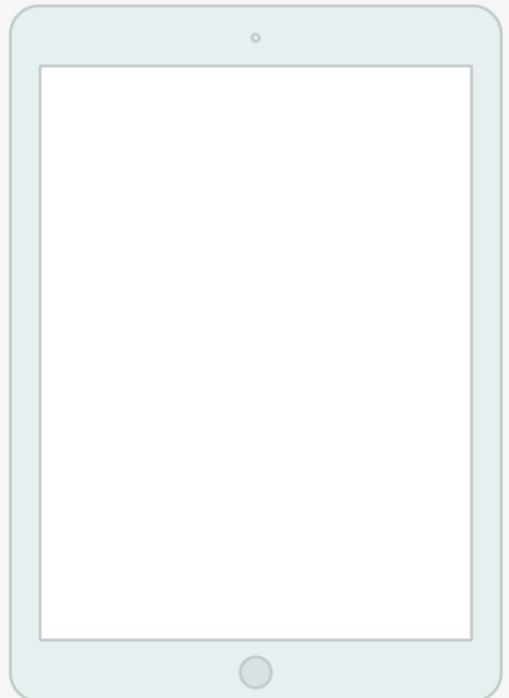
~ /> Kieran Healy \_

Using R and the Tidyverse can be understood within this broader context.

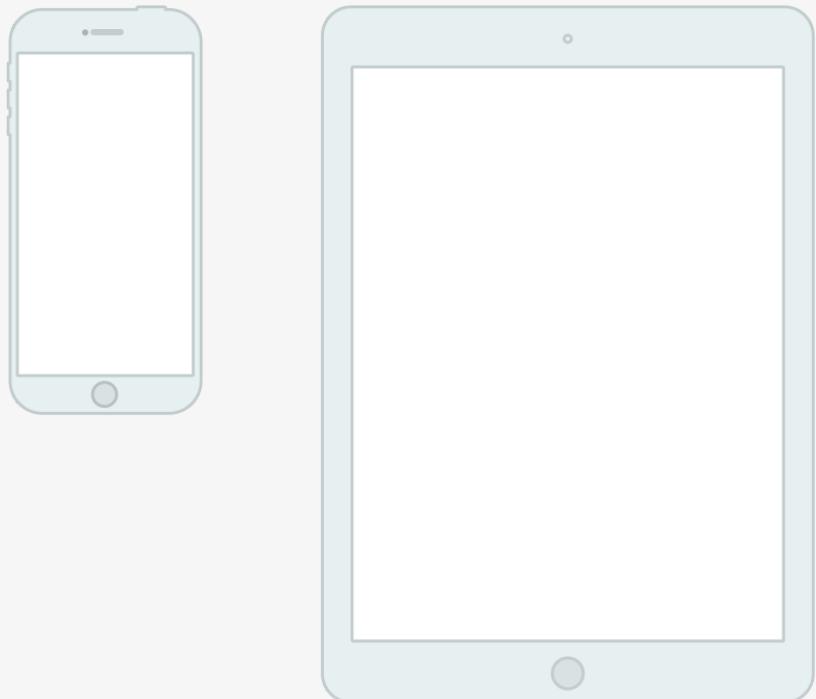
The same principles would apply to, e.g., using Python or similar tools.

# Two revolutions in computing

# Where the action is

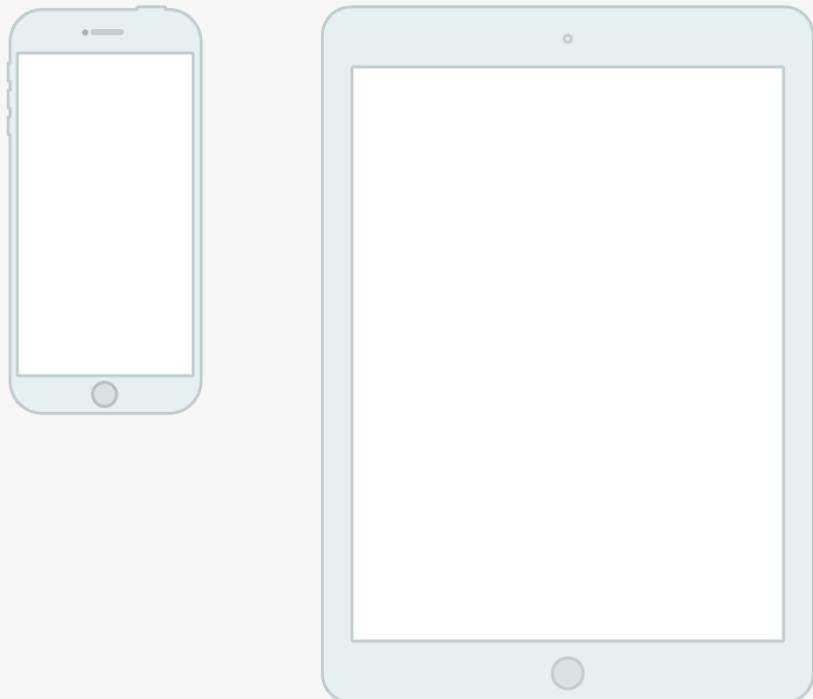


# Where the action is



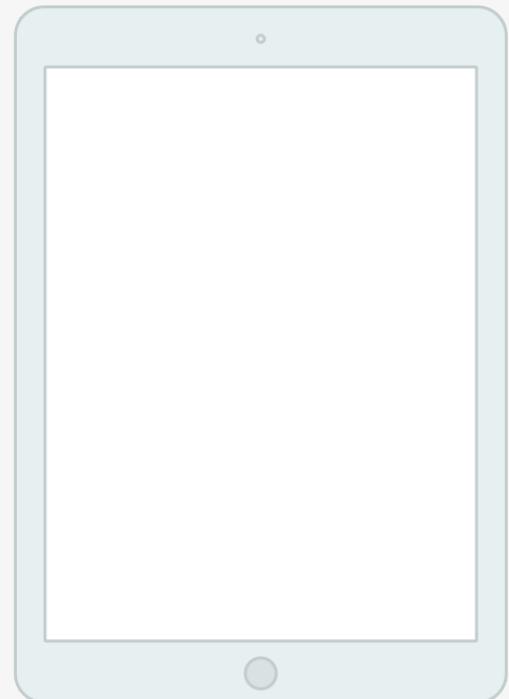
Touch-based user interface

# Where the action is



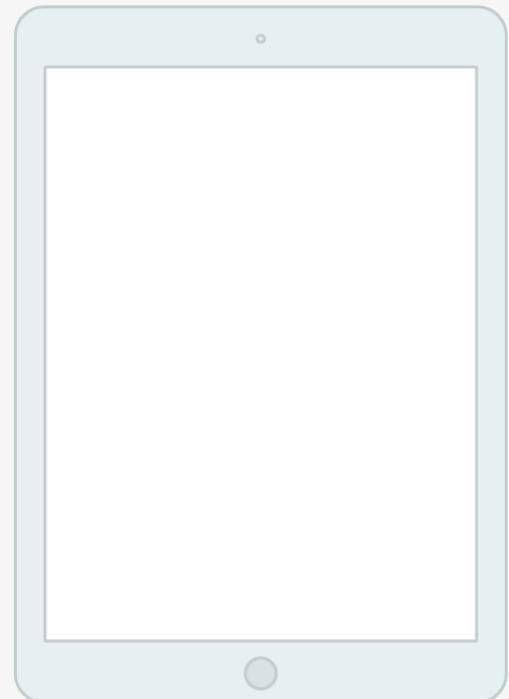
Touch-based user interface  
Foregrounds a single application

# Where the action is



Touch-based user interface  
Foregrounds a single application  
Dislikes multi-tasking\*

# Where the action is



Touch-based user interface  
Foregrounds a single application  
Dislikes multi-tasking\*  
Hides the file system

# \*Multitasking

I mean, “Making different specialized applications and resources work together in the service of a single but multi-dimensional project”, not “Checking Twitter while also listening to a talk and waiting for an update from the school nurse.”

# Where statistical computing lives



# Where statistical computing lives



Windows and pointers



# Where statistical computing lives



Windows and pointers

Multi-tasking, multiple windows



# Where statistical computing lives



Windows and pointers

Multi-tasking, multiple windows

Exposees and leverages the file system



# Where statistical computing lives



Windows and pointers

Multi-tasking, multiple windows

Exposees and leverages the file system

Many specialized tools in concert



# Where statistical computing lives



Windows and pointers

Multi-tasking, multiple windows

Exposees and leverages the file system

Many specialized tools in concert

Underneath, it's the 1970s, UNIX, and the command-line



# Plain-Text Tools for Data Analysis



# Plain-Text Tools for Data Analysis



Better than they've ever been!



# Plain-Text Tools for Data Analysis



Better than they've ever been!

Free! Open! Powerful!



# Plain-Text Tools for Data Analysis



Better than they've ever been!

Free! Open! Powerful!

Friendly community! Many resources!



# Plain-Text Tools for Data Analysis



Better than they've ever been!

Free! Open! Powerful!

Friendly community! Many resources!

But grounded in a UI paradigm that is increasingly far away from the everyday use of computing devices



# Plain-Text Tools for Data Analysis



Better than they've ever been!

Free! Open! Powerful!

Friendly community! Many resources!

But grounded in a UI paradigm that is increasingly far away from the everyday use of computing devices

So why do we use these tools?



The research process is  
*intrinsically messy*

# The research process is *intrinsically messy*

A rough distinction: "Office" vs "Engineering" approaches

# Questions

What is "real" in your project?

What is the final output?

How is it produced?

How are changes managed?

# Different Answers

## In the Office model

Formatted documents are real.

Intermediate outputs get pasted in.

Changes are tracked inside files.

Final output is often in the same format  
you've been working in, e.g. a Word file, or  
perhaps a PDF.

# Different Answers

## In the Office model

Formatted documents are real.

Intermediate outputs get pasted in.

Changes are tracked inside files.

Final output is often in the same format you've been working in, e.g. a Word file, or perhaps a PDF.

## In the Engineering model

Plain-text files are real.

Intermediate outputs are incorporated via code.

Changes are tracked outside files.

Final outputs are assembled programatically and converted to a desired output format.

# Different strengths and weaknesses

Everyone knows Word, Excel, or Google Docs.

"Track changes" is powerful and easy.

Hm, why can't I remember how I made this figure?

Where did this table of results come from?

Paper\_Submission\_Final\_edits\_FINAL\_khcomments-  
1.docx

# Different strengths and weaknesses

Everyone knows Word, Excel, or Google Docs.

"Track changes" is powerful and easy.

Hm, why can't I remember how I made this figure?

Where did this table of results come from?

Plain text is universally portable.

Push button, recreate analysis.

Why can't I make R do this simple thing?

This version control stuff is a pain.

Object of type 'closure' is not subsettable

Paper\_Submission\_Final\_edits\_FINAL\_khcomments-  
1.docx

# Different strengths and weaknesses

Everyone knows Word, Excel, or Google Docs.

"Track changes" is powerful and easy.

Hm, why can't I remember how I made this figure?

Where did this table of results come from?

Plain text is universally portable.

Push button, recreate analysis.

Why can't I make R do this simple thing?

This version control stuff is a pain.

Object of type 'closure' is not subsettable

Paper\_Submission\_Final\_edits\_FINAL\_khcomments-  
1.docx

**Each approach generates solutions to its own problems, too**

# INTO THE KITCHEN



# RStudio is an IDE for R



# An IDE for Meals



# R & RStudio

The screenshot shows the RStudio interface with the following components:

- Top Bar:** RStudio, File, Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools, Window, Help.
- Project Bar:** covdata - main - RStudio, Go to file/function, Addins.
- Code Editor (Left):** covdata.Rmd (R Markdown file). The code includes comments about loading the covdata package and its purpose to make COVID-19 data accessible.
- Environment (Top Right):** Shows the Global Environment and Functions pane.
- Console (Bottom Left):** Displays R startup messages, package loading, and a testthat message.
- Files (Bottom Right):** Shows the file structure of the covdata project, including .github, .gitignore, .Rbuildignore, .Rhistory, \_pkgdown.yml, \_sinewconfig.yml, covdata.Rproj, data, data-raw, DESCRIPTION, inst, LICENSE, LICENSE.md, man, and NAMESPACE files.

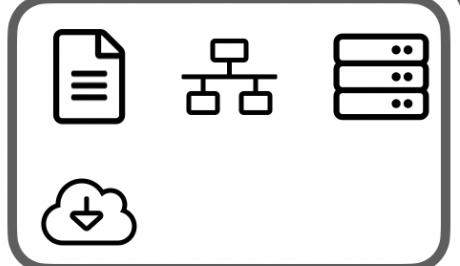
# R & RStudio

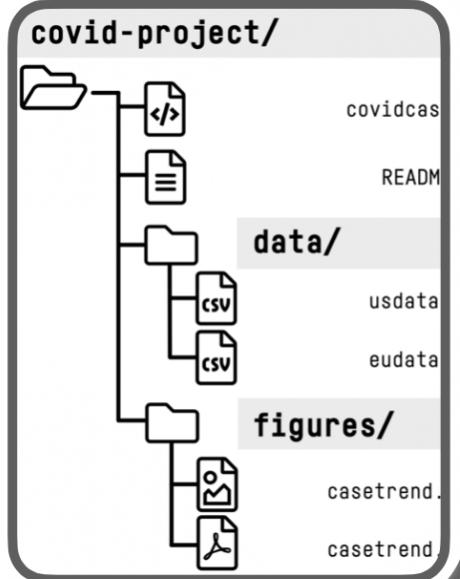
```
# COVID      covidcases.Rmd

## Get data from ECDC
```{r get-data}
covid_raw <- get_ecdc[url]
```

## Get data from the US
```{r get-data}
us_raw <- get_us[url]
```

```

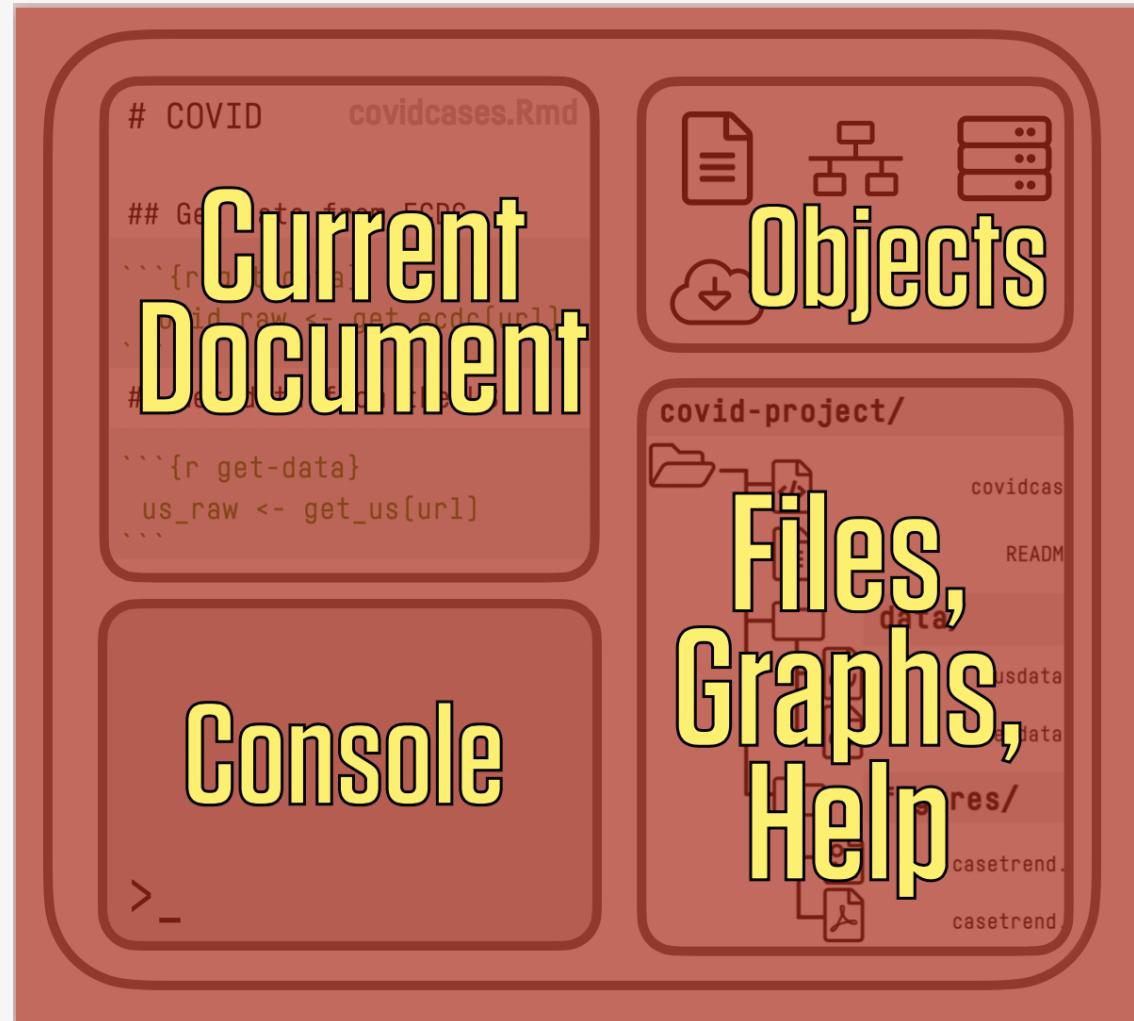




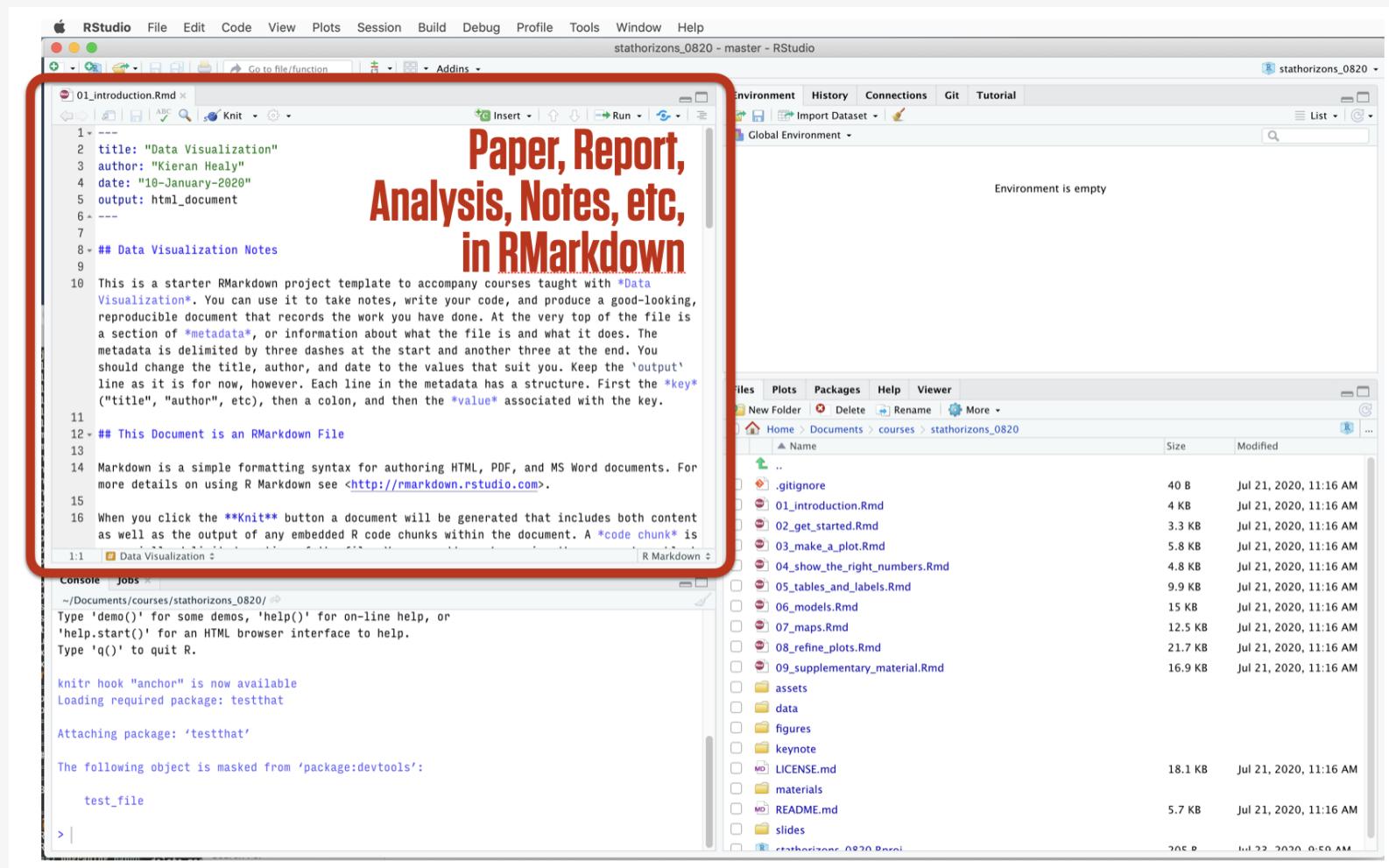
The diagram shows a file tree for a COVID project:

- covid-project/**
  - covidcas** (empty folder)
  - README** (text file)
  - data/** (empty folder)
    - usdata** (empty folder)
      - eudata** (empty folder)
  - figures/** (empty folder)
    - casetrend.** (empty folder)
      - casetrend** (empty folder)

# R & RStudio



# RStudio



# R & RStudio

The screenshot shows the RStudio IDE interface. On the left, the code editor displays a R Markdown file named '01\_introduction.Rmd'. The file contains metadata at the top and then text explaining R Markdown syntax. A red box highlights the 'Console' tab, which shows the R command-line interface with various package loading and help messages. The right side of the interface includes the Environment pane (which is currently empty), the File Browser (listing files in the 'stathorizons\_0820' directory), and the Global Environment pane.

Console: Type or send code here, see results

```
1 ---  
2 title: "Data Visualization"  
3 author: "Kieran Healy"  
4 date: "10-January-2020"  
5 output: html_document  
6 ---  
7  
8 ## Data Visualization Notes  
9  
10 This is a starter RMarkdown project template to accompany courses taught with \*Data Visualization\*. You can use it to take notes, write your code, and produce a good-looking, reproducible document that records the work you have done. At the very top of the file is a section of \*metadata\*, or information about what the file is and what it does. The metadata is delimited by three dashes at the start and another three at the end. You should change the title, author, and date to the values that suit you. Keep the 'output' line as it is for now, however. Each line in the metadata has a structure. First the \*key\* ("title", "author", etc), then a colon, and then the \*value\* associated with the key.  
11  
12 ## This Document is an RMarkdown File  
13  
14 Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see <http://rmarkdown.rstudio.com>.  
15  
16 When you click the \*\*Knit\*\* button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document. A \*code chunk\* is  
17  
18 Data Visualization
```

Console: Type or send code here, see results

# R & RStudio

The screenshot shows the RStudio interface with a project titled "stathorizons\_0820".

- Code Editor:** Displays the file "01\_introduction.Rmd" containing R Markdown code. The code includes metadata (title, author, date, output) and a section for notes.
- Environment:** Shows the Global Environment pane which is currently empty.
- Console:** Displays R session output related to package loading and help documentation.
- Files:** A file browser pane showing the project directory structure. A red box highlights this pane. It lists files like .gitignore, 01\_introduction.Rmd, 02\_get\_started.Rmd, etc., along with their sizes and modification dates.
- Plots, Packages, Help, Viewer:** Buttons for navigating between different RStudio components.

**Project files, Plots, Help**

| Name                          | Size    | Modified               |
|-------------------------------|---------|------------------------|
| ..                            | 40 B    | Jul 21, 2020, 11:16 AM |
| .gitignore                    | 4 KB    | Jul 21, 2020, 11:16 AM |
| 01_introduction.Rmd           | 3.3 KB  | Jul 21, 2020, 11:16 AM |
| 02_get_started.Rmd            | 5.8 KB  | Jul 21, 2020, 11:16 AM |
| 03_make_a_plot.Rmd            | 4.8 KB  | Jul 21, 2020, 11:16 AM |
| 04_show_the_right_numbers.Rmd | 9.9 KB  | Jul 21, 2020, 11:16 AM |
| 05_tables_and_labels.Rmd      | 15 KB   | Jul 21, 2020, 11:16 AM |
| 06_models.Rmd                 | 12.5 KB | Jul 21, 2020, 11:16 AM |
| 07_maps.Rmd                   | 21.7 KB | Jul 21, 2020, 11:16 AM |
| 08_refine_plots.Rmd           | 16.9 KB | Jul 21, 2020, 11:16 AM |
| 09_supplementary_material.Rmd | 18.1 KB | Jul 21, 2020, 11:16 AM |
| assets                        | 5.7 KB  | Jul 21, 2020, 11:16 AM |
| data                          | 205 B   | Jul 22, 2020, 0:50 AM  |
| figures                       |         |                        |
| keynote                       |         |                        |
| LICENSE.md                    |         |                        |
| materials                     |         |                        |
| README.md                     |         |                        |
| slides                        |         |                        |
| stathorizons_0820.Rproj       |         |                        |

# R & RStudio

The screenshot shows the RStudio interface with several panes:

- Code Pane:** Displays the file `01_introduction.Rmd` containing R Markdown code.
- Environment Pane (highlighted by a red box):** Shows the global environment, which is currently empty.
- File Explorer:** Shows the project structure with files like `.gitignore`, `01_introduction.Rmd`, and `02_get_started.Rmd`.
- Console Pane:** Displays R session output, including help documentation and package loading.

A red box highlights the Environment pane, and the text "Inspect objects you create" is overlaid in red.

**Environment pane content:**

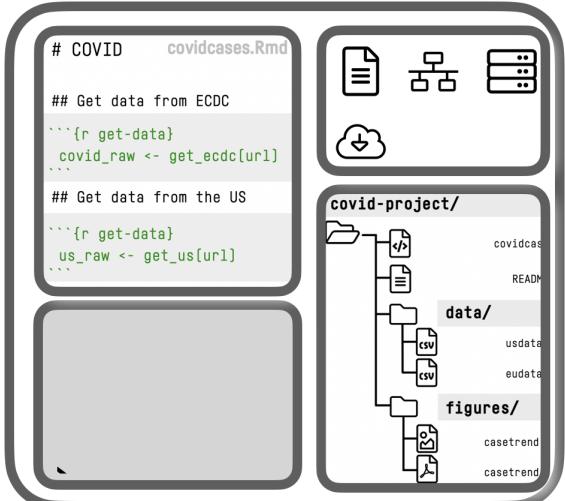
Environment is empty

**File Explorer content:**

| Name                          | Size    | Modified               |
|-------------------------------|---------|------------------------|
| ..                            |         |                        |
| .gitignore                    | 40 B    | Jul 21, 2020, 11:16 AM |
| 01_introduction.Rmd           | 4 KB    | Jul 21, 2020, 11:16 AM |
| 02_get_started.Rmd            | 3.3 KB  | Jul 21, 2020, 11:16 AM |
| 03_make_a_plot.Rmd            | 5.8 KB  | Jul 21, 2020, 11:16 AM |
| 04_show_the_right_numbers.Rmd | 4.8 KB  | Jul 21, 2020, 11:16 AM |
| 05_tables_and_labels.Rmd      | 9.9 KB  | Jul 21, 2020, 11:16 AM |
| 06_models.Rmd                 | 15 KB   | Jul 21, 2020, 11:16 AM |
| 07_maps.Rmd                   | 12.5 KB | Jul 21, 2020, 11:16 AM |
| 08_refine_plots.Rmd           | 21.7 KB | Jul 21, 2020, 11:16 AM |
| 09_supplementary_material.Rmd | 16.9 KB | Jul 21, 2020, 11:16 AM |
| assets                        |         |                        |
| data                          |         |                        |
| figures                       |         |                        |
| keynote                       |         |                        |
| LICENSE.md                    | 18.1 KB | Jul 21, 2020, 11:16 AM |
| materials                     |         |                        |
| README.md                     | 5.7 KB  | Jul 21, 2020, 11:16 AM |
| slides                        |         |                        |
| stathorizons_0820.Rproj       | 205 B   | Jul 22, 2020, 8:50 AM  |

# R & RStudio

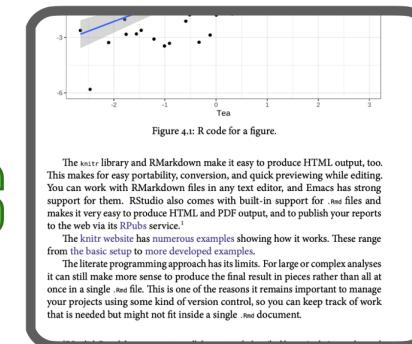
## RStudio



## Drive R

```
if (is.empty.model(mt)) {  
  x <- NULL  
  z <- list(coefficients = if (mlm)  
    matrix(NA_real_, 0,  
    ncol(y)) else numeric(), residuals = y,  
    fitted.values = y *  
    y, weights = w, rank = 0L, df.residual = if  
(!is.null(w)) sum(w !=  
  0) else ny)
```

## Generate Documents



The knitr library and RMarkdown make it easy to produce HTML output, too. This makes for easy readability, conversion and quick previewing while editing. You can work with RMarkdown files in any text editor, and Emacs has strong support for them. RStudio also comes with built-in support for `.Rmd` files and makes it very easy to produce HTML and PDF output, and to publish your reports to the web via its RPubs service.<sup>1</sup>

The Knitr website has numerous examples showing how it works. These range from the basic setup to more developed examples.

The literate programming approach has its limits. For large or complex analyses it can still make more sense to produce the final result in pieces rather than all at once in a single `.Rmd` file. This is one of the reasons it remains important to manage your projects using some kind of version control, so you can keep track of work that is needed but might not fit inside a single `.Rmd` document.

## View & Manage Environment

| Name                               | Type                              | Value                                 |
|------------------------------------|-----------------------------------|---------------------------------------|
| <input checked="" type="radio"/> p | list [9] (S3: gg, ggplot)         | List of length 9                      |
| <input type="radio"/> data         | list [1704 x 6] (S3: tbl_df, tbl) | A tibble with 1704 rows and 6 columns |
| <input type="radio"/> layers       | list [0]                          | List of length 0                      |
| <input type="radio"/> scales       | environment [1] (S3: ScalesList)  | <environment: 0x7f0f08c1e010>         |
| <input type="radio"/> mapping      | list [3] (S3: unval)              | List of length 3                      |
| <input type="radio"/> theme        | list [0]                          | List of length 0                      |
| <input type="radio"/> coordinates  | environment [5] (S3: CoordCa)     | <environment: 0x7f0f08c27b40>         |
| <input type="radio"/> facet        | environment [2] (S3: FacetNul)    | <environment: 0x7f0f08c55210>         |
| <input type="radio"/> plot_env     | environment [6]                   | <environment: R_GlobalEnv>            |
| <input type="radio"/> labels       | list [3]                          | List of length 3                      |