

# Overview

*Data Wrangling, Session 1*

Kieran Healy  
Code Horizons

July 22, 2024

# Housekeeping

**10:30am till 12:30pm US EST each day**

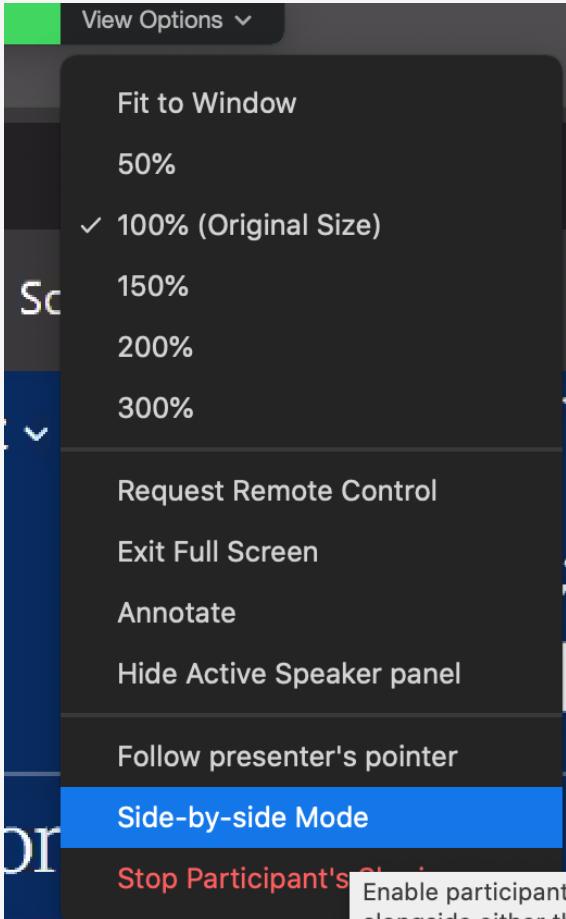
**1:30pm to 3:00pm US EST each day**

**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.

# My Setup and Yours

**Talking, Slides, and Live-Coding in RStudio**

**Follow along with RStudio yourself if you can**

**The course packet is also an RStudio project and  
the place for your notes**

# Goals for this first session

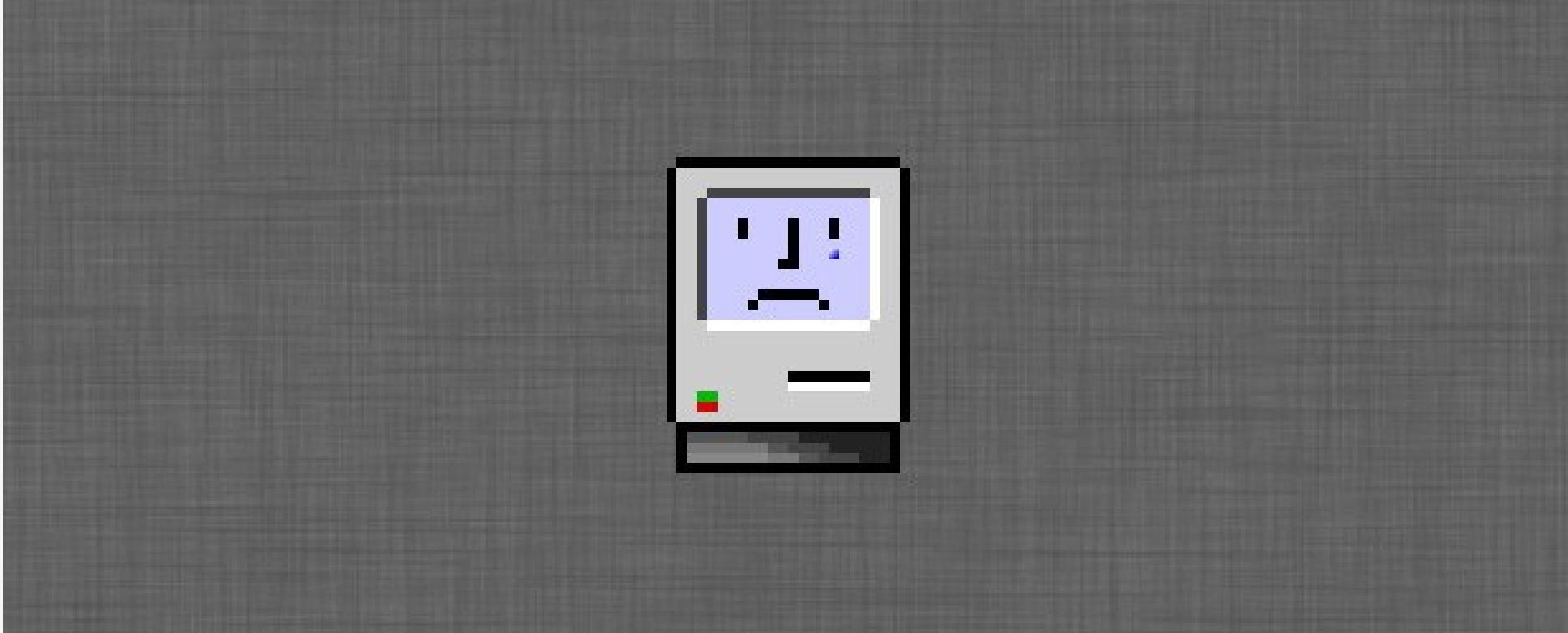
Some big-picture motivation & perspective

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



Sad Mac

# Can we make it **fun**?



No.

Fun data wrangling

# Can we make it **fun**?



Fun data wrangling

**No.**

↔ Not *this* much fun, at  
any rate

# OK but can we eliminate frustration?



Also no.

Frustration-free data wrangling

# OK but can we eliminate frustration?



Frustration-free data wrangling

**Also no.**  
(Sorry.)

However, we *can*  
make it **work**

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

# We take a broadly *Plain Text* approach

The Plain Person's Guide to  
**Plain Text**  
**Data Science**

~ / > Kieran Healy \_

The plain person's guide

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



iPhone and iPad

Touch-based user interface

Foregrounds a single application

Dislikes multi-tasking\*

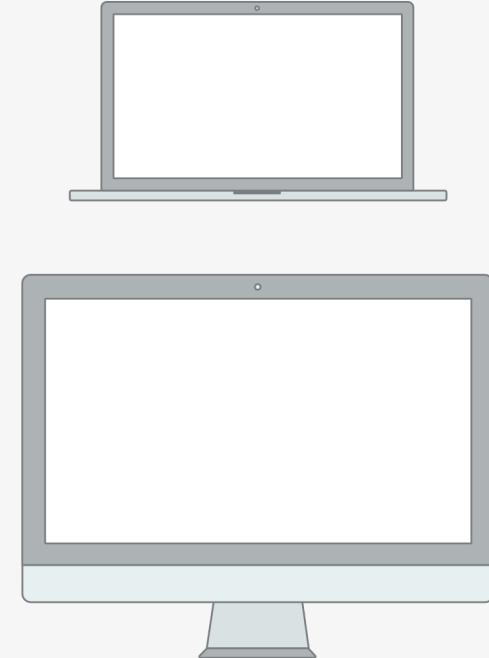
Hides the file system

“Laundry basket” model of where  
things are

# \*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



Desktop and laptop

Windows and pointers.

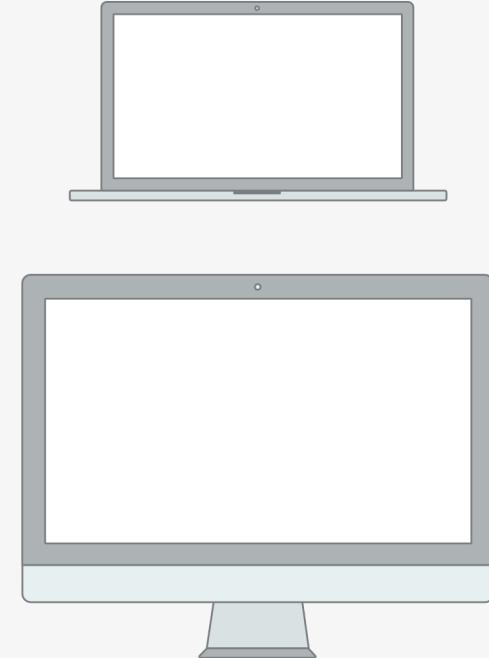
Multi-tasking, multiple windows.

Exposes 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



Desktop and laptop

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?

Because the  
research process  
is *intrinsically*  
*messy*

# “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

## Office model

Formatted documents are real.

Intermediate outputs are cut and pasted into documents.

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

## Office model

Formatted documents are real.

Intermediate outputs are cut and pasted into documents.

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.

## Engineering model

Plain-text files are real.

Intermediate outputs are produced via code, often inside documents.

Changes are tracked outside files.

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

# Different strengths and weaknesses

## Office model

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\_Final\_edits\_FINAL\_kh-1a.docx

# Different strengths and weaknesses

## Office model

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\_Final\_edits\_FINAL\_kh-1a.docx

## Engineering model

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

Each approach generates solutions to  
its own problems

# Into the Kitchen



Studio<sup>®</sup>

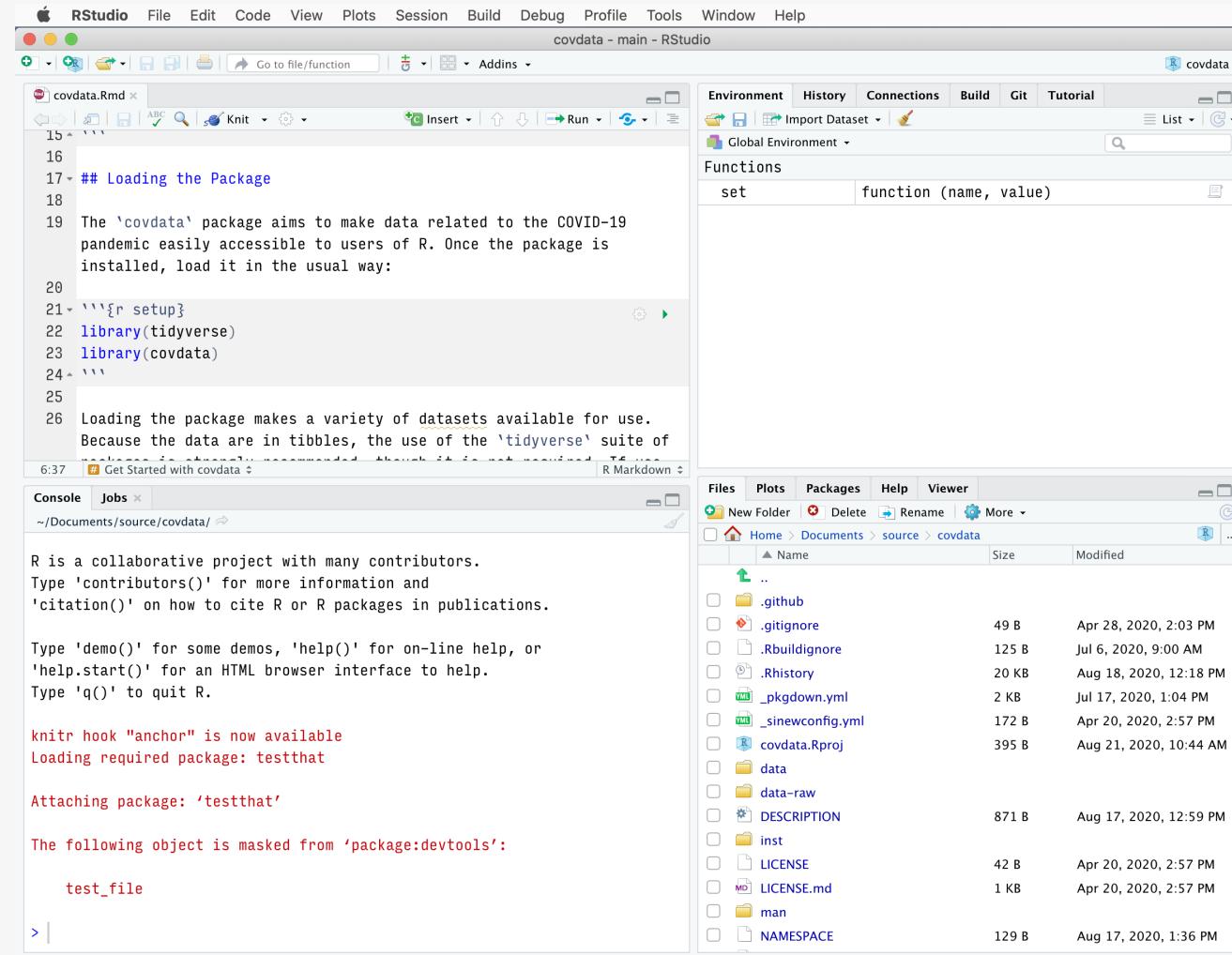
# RStudio is an IDE for R



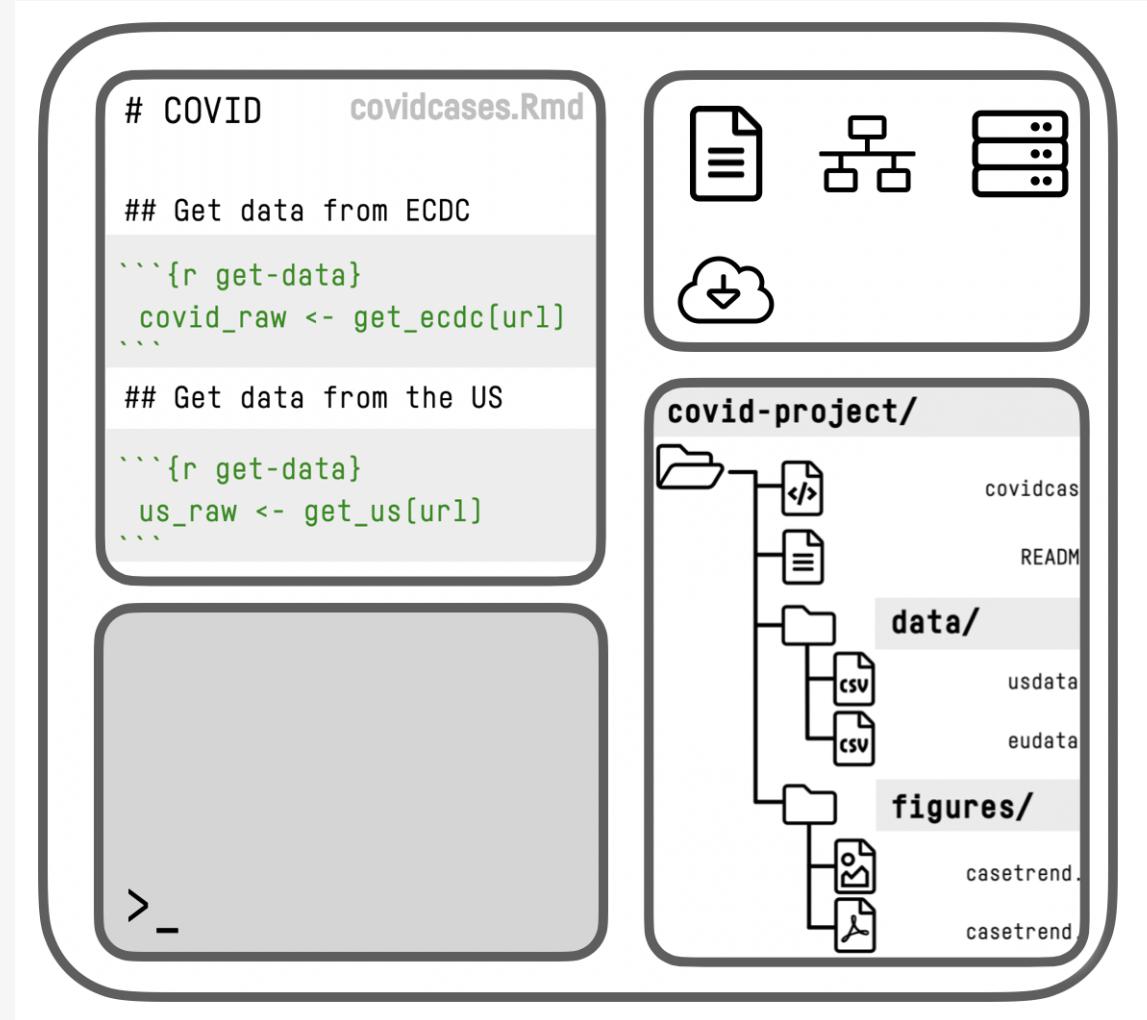
# A kitchen is an IDE for Meals



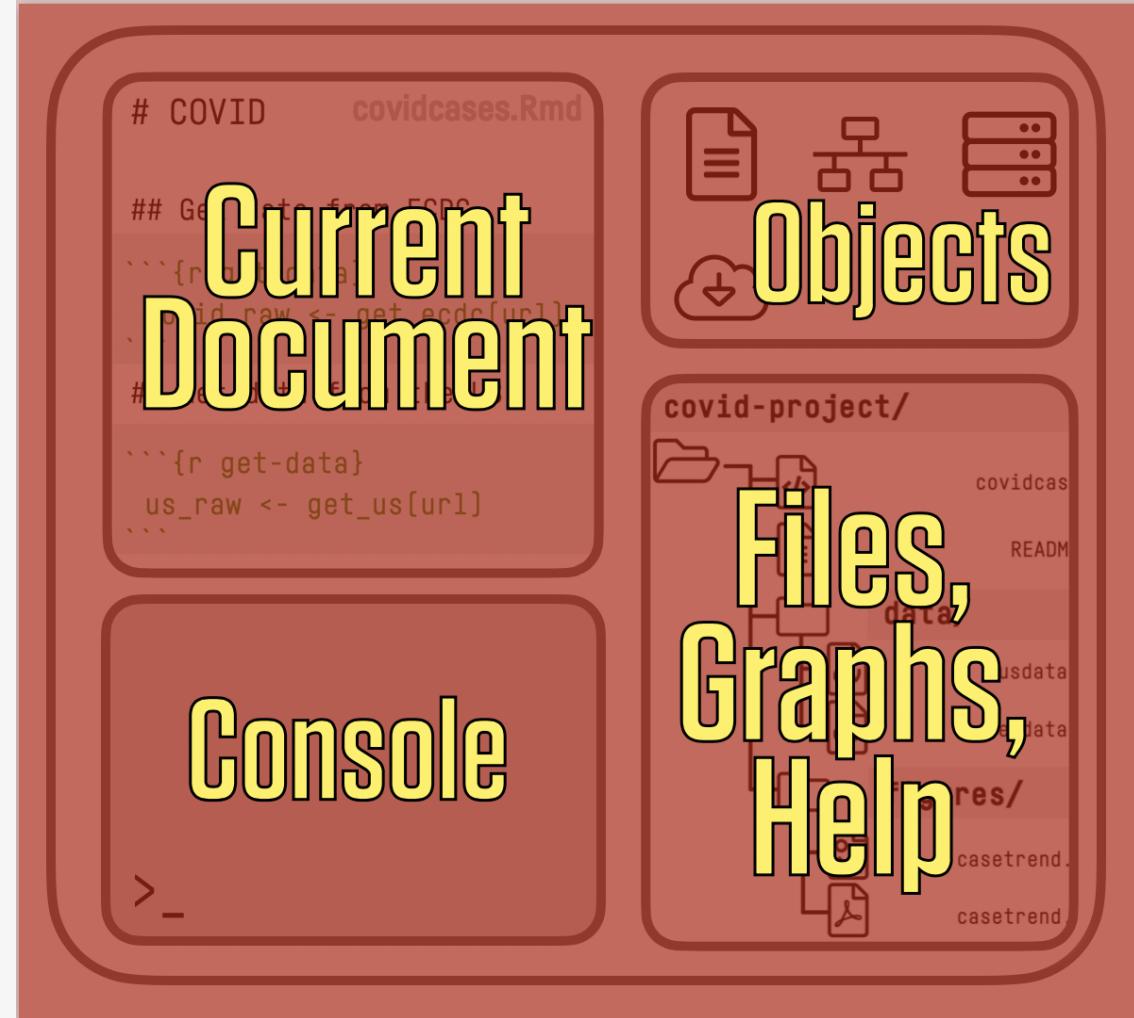
# R and RStudio



# R and RStudio



# R & RStudio



# R & RStudio

The screenshot shows the RStudio interface with a project titled "stathorizons\_0820". The left pane displays the RMarkdown file "01\_introduction.Rmd". The code includes metadata at the top and sections for "Data Visualization Notes" and "This Document is an RMarkdown File". A large red box highlights the preview pane, which shows a stylized title "Paper, Report, Analysis, Notes, etc, in RMarkdown" in red. The right pane shows the "Environment" tab with the message "Environment is empty". Below the tabs is a file browser showing the project structure. The bottom pane is the "Console" showing R session output.

```
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  
11 Visualization*. You can use it to take notes, write your code, and produce a good-looking,  
12 reproducible document that records the work you have done. At the very top of the file is  
13 a section of *metadata*, or information about what the file is and what it does. The  
14 metadata is delimited by three dashes at the start and another three at the end. You  
15 should change the title, author, and date to the values that suit you. Keep the 'output'  
16 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.  
17  
18 ## This Document is an RMarkdown File  
19  
20 Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For  
21 more details on using R Markdown see <http://rmarkdown.rstudio.com>.  
22  
23 When you click the **Knit** button a document will be generated that includes both content  
24 as well as the output of any embedded R code chunks within the document. A *code chunk* is  
25  
1:1 Data Visualization
```

Console

```
~/Documents/courses/stathorizons_0820/
```

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.

knitr hook "anchor" is now available  
Loading required package: testthat

Attaching package: 'testthat'

The following object is masked from 'package:devtools':

test\_file

>

Environment History Connections Git Tutorial

Global Environment

Environment is empty

Files Plots Packages Help Viewer

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, 0:50 AM

# R & RStudio

The screenshot shows the RStudio IDE interface. The main window is divided into several panes:

- Code Editor:** Displays the content of the file `01_introduction.Rmd`. The code includes metadata at the top and a section titled `## Data Visualization Notes`. The notes explain the structure of R Markdown files, mentioning metadata (three dashes at the start and end), reproducibility, and the generation of HTML, PDF, or MS Word documents.
- Environment:** A pane showing the global environment, which is currently empty.
- Files:** A sidebar showing the project structure under `stathorizons_0820`. It lists various R Markdown files (e.g., `01_introduction.Rmd`, `02_get_started.Rmd`, etc.) and other project files like `LICENSE.md` and `README.md`.
- Console:** A pane at the bottom left showing the R console output. It includes messages about the `knitr` hook, package loading, and attaching packages. A red box highlights this pane with the text "Console: Type or send code here, see results".

# R & RStudio

The screenshot shows the RStudio interface with a project titled "stathorizons\_0820". The left pane displays an R Markdown file named "01\_introduction.Rmd" containing introductory text about R Markdown. The right pane shows the "Environment" tab, which is currently empty. A red box highlights the "Files" tab in the bottom navigation bar of the right sidebar, which lists all files in the project directory. The "Project files, Plots, Help" section at the bottom right is also highlighted.

01\_introduction.Rmd

```
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  
11 Visualization*. You can use it to take notes, write your code, and produce a good-looking,  
12 reproducible document that records the work you have done. At the very top of the file is  
13 a section of *metadata*, or information about what the file is and what it does. The  
14 metadata is delimited by three dashes at the start and another three at the end. You  
15 should change the title, author, and date to the values that suit you. Keep the 'output'  
16 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.  
17  
18 ## This Document is an RMarkdown File  
19  
20 Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For  
21 more details on using R Markdown see <http://rmarkdown.rstudio.com>.  
22  
23 When you click the **Knit** button a document will be generated that includes both content  
24 as well as the output of any embedded R code chunks within the document. A *code chunk* is  
25
```

Console

```
~/Documents/courses/stathorizons_0820/   
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.  
  
knitr hook "anchor" is now available  
Loading required package: testthat  
  
Attaching package: 'testthat'  
  
The following object is masked from 'package:devtools':  
  
 test_file  
  
>
```

Environment

Global Environment

Environment is empty

Files

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		
data		
figures		
keynote		
LICENSE.md	5.7 KB	Jul 21, 2020, 11:16 AM
materials		
README.md		
slides		
stathorizons_0820.Rproj	205 B	Jul 22, 2020, 8:50 AM

Project files, Plots, Help

# R & RStudio

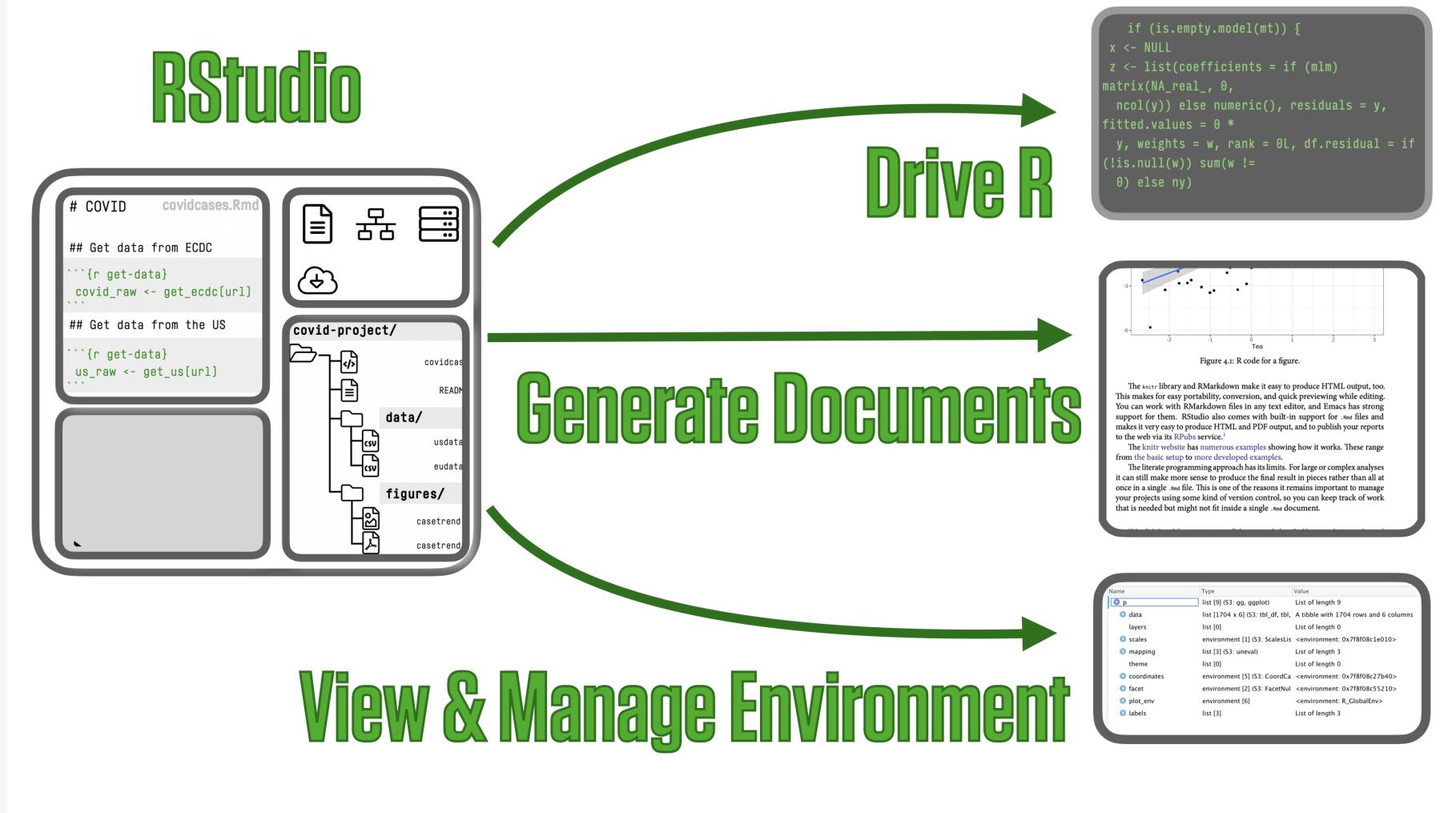
The screenshot shows the RStudio interface on a Mac OS X system. The main window displays an R Markdown file named '01\_introduction.Rmd'. The code in the file includes metadata (title, author, date, output) and introductory text about R Markdown. The 'Console' tab at the bottom shows R session logs, including the loading of the 'testthat' package and the creation of a 'test\_file'. To the right of the code editor is the 'Environment' pane, which is currently empty. A large red box highlights this pane. Below the environment pane is a file browser showing a directory structure for a project named 'stathorizons\_0820'. The browser lists files like .gitignore, various RMD files, and supplementary material. A red box also highlights the text 'Inspect objects you create' overlaid on the empty environment pane.

Environment is empty

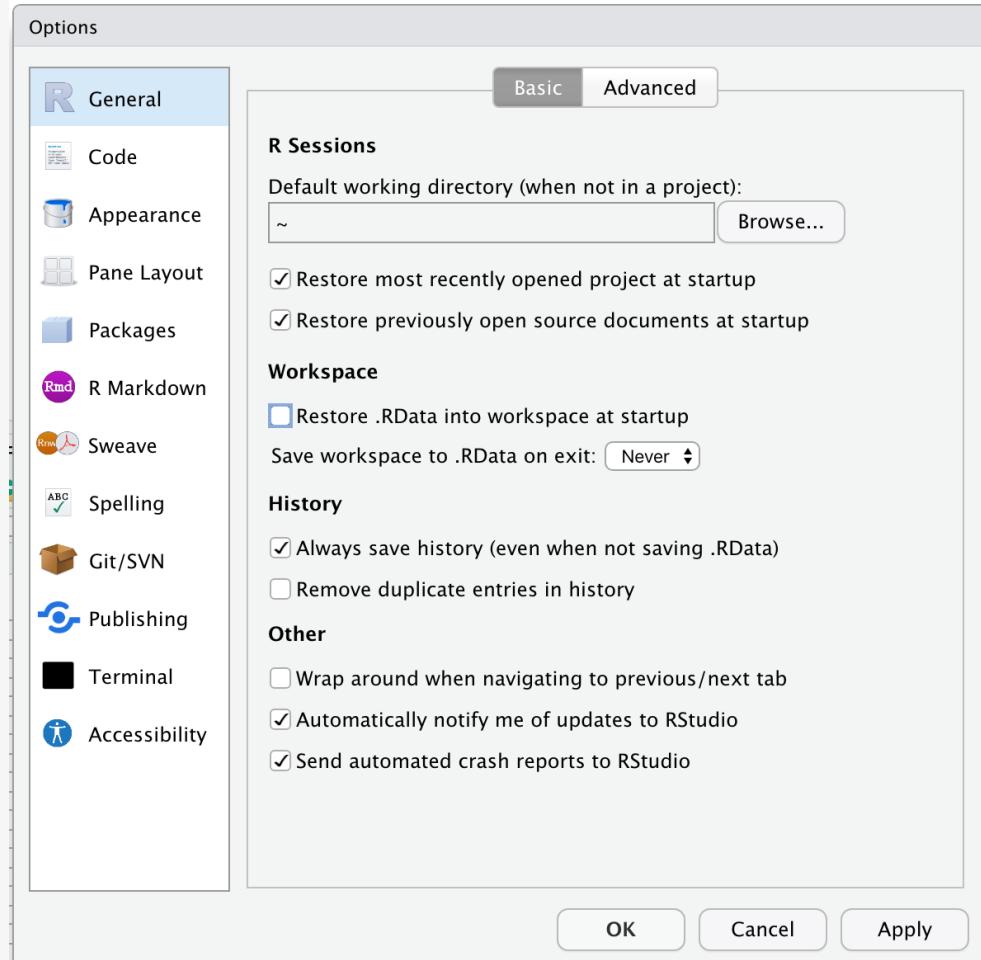
Inspect objects you create

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, 0:50 AM

# R & RStudio



# Your code is what's real in your project



# Consider not showing output inline

