

Sharing your work

A reproducible data science workflow
using open source tools



paradigm data group

Sharing your work

A reproducible data science workflow using open source tools

Martin Frigaard

This book is for sale at <http://leanpub.com/showingyourwork>

This version was published on 2019-08-08



This is a [Leanpub](#) book. Leanpub empowers authors and publishers with the Lean Publishing process. [Lean Publishing](#) is the act of publishing an in-progress ebook using lightweight tools and many iterations to get reader feedback, pivot until you have the right book and build traction once you do.

© 2019 Martin Frigaard

Contents

Intro: Why we wrote this	1
Our goal for anyone reading this book	2
Who this book is for	2
What this book covers	3
What this book doesn't cover	3
How this book is structured	3
What you'll walk away with	5
Part 1: ‘Good enough’ data skills	7
Why ‘good enough?’	7
How to share your work	7
‘Good enough’ communication	8
Get ‘good enough,’ then go for more if you need it	9
Part 2: “Have a workflow.”	10
Principle 1: Use open-source software	10
Principle 2: Write code	14
Principle 3: Document everything in plain text	17
Additional reasons for using R & RStudio	22
Part 3: Setting up your data science project	25
Example: FiveThirtyEight’s 2019 Presidential debate project	25
The Command line: Unix and Windows	33
Good enough command-line tools	35
Command line recap	48
Part 4: Keep track of changes with version control	49
Tracing our steps	49
Git	51
Setting up Git	52
Part 5: RStudio.Cloud	61
What have we done so far?	61
Recap on RStudio panes	71
Where should we write code?	72

CONTENTS

Rmarkdown	73
Documenting our projects	80
Import	81
Document changes with Git	84
Wrangle	88
Visualize	92
Knitting RMarkdown files	94
Part 6: Putting your project on Github	100
Push the changes to Github	100
Sharing our work (with Github pages)	104
Conclusion	110
Appendix	111
Intro	111
Chapter 1	111
Chapter 2	111
Chapter 3	112
Chapter 4	112
Chapter 5	112
Chapter 6	112

Intro: Why we wrote this

It seems like everywhere we look now, people are using data in beautiful and surprising ways to present their positions or shed light on new topics. We remember the first time we saw the impact data can have on storytelling. Hans Rosling, the Swedish physician/statistician, gave a [talk displaying the gapminder dataset¹](#). Rosling perfectly paired his enthusiasm with a brilliant display on the screen. As he spoke, over a dozen colorful circles slid across the projection. His [Ted profile description²](#) is perfect, “In Hans Rosling’s hands, data sings.”

Today, most primary sources of media use data as part of their evidence base. Check out the interactive data visualizations on the [UpShot³](#) at the New York Times, the visual journalism data projects at the [BBC⁴](#), or the daily graphs in the [Economist⁵](#).

The massive amounts of data available have spawned new forms of media. Nate Silver’s blog covering elections and politics has grown into multiple projects on [fivethirtyeight⁶](#). The [Pudding⁷](#) is an example of an online data journalism site that covers non-conventional sources of data. [Vox⁸](#) recently won an award for producing a [graph⁹](#) that communicates a topic that pundits could’ve debated endlessly.

Now that we’ve shown you all this cool stuff, we want to tell you why we wrote this book,

****“You’ve found something cool on the Internet, but you have no idea what it took to make it.” ***

There are a ton of really great resources on the Internet right now for learning data science (see [here¹⁰](#), [here¹¹](#), and [here¹²](#)). Many of these courses are fantastic—they can teach you programming languages, website design, database management, statistics, and machine learning. But we sometimes found the sheer volume of these courses can be overwhelming for audiences who are wanting to understand how these technologies fit together.

We decided to take a step back and write a book that describes a data science workflow, or *shows how these tools work together*. We’ll show you how R, RStudio, Git, & Github can be used to create elegant yet durable data analysis projects.

We chose to center this book around a particular use case, so there will be code files and tools we’ll use that are specific to this project. But we’ve chosen not to spend too much time on the content of

¹<https://www.youtube.com/watch?v=hVimVzgtD6w>

²https://www.ted.com/speakers/hans_rosling

³<https://www.nytimes.com/section/upshot>

⁴<https://www.bbc.com/news/world-32209370>

⁵<https://www.economist.com/node/21011894>

⁶<https://projects.fivethirtyeight.com/>

⁷<https://pudding.cool/>

⁸<https://www.vox.com/>

⁹<https://www.vox.com/policy-and-politics/2018/9/28/17914308/kavanaugh-ford-question-dodge-hearing-chart>

¹⁰<https://www.coursera.org/learn/r-programming>

¹¹<https://www.edx.org/learn/r-programming>

¹²<https://www.udacity.com/course/data-analysis-with-r--ud651>

these files (we've documented them you want to look into the details). Instead, We're going to focus more on the "high level" ideas because these are topics you can take with you to your next project. We also encourage you to consult the articles and resource we've recommended throughout each chapter for more materials on each topic.

Our goal for anyone reading this book

We want to show you how to 1) take something neat you found on the Internet, 2) figure out what went into making it, and 3) see if you can reproduce the result.

We plan to include enough information to get you up and running and at the same time, not overwhelm you. If you've already Googled "Getting started in data science," you know there are a *ton* of resources. Figuring out where to start can feel like trying to get a drink of water from a fire hose.

Along the way, we will also cover some practical principles of programming, command-line tools, project file organization, and a few computer science topics.

Who this book is for

We've tried to keep the materials accessible to a broad audience, but we understand there are few useful data analysis texts written for everyone. Data tends to be very specific to the field they come from, and it's hard to find data that gets everyone excited. To try and help address this issue, we use data from multiple sources (Google trends, Twitter, Wikipedia, and Googlesheets).

We focus on the workflow and tools.

The next chapters outline a 'one-stop-shop' toolset that you can learn quickly and readily re-use (because we know your time is limited).

Technical assumptions

The reader we had in mind while writing this book was someone who,

1. Uses a computer every day at work
2. Understands how to navigate a web browser (Chrome, Safari, Firefox, etc.)
3. Has worked with a word processor (like Microsoft Word, Google Docs, or Apple Papers)

If you're an accountant, scientist, analyst, journalist, grad student, product manager, or decision-maker, this book is for you.

What this book covers

We will be covering R, RStudio, Git, and Github. We use these tools daily now, but we began our careers in other statistical programs (SPSS, Stata, SAS). We abandoned those tools (we know your pain) because of the sheer number of tasks we can accomplish, and that's what makes us recommend this toolkit to you. We've also reached out to our colleagues and included their lessons and insights.

What this book doesn't cover

We also understand there are alternative approaches to accomplishing the same goal, and we will try to mention these examples wherever possible. Jupyter Lab and Jupyter Notebooks, for example, offer reproducible scientific programming environments that can accomplish many of the same objectives we'll tackle in this book. However, we still think there are reasons you should use RStudio + Github instead, and we will outline these in the following chapters.

How this book is structured

We structured this book somewhat like an Army Field Manual, which means each topic was chosen using the following criteria:

- (a) *Relative importance*. Which activities contribute most to successful training?

This book contains **brief descriptions** of the tools we recommend, with **diagrams and figures** outlining how they work, and **examples** for using them.

- (b) *Need*. Which training activities will benefit the most from guidance? Which activities have received little attention in the past or which have previously required improvement?

We'll expand on a few tools we felt are harder to grasp (Git and version control). We will also go over topics typical college courses overlook or neglect (file naming, project organization).

- (c) *Time*. How much time is available? Which activities can be effectively taught in that time?

Time is the real enemy of any data project. All computational work comes down to keystrokes and neurons. This book is trying to narrow the gap between 1) seeing a data product (neurons) and 2) translating what you see into commands on a computer that can be used to recreate that product.

- (d) *Personnel.* What are the known or suspected levels of expertise among individuals receiving training?

We assume everyone reading this text has very little exposure to the tools we'll be covering (R, RStudio, Git, or Github). We do expect you are comfortable using a computer.

The secret to the Army's training abilities is the Field Manual (FM). Army FMs are amazing—they cover almost any topic you can imagine and are well illustrated. For example, watch this video of the drill and ceremony movement called the “counter column”¹³.

Now, look at the same thing in a figure.

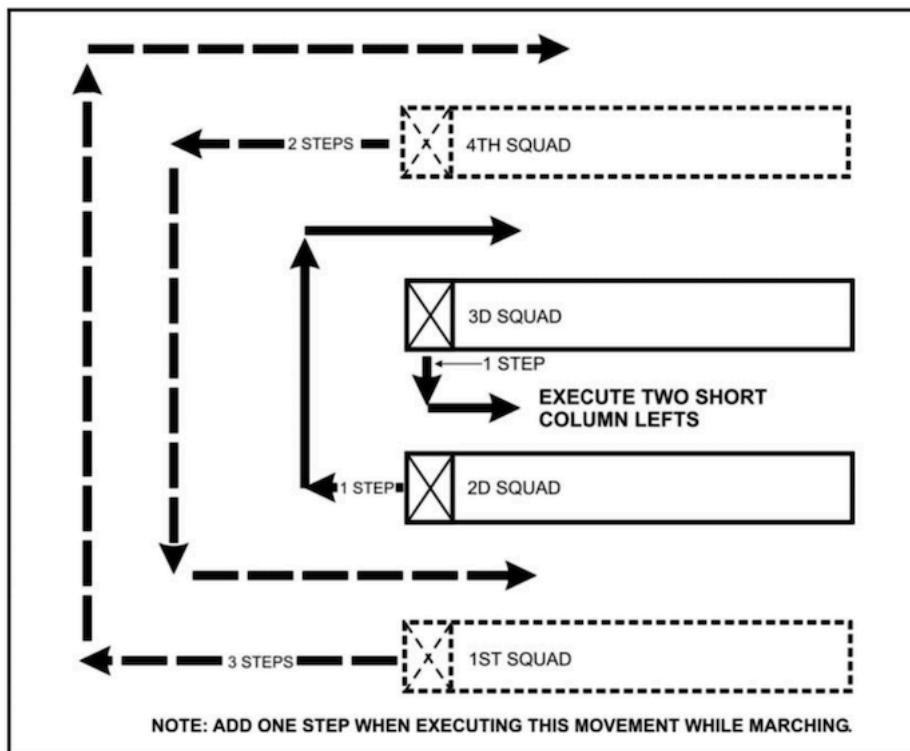


Figure 7-2. Counter-Column March at the Halt.

As you can see, executing a counter column is complicated. But the Army has taught hundreds of thousands of soldiers on this movement for decades. How? They give soldiers a field manual (FM 22-5) to read and dedicated time to practice.

The strength of the FMs is how they present information: they gave the material in everyday language (usually between a 6th-8th-grade reading level) with an emphasis on diagrams, pictures, and simple drawings.

¹³<https://www.youtube.com/watch?v=EgeZl9UOJ0I>

We've found so much data science and statistical information on the internet has a ton of acronyms, jargon, and equations. We've actively avoided using technical verbiage, and focused on using figures and graphics.

"...there are lots of other books that explain what things are called. This book explains what they do."

The quote above comes from Randall Munroe, author of the [xkcd¹⁴](#) comic. In his book "[Thing Explainer](#)"¹⁵, Munroe uses pictures and plain language to describe multiple complex systems (rocket ships, the periodic table, laptops, etc.).

The subtitle of "*Thing Explainer*" is *Complicated stuff in simple words*, which is what we're trying to replicate here. Wherever possible, we've dropped unnecessary technical jargon and spelled out any acronyms.

What you'll walk away with

You will have a working project (cool visualizations, lots of code, data) a ton of resources, and a book for reproducing this process again.

Language and style guide

We use the plural 'we' throughout the book based on the [excellent advice¹⁶](#) from Donald Knuth, Tracy Larabee, and Paul Roberts, "*think of a dialog between author and reader..*"

As with most written works, the topics in this book are the result of many conversations, emails, comment threads, and communications that could not have happened in isolation. We want to thank everyone who's contributed to these ideas over the years.

The text uses the following style guide:

this is code.

1 # this a code chunk

some quoted text

click on hyperlinks¹⁷

plain text for our thoughts

¹⁴<https://xkcd.com/>

¹⁵<https://xkcd.com/thing-explainer/>

¹⁶<http://www.econ.uiuc.edu/~econ508/Papers/mathwriting.pdf>

¹⁷

Learn more

- Practical Data Science for Stats¹⁸ is a resource you should bookmark in your browser. The articles in this collection will come up again in future sections, but we found we use these resources so much it's nice to have them somewhere handy.
- The R for data science community¹⁹ and R for Data Science²⁰ book are excellent resources to help you started.
- Collaboration and reproducibility - there's a direct connection between collaboration and reproducibility. The better your collaborators can reproduce your work, the better they'll understand your results.
- We recommend RStudio and Github for anyone looking to get started with data science, visualization, reproducible reporting, dashboards development, or website/blog creation. By suggesting these particular tools, we're not saying there aren't other ways or workflows capable of accomplishing the same activities. These are the tools we've found success with, so they're what we recommend.

¹⁸<https://peerj.com/collections/50-practicaldatascistats/>

¹⁹<https://www.rfordatasci.com/>

²⁰<https://r4ds.had.co.nz/>

Part 1: 'Good enough' data skills

In 2016, [The Carpentries²¹](#)—a non-profit organization that teaches coding and data science skills to researchers—published an article titled, “[Good Enough Practices for Scientific Computing](#)²²”.

The article above has tons of great information, but it was aimed at “*researchers who are working alone or with a handful of collaborators on projects lasting a few days to several months.*” We thought the article’s information was too useful *not* to share with as many people as possible. This book essentially attempts to extend the advice in paper to analysts, journalists, grad students, and non-technical audiences.

Why ‘good enough’?

A constant theme that runs through this book is being able to do “good enough” work. We’re going to introduce you to a lot of technology in the next few pages, and we want you to set reasonable expectations for using these tools.

There are many courses, tutorials, and resources that promise you ‘expert training’ in all of the tools we’re covering. You don’t have time for that—you need enough knowledge to see cool stuff and recreate it. In the beginning, you need to focus on doing good enough work, and then sharing your work. The feedback and input you get from developing things in the open are what builds your skillset. Expertise is something you can attain with experience (if that’s your goal), but in many cases, being good enough will get lots of work done.

How to share your work

“*Your work should speak for itself...*” - author unknown

There is a sea of information on the Internet, and that means everyone is competing for everyone else’s attention. You want to share whatever you’re doing (writing code, building graphs, creating apps, etc.) so it’s discoverable. This way, if a future collaborator, prospective employer, or up-and-coming analyst is looking around for cool stuff on the Internet, they’ll see what you’ve been doing.

²¹<https://carpentries.org/>

²²<https://swcarpentry.github.io/good-enough-practices-in-scientific-computing/>

Make cool sh!t, then share it

We’ll introduce you to a few data scientists and journalists who are excellent communicators of their work. These examples use the Internet as a tool to engage with broader audiences, create better tools for doing science, document some of their daily struggles/successes.

Our first example, [Lucy D’Agostino McGowan²³](#) is a post-doc at Johns Hopkins Bloomberg School of Health. She maintains a [blog²⁴](#), publishes [ebooks²⁵](#), has [online courses²⁶](#), and also attempts to create a [real BB-8²⁷](#).

Ricardo Bion is the [Data Science manager at Airbnb²⁸](#). He [publishes papers²⁹](#) on using R in their business setting, gives [webinars³⁰](#) on how to use modeling to make business decisions, and [writes articles³¹](#) on workflow practices that contribute to success in their data science teams.

Or take [Amber Thomas³²](#), the Senior Journalist-Engineer at the [Pudding³³](#). She writes [tutorials³⁴](#), [maintains a blog³⁵](#), and spends her time “hanging upside down in aerial silks.”

All of these people have done two things very well:

1. **Created good work:** All created projects across multiple mediums, websites, and platforms
2. **Shared with as many people as possible:** these examples made their projects discoverable and collaborated with the data science community (by putting their work online for people to find)

Of course, they also had to know their subject areas, and have something worth sharing. But they didn’t wait until their work was perfect, or until they were done with their careers to share and get feedback. They started engaging with people while they were working to show how their work gets done.

‘Good enough’ communication

It’s important to remember that whenever you’re trying to communicate something, you’re convincing your audience that they should be paying attention to you instead of everything else. Today, there’s a lot more ‘everything else.’

²³<https://www.lucymcgowan.com/>

²⁴<https://livefreeordichotomize.com/>

²⁵<https://leapub.com/ggplot2in2>

²⁶<https://leapub.com/u/lucymcgowan>

²⁷<https://magazine.amstat.org/blog/2017/11/01/lucy-dagostino-mcgowan-and-ryan-jarrett/>

²⁸<https://t.co/EaT2pX2wWm?amp=1>

²⁹<https://peerj.com/preprints/3182/>

³⁰<https://www.rstudio.com/resources/videos/airbnb/>

³¹<https://medium.com/airbnb-engineering/using-r-packages-and-edition-to-scale-data-science-at-airbnb-906faa58e12d>

³²<https://twitter.com/puddingviz>

³³<https://pudding.cool/>

³⁴<https://pudding.cool/process/flexbox-layout/>

³⁵<https://amber.rbind.io/>

Approaching communication this way puts you in the mind of your audience and keeps you asking, “*why would they want to know this?*” Try to keep the value of what you’re saying visible to your audience.

Avoid technical jargon & acronyms

“You must learn to talk clearly. The jargon of scientific terminology which rolls off your tongues is mental garbage.” - Martin H. Fischer

The most substantial barrier to understanding new disciplines or technologies is getting a handle on their jargon. Because this book sits at the intersection of computer science, statistics, and web technologies, all the new vocabulary can often seem like learning a foreign language.

Wherever possible, we’ll do our best to clear up or define any terms related to computer science, data management system, web technology, or statistics. To maximize the power of the tools in this text, it will help to know a little about their history, so we’ll also cover some background.

Communication takes practice, but it's worth it!

No one is born with an ability to write well—it takes a lot of practice and feedback. The more you communicate with different audiences, the better you’ll get at finding an ability to convey why what you have to say is essential.

When we were kids in math class, most teachers required us to show how we got the answers (“show your work”). Teachers told us ‘show our work’ so they could follow our thought processes through a problem, and see where our thinking was incomplete or mistaken.

If you’re regularly sharing your work, people can follow your line of thinking as you progress throughout your project. More importantly, people who find your work will give you feedback and improve your ideas.

Get 'good enough,' then go for more if you need it

What is ‘good enough’? We think being ‘good enough’ means reading about a tool or technology and being capable of distinguishing it from magic. Good enough also means seeing a chart or graph on the Internet and knowing how to evaluate its contents. Or imagining something that matters in your business or personal life, and then devising a way to count it.

Feedback

We sincerely hope you’ll find this information useful and give us feedback at mfrigaard@paradigmdata.io or pspangler@paradigmdata.io.

Part 2: “Have a workflow.”

A **workbench** is a place to keep and organize tools, and a **workflow** is how we combine these tools to get things done. This chapter will cover the workbench we use and the three guiding principles of the workflow we recommend.

1. Use free open source software
2. Write code
3. Document everything in plain text

Principle 1: Use open-source software

All of the tools in this book are available open-source and available free of charge. Just as a point of reference, the cost of a subscription to SPSS at the time of this writing is \$99.00 per user per month. Stata is \$595 per year or \$1,595 for a perpetual license. There are educational discounts available, but this cost is not offset by much when we take into account the rising price of tuition.

A more important reason we recommend open source tools are the communities that we got access to when we started using them. By entering the universe of open source software, we got to take advantage of seeing problems solved in the open. We also found people like us, grappling with the same issues, and it’s hard to overstate the benefit of this shared camaraderie.

The final reason is philosophical: we all benefit from using open source tools and sharing improvements on them together. The ‘four freedoms’ of open source software³⁶ captures this sentiment below.

Freedom 0: The freedom to run the program as you wish, for any purpose.

Freedom 1: The freedom to study how the program works, and change it, so it does your computing as you wish.

Freedom 2: The freedom to redistribute copies so you can help your neighbor.

Freedom 3: The freedom to distribute copies of your modified versions to others. By doing this, you can give the whole community a chance to benefit from your changes.

We’ve displayed some examples of open source tools for data management, statistics, and communication in the image below:

³⁶<https://www.gnu.org/philosophy/free-sw.html>

Open source tools



Follow the links below to learn more.

- Git³⁷
- Github³⁸
- Linux³⁹
- MySQL⁴⁰
- Netlify⁴¹
- Python⁴²
- R⁴³
- RStudio⁴⁴

³⁷<https://git-scm.com/>

³⁸<https://github.com/>

³⁹<https://www.linux.org/>

⁴⁰<https://www.mysql.com/>

⁴¹<https://www.netlify.com/>

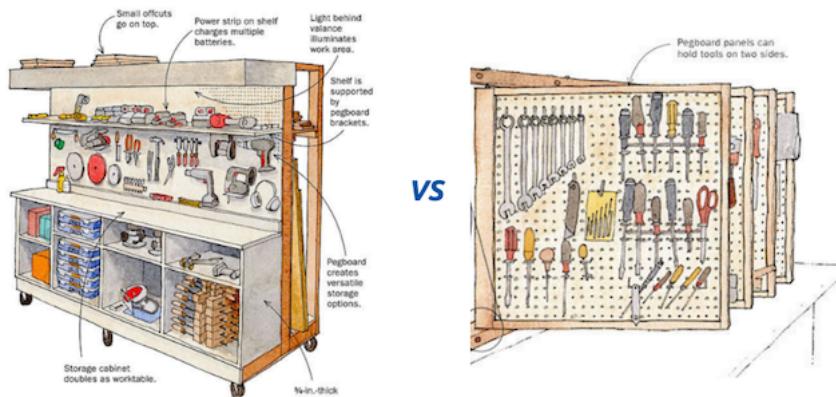
⁴²<https://www.python.org/>

⁴³<https://www.r-project.org/>

⁴⁴<https://www.rstudio.com/>

The integrated development environment (aka the workbench)

An [integrated development environment](#)⁴⁵ (IDE) is an application typically used by programmers to build and test software. We find it helpful to think of a woodworkers workbench as an analogy. The workbench on the left is an example of a simple rolling cart design (for people with minimal garage space).



As we can see, this workbench is efficiently designed to keep essential tools for the job within arms reach, and it uses storage space efficiently. IDE design follows these same principles. However, we also know some models are better than others. For example, consider the design of the workbench on the right. By making the panel positions adjustable, the workbench allows for easier access to more tools. Depending on the job, woodworkers can customize the panel arrangements with "*the simple pin that allows the rack's various faces to swing left/right for access to either side..*"

These examples illustrate how differences in the design of a workbench can have a meaningful impact on levels of productivity. Well-designed workbenches give us access to more tools without making these tools more difficult to find.

Our workbench

We recommend choosing a workbench that minimizes the number of additional applications you'll need to have open to get work done. We've found we can use [RStudio](#)⁴⁶ for ~90% of our daily work (*they're not paying us to say this*). RStudio gives us access to all the tools we need in the same place.

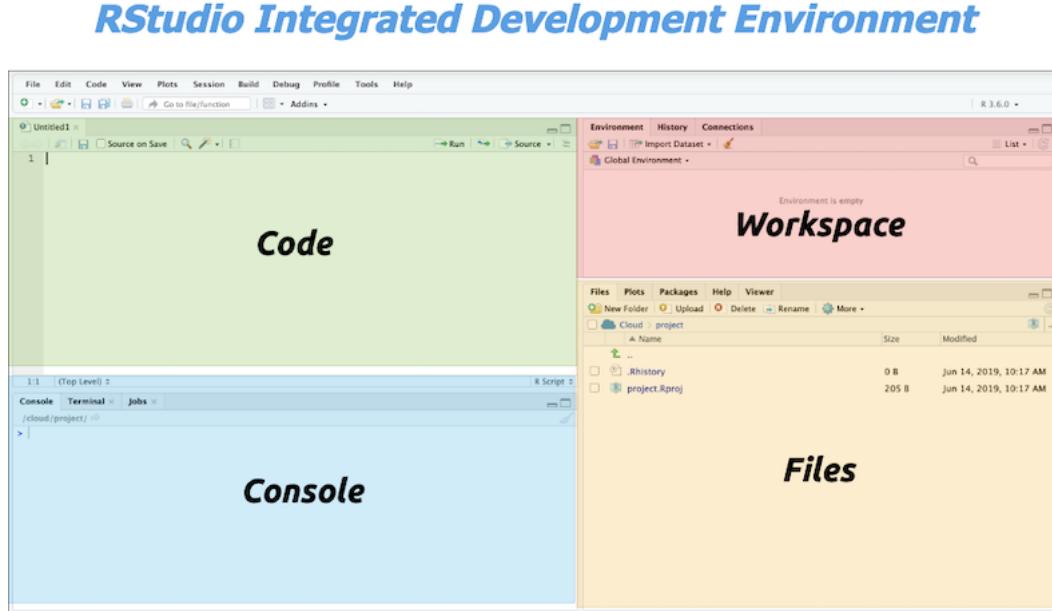
RStudio has four primary panes, each serving a specific function.

- the **Code or Source** pane is where we can document both human-readable and computer-readable text
- the **Workspace** holds the data, functions, and other data analysis objects
- the **Console** displays the results from our written code (and allows us to enter commands directly)

⁴⁵https://en.wikipedia.org/wiki/Integrated_development_environment

⁴⁶<https://www.rstudio.com/products/RStudio/>

- the **Files** gives us access to the happenings outside the RStudio environment (imported raw data, exported results, etc.)



Just like any workbench, we need to fill RStudio with tools we need for the job. RStudio plays well with many open-source software tools, but or now we're going to focus on R and Git.

What is R?

R is a free statistical modeling software⁴⁷ application and language. If you are using the desktop application, follow the links below to install R.

Installing R & RStudio

1. First, you'll need to download and install R from [CRAN⁴⁸](https://www.r-project.org/).
2. Second, download and install [RStudio⁴⁹](https://www.rstudio.com/products/rstudio/download/), the integrated development environment (IDE) for R

Use RStudio in the browser

An alternative to downloading and installing R and RStudio is using [RStudio.Cloud⁵⁰](https://www.rstudio.cloud/) which operates entirely in your browser. You'll need to sign up for RStudio.cloud for free using your Google account or email address, but we recommend using a Github account. You can create a Github account [here⁵¹](https://github.com/join).

⁴⁷<https://www.r-project.org/>

⁴⁸<https://cran.r-project.org/>

⁴⁹<https://www.rstudio.com/products/rstudio/download/>

⁵⁰<https://rstudio.cloud/>

⁵¹<https://github.com/join>

What is Git?

Git⁵² is a version control system (VCS). VCSs are used to track changes to projects with code. You can read more about Git in their online [text here⁵³](#).

What is Github?

Github⁵⁴ is the web-based hosting service for Git. You should set up a free account with Github [here⁵⁵](#).

What do Git/Github do?

We will cover more on Git/Github in later sections, but for now, know these tools will allow us to keep track of changes to our projects over time.

Note: You should explore different IDE's on your own – you'll see there are many options, both paid and unpaid. We're confident you'll see RStudio is well suited to handle more than most of the things you'll want to accomplish.

Open-source software bonus: As mentioned previously, you'll also find a massive network of support on [Stackoverflow⁵⁶](#), [RStudio Community⁵⁷](#), and [Google Groups⁵⁸](#).

Principle 2: Write code

Usually, people interact with their computers using point-and-click [graphical user interfaces⁵⁹](#) or GUIs (pronounced ‘gooey’). GUIs are quick and easy to learn because their design environment usually mimics an actual physical space (i.e., desktops, folders, or documents). GUIs are a mostly positive development because designing software with a more [user-centered design⁶⁰](#) is one of the main reasons technology adoption has been on the rise for the past 20+ years.

User-centered-designed software includes most of the point-and-click operating systems and applications. These programs give the user the ability to click through a predetermined list of options and procedures using their mouse or track-pad. Two examples of statistical programs designed around the GUI are below.

⁵²<https://git-scm.com/>

⁵³<https://git-scm.com/book/en/v2>

⁵⁴<https://github.com/>

⁵⁵<https://github.com/join?source=header>

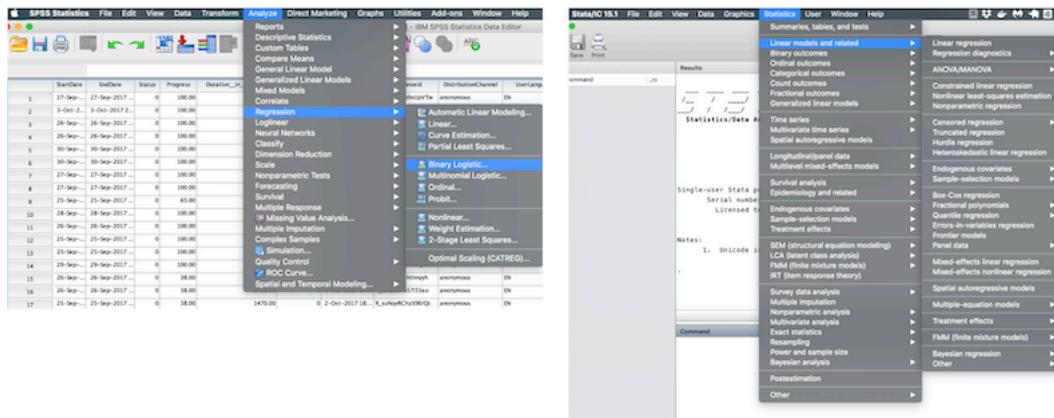
⁵⁶<https://stackoverflow.com/questions/tagged/r>

⁵⁷<https://community.rstudio.com/>

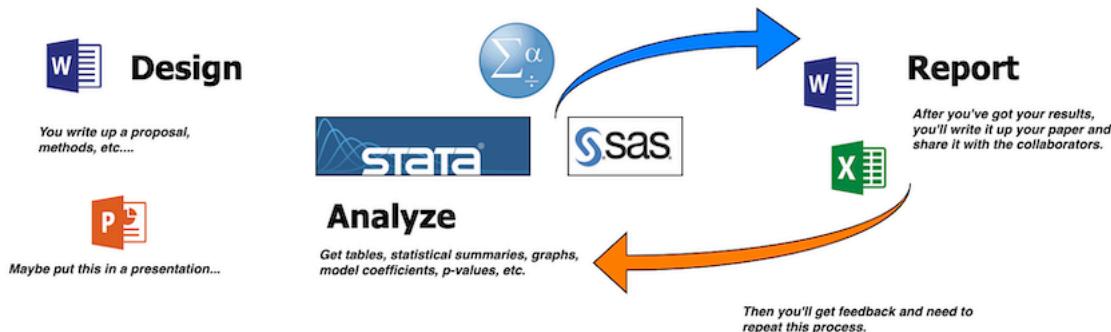
⁵⁸<https://groups.google.com/forum/#!forum/r-help-archive>

⁵⁹https://en.wikipedia.org/wiki/Graphical_user_interface

⁶⁰https://en.wikipedia.org/wiki/User-centered_design



As we can see, these programs provide an extensive list of options to choose from. However, we think there are times when we should resist the temptation to abstract away some of life's complexity, and data science is one of them. Using applications like these encourage a copy + paste workflow like the one below.



If you follow the arrows in the image, you can see we'll be required to go *back* through the same elaborate workflow whenever we receive feedback or input to our project. Each time will take just as long as the first.

We recommend an alternative to the copy + paste workflow based on the activities of a modern scientist from Jeff Leek we outlined in Chapter 1:

1. Develop code in the open
2. Publish data and code open source
3. Post preprints of your work
4. Submit and review for traditional journals
5. Blog or use social media to critique published work

Two things should stand out from the list above: First, modern science is mostly writing. Second, some of that writing is code (i.e., programming). It's, for this reason, we recommend adopting a workflow based on "*writing code*" wherever possible. We are aware that 'writing code' means being able to type, which might be daunting for people who struggle on a keyboard. We recommend

practicing this skill (there are plenty of great apps out there to help!) because typing is an unavoidable necessity for using a computer.

Software as a tool, not a solution

We think of data science software as the tools that help us gain a deeper understanding of the world. We don't think of software as an alternative to thinking or expect a software tool to do our thinking for us. We also prefer tools that we can reuse on future projects.

As we stated in chapter 1, the scientific method is a process. Software is a tool to help move that process along faster. Tools that improve our understanding shorten the distance between questions and answers, but doesn't leave out any crucial details.

For this reason, we don't recommend relying too heavily on point-and-click proprietary software applications (SPSS, Stata, SAS, etc.). These GUI's make it hard to keep track of what we've clicked on, the order we've clicked on them in, and can oversimplify or obfuscate what's going on.

What do we mean by 'a workflow'?

A workflow is a set of steps that can be used repeatedly to answer any question we might encounter in our work.

The quote below is from an interview with Andrew Gelman, a statistician from Cornell, who is an author on the excellent blog [Statistical Modeling, Causal Inference, and Social Science](#)⁶¹.

Question: "I'm wondering how you, as an educator and statistician, would like to see statistical and data literacy change in general for a general population?"

Answer: "... I've come to realize that a lot of people don't even know what they did. People don't have a workflow, they have a bunch of numbers, and they start screwing around with the numbers and putting calculations in different places on their spreadsheet, and then at the end, they pull a number out and write it down and type it into their report." - [Andrew Gelman](#)⁶²

As the quote above illustrates, ***how you got an answer is just as relevant as the answer you got.*** The tools we provide in this text give us a start-to-finish chain of documentation from question to solutions.

Data science jobs need a particular set of tools, and a workbench to organize these tools. To manage our data science projects, we'll need a workflow that gives us the ability to 1) document our intentions and, 2) write code that can translate our plans into something a computer can execute. These two points bring us to our next topic: plain text. As we'll discover, plain text files are a great way to accomplish these tasks.

⁶¹<https://statmodeling.stat.columbia.edu/>

⁶²<https://soundcloud.com/dataframed/election-forecasting-polling>

Principle 3: Document everything in plain text

In [The Pragmatic Programmer](#)⁶³, authors Hunt and Thomas advise ‘*Keep[ing] Knowledge in Plain Text*’. This sentiment has been repeated [here](#)⁶⁴, [here](#)⁶⁵, and [here](#)⁶⁶.

We recommend keeping files, notes, and any pertinent documentation about our project in plain text files. The reasons for this will become more apparent as we move through the example, but I wanted to outline a few here:

- plain text lasts forever (files written 40 years ago are still readable today)
- plain text can be *converted* to any other kind of document
- plain text is searchable (ctrl+F or cmd+F allows us to find keywords or phrases)

We'll also cover why you might want switch over to a plain text editor if you're currently using Google Docs, Apple Papers, or Microsoft Word.

Wait—why would I change what I'm doing if it works?

We get it—change is difficult, and if you've got a working ecosystem of software that keeps you productive, don't abandon it. However, you should be aware of these technologies and recognize that people using them will be adapting *their* workflows to collaborate with you.

We covered the problems with a copy+paste workflow previously, but there are additional reasons to avoid this toolset:

1. It's not reproducible
2. It's not logical or necessarily honest to separate computation from the analysis or presentation
3. It's error-prone

If we've convinced you RStudio is a flexible and adaptable tool, but you're still unclear on what makes a file ‘plain text,’ we'll cover that next. But first, we need to talk about what *isn't* a plain text file.

What *isn't* plain text

Non-plain text files are usually called binary (i.e., files with binary-level compatibility) need special software to run on our computer. The language below is a handy way to think about these files:

“Binary files are *computer-readable but not human-readable*⁶⁷”

⁶³<https://www.amazon.com/Pragmatic-Programmer-Journeyman-Master/dp/020161622X>

⁶⁴<https://simplystatistics.org/2017/06/13/the-future-of-education-is-plain-text/>

⁶⁵<https://richardlent.github.io/post/the-plain-text-workflow/>

⁶⁶<http://plain-text.co/index.html#introduction>

⁶⁷https://www.webopedia.com/TERM/B/binary_file.html

What is plain text

So if binary files aren't plain text, what is a plain text file? The language from the [Wikipedia](#)⁶⁸ description is helpful here:

When opened in a text editor, plain text files display computer and human-readable content.

The last bit is the most crucial distinction—**human-readable vs. computer-readable**. In this manual, we'll point out which files are binary and which are plain text. Generally speaking, plain text files can be opened using a text editor or viewed with a command-line tool. Examples of text editors include [Atom](#)⁶⁹, [Sublime Text](#)⁷⁰, and [Notepad++](#)⁷¹.

Markdown & Rmarkdown

Markdown⁷² files (.md) are common type of plain text files. Markdown is a ‘lightweight markup language,’ which means it’s easy for humans to read, and computers can convert it to HyperText Markup Language (HTML). Markdown allows for some formatting options to aid with communication (see below)

```

1  <!-- comments -->
2
3  normal text
4
5  *italic*
6
7  **bold**
8
9  > quote
10
11 `code`
12
13 # h1
14 ## h2
15 ### h3
16 ##### h4
17 ###### h5
18 ####### h6

```

⁶⁸https://en.wikipedia.org/wiki/Text_file

⁶⁹<https://atom.io/>

⁷⁰<https://www.sublimetext.com/>

⁷¹<https://notepad-plus-plus.org/>

⁷²<https://en.wikipedia.org/wiki/Markdown>

To learn more, see [Markdown Syntax Documentation⁷³](#) on John Gruber’s site).

We recommend learning markdown before any other programming language because it’s the lingua franca for asking questions. [Stackoverflow⁷⁴](#), [RStudio community⁷⁵](#), [Reddit⁷⁶](#), [Github⁷⁷](#), and many other sites use markdown to post questions and answers. We

recommend experimenting with [StackEdit⁷⁸](#), a browser-based markdown editor that gives user’s the ability to write in markdown and see the syntax rendered as HTML.

RStudio has an extension of markdown, [RMarkdown⁷⁹](#). Using RMarkdown in RStudio allows for a genuinely reproducible workflow: we’re able to write our thoughts, code, display results, and then share everything in multiple outputs.

⁷³<https://daringfireball.net/projects/markdown/syntax>

⁷⁴<https://stackoverflow.com/>

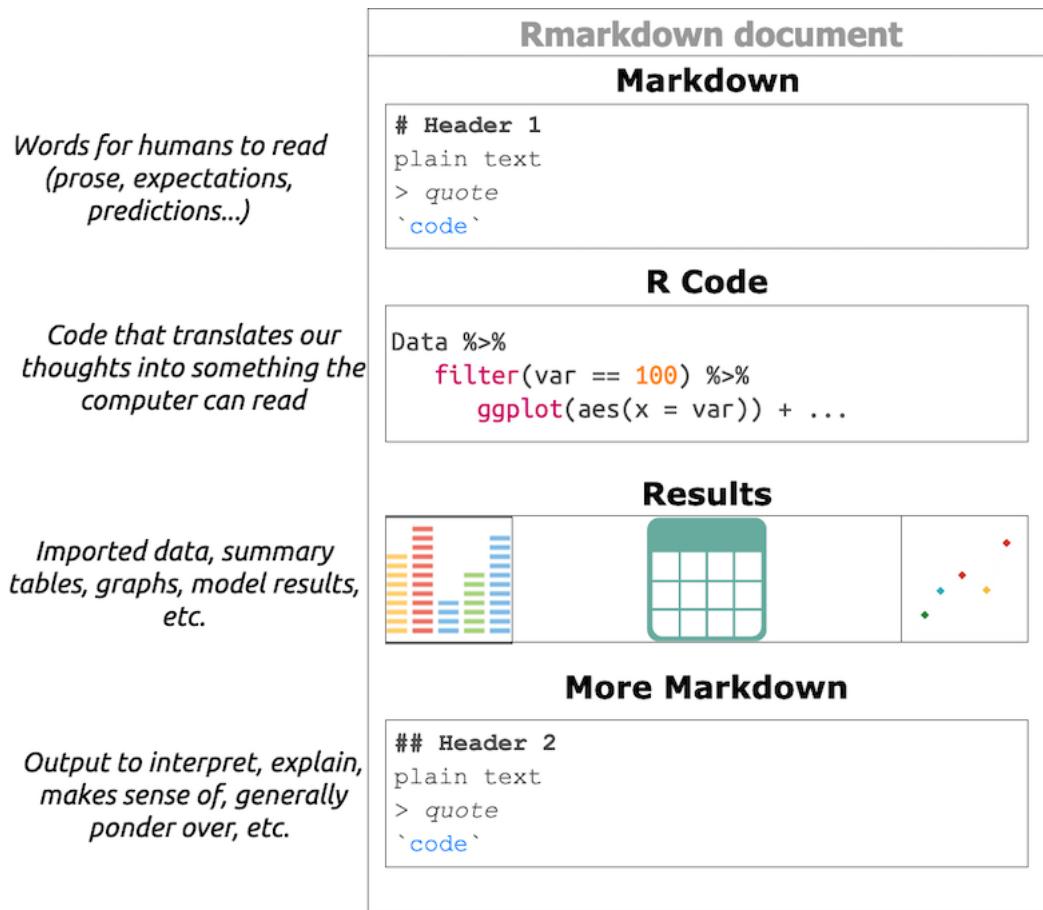
⁷⁵<https://community.rstudio.com/>

⁷⁶<https://www.reddit.com/>

⁷⁷<http://github.com/>

⁷⁸<https://stackedit.io/app#>

⁷⁹<https://rmarkdown.rstudio.com/>



document in RStudio.

⁸⁰<https://bookdown.org/yihui/rmarkdown/>

I recommend reading up on R and RMarkdown because of how many different outputs this combination can be used to produce (.pdf, .docx, and .html). Consult the [R Markdown: The Definitive Guide⁸⁰](#) for more information. The image below is an output from an .Rmd

Rmarkdown Document

R Markdown

This is an R Markdown document. 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>.

Load the tidyverse.

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. You can embed an R code chunk like this:

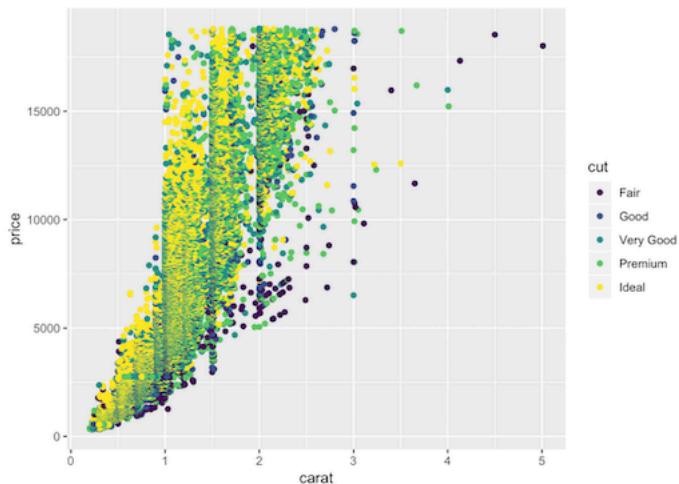
```
ggplot2::diamonds %>% glimpse(78)
```

```
## Observations: 53,940
## Variables: 10
## $ carat    <dbl> 0.23, 0.21, 0.23, 0.29, 0.31, 0.24, 0.24, 0.26, 0.22, 0.23, ...
## $ cut       <ord> Ideal, Premium, Good, Premium, Good, Very Good, Very Good, ...
## $ color     <ord> E, E, E, I, J, J, I, H, E, H, J, J, F, J, E, E, I, J, J, ...
## $ clarity   <ord> SI2, SI1, VS1, VS2, SI2, VVS2, VVS1, SI1, VS2, VS1, SI1, VS...
## $ depth     <dbl> 61.5, 59.8, 56.9, 62.4, 63.3, 62.8, 62.3, 61.9, 65.1, 59.4, ...
## $ table     <dbl> 55, 61, 65, 58, 58, 57, 57, 55, 61, 61, 55, 56, 61, 54, 62, ...
## $ price     <int> 326, 326, 327, 334, 335, 336, 336, 337, 337, 338, 339, 340, ...
## $ x         <dbl> 3.95, 3.89, 4.05, 4.20, 4.34, 3.94, 3.95, 4.07, 3.87, 4.00, ...
## $ y         <dbl> 3.98, 3.84, 4.07, 4.23, 4.35, 3.96, 3.98, 4.11, 3.78, 4.05, ...
## $ z         <dbl> 2.43, 2.31, 2.31, 2.63, 2.75, 2.48, 2.47, 2.53, 2.49, 2.39, ...
```

Including Plots

You can also embed plots, for example:

```
diamonds %>%
  ggplot(aes(x = carat,
             y = price,
             color = cut)) + geom_point()
```



Python & Jupyter vs. R & RStudio

Python is a neat language and a great tool to combine with R.

It's also helpful to know a little Python even if we're primarily working in R, because the benefits of being multilingual extend beyond just spoken languages, too.

We recommend R/RStudio because we wrote this book for people who have a data file and specific questions (or general curiosity). Thus, the entry point for our audience into data science is *with data they need to analyze*, and this is what R was made to do.

Additional reasons for using R & RStudio

Below are a few more reasons you should consider using R/RStudio in case you're still on the fence.

You can focus on your work

'The only factor becoming scarce in a world of abundance is human attention' – Kevin Kelly in Wired

We recommend R/RStudio because of the time saved by switching between software applications. For example, when I was in graduate school, I'd have *a minimum of five applications open* to do data analysis. I would be using Word to write, Stata for statistics, Excel to create tables, the browser for internet research, and Adobe Acrobat for reading PDFs. That means I needed to learn five different GUIs, each with their design characteristics.

Each software application cost me valuable neurons whenever I had to switch between them (read more about attentional residue in the footnotes). With R/RStudio, I cut this number to two (RStudio and the browser).

RStudio gives us a better mental model

The third reason is the design of the IDE itself. RStudio is a complementary cognitive artifact, something described in [this article from David Krakauer⁸¹](#),

"[complementary cognitive artifacts are] certainly amplifiers, but in many cases, they're much, much more. They're also teachers and coaches...Expert users of the abacus are not users of the physical abacus—they use a mental model in their brain. And expert users of slide rules can cast the ruler aside having internalized its mechanics. Cartographers memorize maps, and Edwin Hutchins has shown us how expert navigators form near symbiotic relationships with their analog instruments."

These are in contrast to competitive cognitive artifacts, which is what a GUI does.

"In the case of competitive artifacts, when we are deprived of their use, we are no better than when we started. They're not coaches and teachers—they are serfs."

dialog boxes, and point-and-click options.

RStudio does not remove the complexity of doing data analysis, writing blog posts, building applications, debugging code, etc. Instead, it creates an environment where you can do each of these tasks without having them abstracted away from you into drop-down menus,

There have been considerable efforts from the scientists at RStudio to create an environment

⁸¹<http://nautil.us/blog/will-ai-harm-us-better-to-ask-how-well-reckon-with-our-hybrid-nature>

and ecosystem of tools (called packages) to make data analysis less painful (and even fun). We’re confident you’ll find it helps you think about the inputs and outputs of your work in productive and creative ways.

Part 3: Setting up your data science project

"If you can't describe what you are doing as a process, you don't know what you're doing" - W. Edwards Deming

In the last chapter, we recommended a workbench (RStudio) and a set of tools (R, Git, Github). Now we'll use an example project to show how combining these tools create a durable and adaptive workflow. We want to get started with an example early because having a job to do allows us to cover project organization.

Our statistical coursework never covered the details of setting up a project (and we often marvel at how much time we wasted trying to find our files). The way we set our projects up—how we organize files and folders—will directly contribute to our ability to be productive. You've probably discovered it's hard to get things done in a messy office? Well, it will be hard to do data science if we don't organize our files in a logical way that helps us get things done.

Example: FiveThirtyEight's 2019 Presidential debate project

I read something on the internet, got curious, and decided I wanted to dig a little deeper.

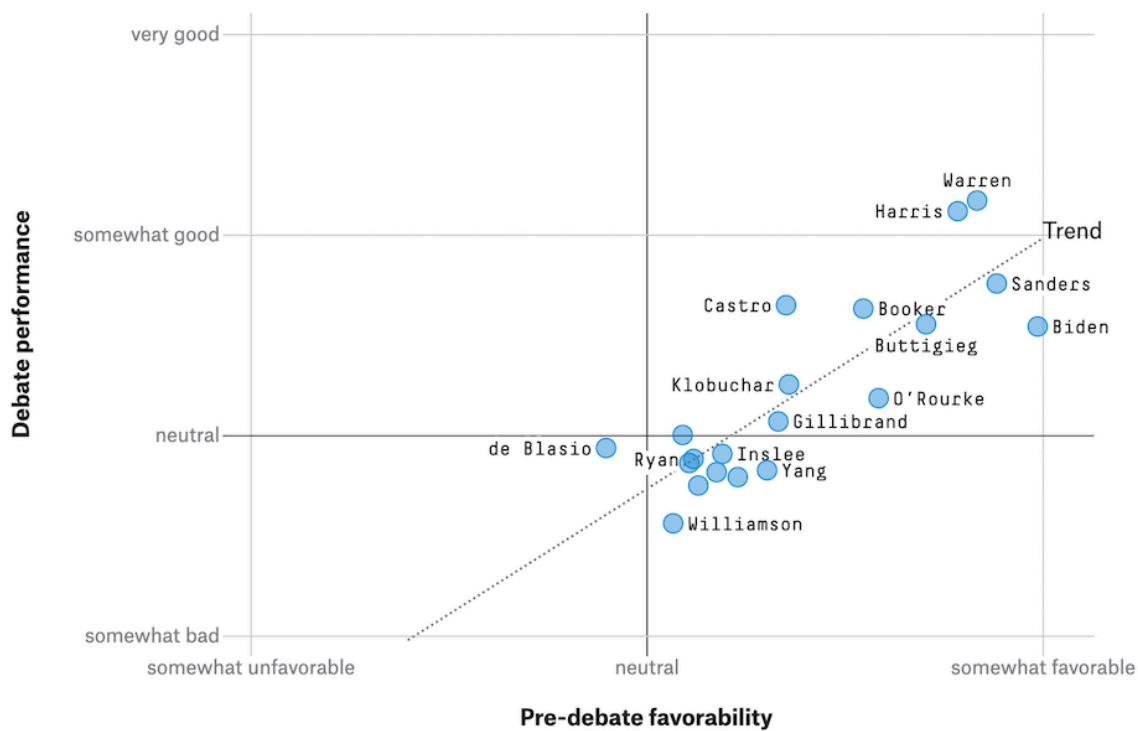
The scenario we've described above might seem vague, but we want to show the power of these tools. We've also found some of the most exciting data projects are born from basic curiosity.

In this case, let's imagine we read something about the first round of the [2019 Democratic Presidential Debates⁸²](#), but we missed all the news coverage.

We stumbled across an article on the data journalism website [fivethirtyeight⁸³](#), and it displayed an image showing the relationship between how voters felt about the candidates before the debates, and how the candidates did in the debate.

⁸²https://en.wikipedia.org/wiki/2020_Democratic_Party_presidential_debates_and_forums

⁸³<https://projects.fivethirtyeight.com/democratic-debate-poll/>



source: <https://projects.fivethirtyeight.com/democratic-debate-poll/>

Wanting to be informed citizens—and knowing how to collect and analyze data—we decide to investigate how each candidate performed using various sources of data.

Data journalism

Journalists are a bit like statisticians in the sense that both get to “play in everyone’s backyard”⁸⁴. Data journalists explicitly combine analysis and communication skills. Marrying these two skills makes data journalism an extraordinary place to look for tools and methods to adapt to different projects.

The best data journalism projects combine the rigor of numbers and math with an ability to **write something people want to read**. Data journalists like [Aleszu Bajak](#)⁸⁵, [Andrew Flowers](#)⁸⁶, and [Andrew Ba Tran](#)⁸⁷ have been hugely influential in introducing R as a tool in the newsroom.

Another reason journalism is an excellent resource for sharing your work is that journalists are trained to view the world differently than typical scientists or analysts. As the NBC investigative reporter [Andy Lehren](#)⁸⁸ describes in the text [Digital Investigative Journalism](#)⁸⁹,

⁸⁴<https://www.nytimes.com/2000/07/28/us/john-tukey-85-statistician-coined-the-word-software.html>

⁸⁵<https://twitter.com/aleszubajak>

⁸⁶<https://twitter.com/andrewflowers>

⁸⁷<https://twitter.com/abtran>

⁸⁸<https://twitter.com/lehrennbc>

⁸⁹<https://www.springer.com/gp/book/9783319972824>

“Journalists can approach data differently than those more trained in computer sciences. Take, for instance, matching databases. Traditional IT managers compare data sets that were designed to talk with each other. Journalists may wonder if the payroll list of school teachers includes registered sex offenders.”

Journalists can communicate *why something matters*, which is a great skill to hone. Explaining a data project to someone with zero domain expertise (or data science knowledge) is a great way to practice your communication skills. It’s also a great way to make sure you’re thinking about an audience for the project.

A collection of modern data sources

To demonstrate how powerful R and RStudio can be, we are going to combine data from four different sources. Each source represents a different way to access data in using R and RStudio.

- 1) The [gtrendsR⁹⁰](#) package for R gives us access to Google search terms and trends. We’re going to import data from Google searches before and after the night of the debates.
- 2) [rtweet⁹¹](#) package in R can be used to download Twitter data but takes a few steps to get set up. Fortunately, we’ve written a tutorial [here⁹²](#) and the package has excellent documentation (see [here⁹³](#) and [here⁹⁴](#)).
- 3) There is a [Wikipedia⁹⁵](#) page dedicated to the debates. We’ll be scraping the tables with airtime for a candidate using the [xml2⁹⁶](#) and [rvest⁹⁷](#) packages.
- 4) Finally, we also have some data from voters on how they felt about each democratic candidate going into the debates. These data are in a [Google Sheet⁹⁸](#), and we’ve used the [datapasta⁹⁹](#) package and copy + paste these data into R. Another option is the [googlesheets4¹⁰⁰](#) package in R (*you will need to copy this sheet into your Google drive to get this data set*).

Step 1: Github

In this example, we will be using an RStudio.Cloud environment to perform the analyses. All of these steps can be accomplished using the RStudio IDE on your local desktop, too.

⁹⁰<https://github.com/PMassicotte/gtrendsR>

⁹¹<https://rtweet.info/>

⁹²<http://www.storybench.org/get-twitter-data-rtweet-r/>

⁹³<https://rtweet.info/articles/auth.html>

⁹⁴<https://rtweet.info/articles/intro.html>

⁹⁵https://en.wikipedia.org/wiki/2020_Democratic_Party_presidential_debates_and_forums

⁹⁶<https://cran.r-project.org/web/packages/xml2/index.html>

⁹⁷<https://rvest.tidyverse.org/>

⁹⁸<http://bit.ly/2YEVASu>

⁹⁹<https://cran.r-project.org/web/packages/datapasta/README.html>

¹⁰⁰<https://googlesheets4.tidyverse.org/>

Head over to [Github](#) and sign up¹⁰¹ for a free account.

The screenshot shows the GitHub 'Join GitHub' page. At the top, there's a navigation bar with links for 'Why GitHub?', 'Enterprise', 'Explore', 'Marketplace', 'Pricing', a search bar, and a 'Sign in' button. Below the navigation is a large heading 'Join GitHub' with the subtext 'The best way to design, build, and ship software.' The page is divided into three main sections: 'Step 1: Set up your account', 'Step 2: Choose your subscription', and 'Step 3: Tailor your experience'. The first section contains fields for 'Username', 'Email address', and 'Password'. The second section lists 'Unlimited public repositories' and 'Unlimited private repositories'. The third section highlights 'Limitless collaboration', 'Frictionless development', and 'Open source community'. A large 'Join Github' button is centered at the bottom of the form.

After you've completed the necessary forms (*remember you only need a free account!*), you should see a page with a message telling you “You don’t have any public repositories yet”.

The screenshot shows a GitHub profile overview page. At the top, there are tabs for 'Overview', 'Repositories 0', 'Projects 0', 'Stars 0', 'Followers 1', and 'Following 0'. Below the tabs, there's a section for 'Popular repositories' and a prominent message: 'New Github account with no repositories'. A large box below states 'You don't have any public repositories yet.' At the bottom, there's a link 'Github profile overview'.

Step 2: RStudio.Cloud

We will eventually create our repositories, but for now, let's head over and use our Github account to sign in to RStudio.Cloud¹⁰². After we're all signed in, we will see the screen below:

¹⁰¹<https://github.com/join>

¹⁰²<https://rstudio.cloud/>

The screenshot shows the RStudio.Cloud interface. On the left, there's a sidebar titled 'Your Workspaces' with options like 'New Space', 'Your workspaces', 'Learn', 'Guide', 'What's New', 'Primer', 'Cheat Sheets', and 'Feedback and Questions'. The main area is titled 'Your Projects' with a 'Create new projects' button. To the right, there's an 'Options' panel with 'Search Projects' and 'Sort Projects' (set to 'By name'). Below that is a 'Capacity' section stating 'This is your personal workspace, where you can create a virtually unlimited number of projects.' A 'Resources to learn more' box contains a link to the 'Guide'.

RStudio.Cloud environment

We've outlined the various resources, projects, and workspaces in the image above (we will go over each in more detail in a later section).

Setting up the panes in your workbench

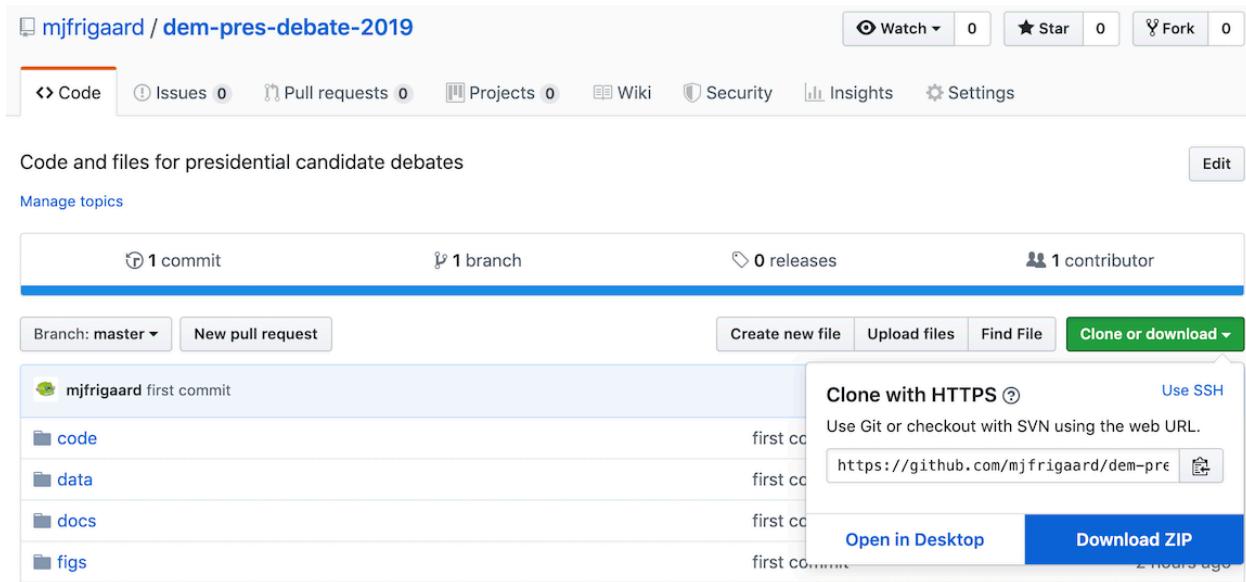
RStudio displays three panes after opening: the **Console** pane, the **Files** pane, and the **Environment** pane. If we open a new file (.R, .txt or .Rmd), RStudio opens the **Source** pane. This is the standard four-pane view for using RStudio (shown below).

The screenshot shows the RStudio interface with four panes. The top-left pane is the 'Console' showing R version 3.6.0 output. The top-right pane is the 'Environment' pane showing 'Global Environment' with 'Environment is empty'. The bottom-left pane is the 'Source' pane with an 'R Script' tab open. The bottom-right pane is the 'Files' pane with tabs for 'Plots', 'Packages', 'Help', and 'Viewer'.

For now, we are going to download a repository from Github and open it in RStudio.Cloud.

Step 3: Download a repository from Github

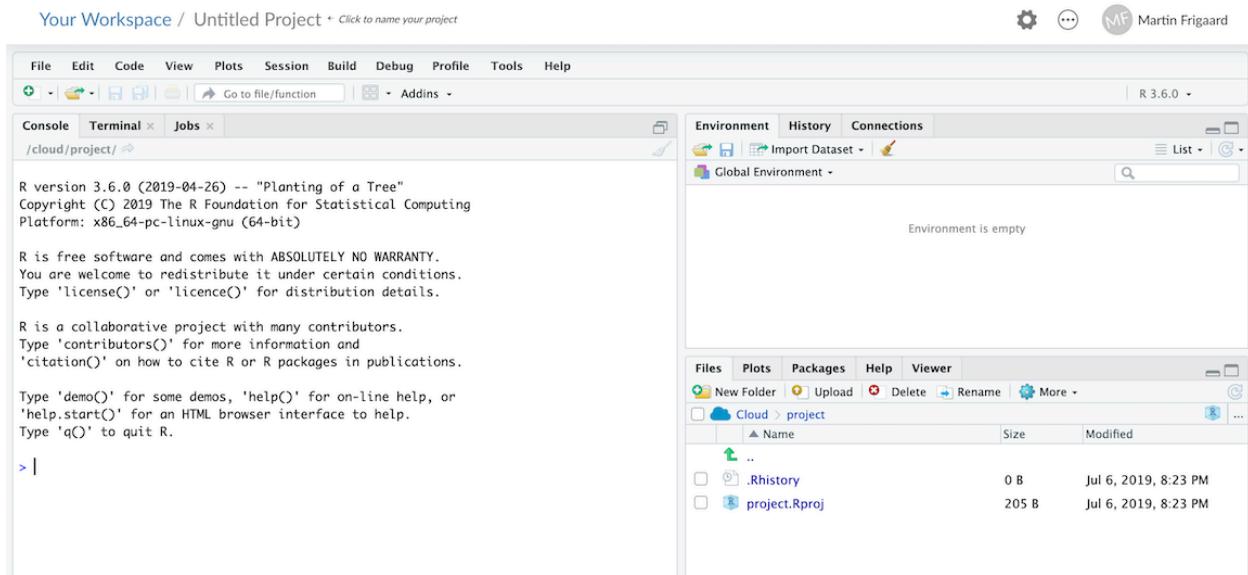
Most of the repos on Github are free for us to download and use. We can do this by clicking on the green **Clone or download** button and click **Download ZIP**.



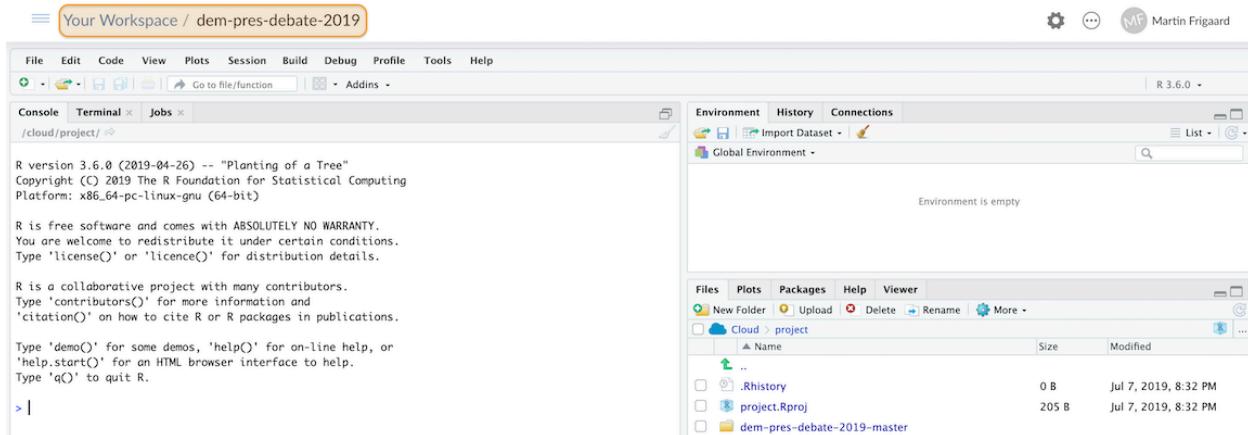
Pick a location on your computer to put your project and download the zipped Github folder.

Step 4: Upload files into RStudio.Cloud

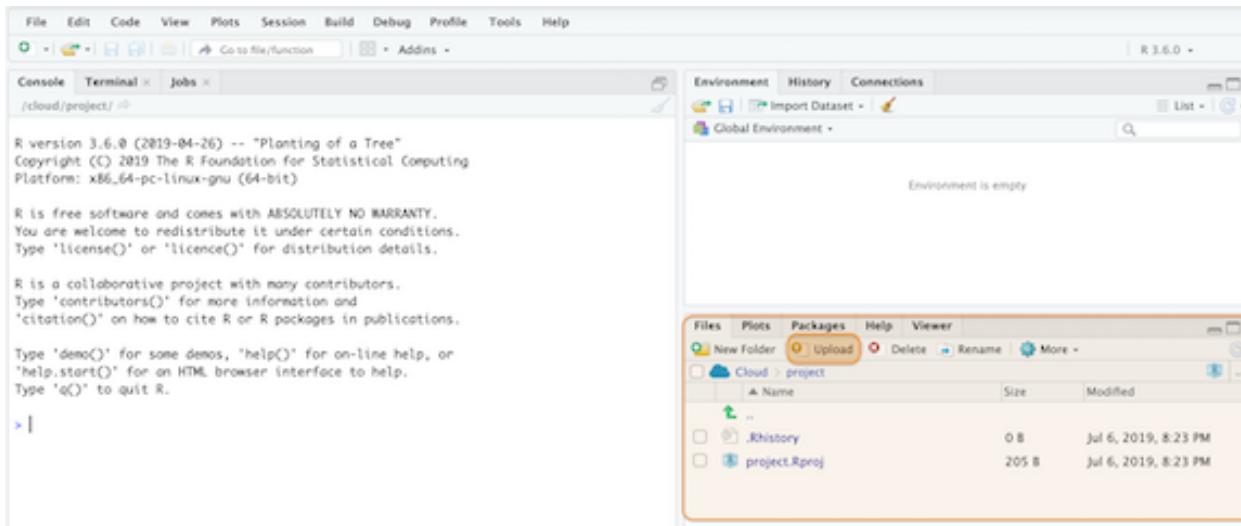
Back in the RStudio.Cloud browser, we're going to click on the **New Project** button. It should display the RStudio IDE in the browser like the image below:



We are going to change the name of this **Untitled** project to **dem-pres-debate-2019**. The results should look like the image below:



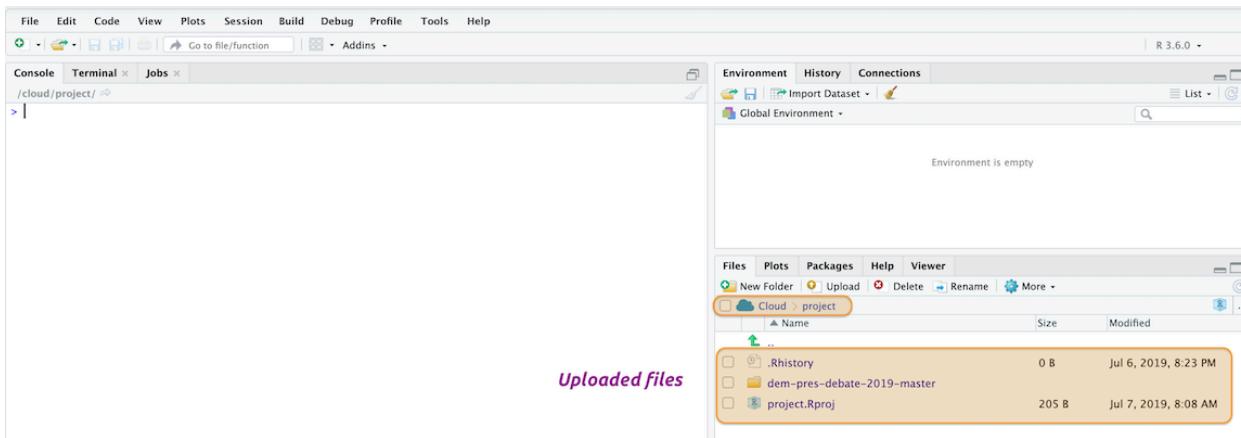
Look at the **Files** pane in the lower right corner and click on the **Upload** button, then click on **Choose files** and locate the recently downloaded zipped Github folder. Upload this file into the RStudio.Cloud project workspace.



upload button

Accessing files in RStudio.Cloud

Uploading these files might take some time, but when everything is in RStudio.Cloud, we'll see the `dem-pres-debate-2019-master` folder in the **Files** pane.



the newly uploaded files

I unzipped the file we uploaded and created a folder called `dem-pres-debate-2019-master`. Unfortunately, it put this folder *inside* the `cloud/project` folder. We wanted to upload the *contents* of the `dem-pres-debate-2019-master` file into the `cloud/project` folder (and not the folder itself).

```
Cloud/project/ # here ←
    └── dem-pres-debate-2019-master # move these contents...
        └── project.Rproj
```

We are going to use this opportunity to introduce a few command-line tools. To do this, we'll be working from the **Terminal** pane in RStudio.Cloud. The next session will be a quick 'crash course' on operating systems, some **Terminal** commands, and how they work together.

The Command line: Unix and Windows

In 2007, Apple released its [Leopard¹⁰³](#) operating system that was the first to adhere to the [Single Unix Specification¹⁰⁴](#). I only introduce this bit of history to help keep the terminology straight. macOS and Linux are both Unix systems, so they have a similar underlying architecture (and philosophy). Most Linux commands also work on macOS.

Windows has a command-line tool called Powershell, but this is not the same as the Unix shells discussed above. The differences between these tools reflect the differences in design between the two operating systems. However, if you're a Windows 10 user, you can install a [bash shell command-line tool¹⁰⁵](#).

Command-line interfaces

The [command line interface¹⁰⁶](#) (CLI) was the predecessor to a GUI, and there is a reason these tools haven't gone away. CLI is a text-based screen where users interact with their computer's programs, files, and operating system using a combination of commands and parameters. This basic design might make the CLI sound inferior to a trackpad or touchscreen, but after a few examples of what's possible from on the command-line and you'll see the power of using these tools.

"What am I getting out of this?"

That's a fair question—being able to use the command line gives us more 'under-the-hood' access to any computer. We can use the command line to navigate a computer's directories (folders), install new programs or libraries, and track changes to files. It might seem clunky and ancient, but people keep this technology around because of it's 1) specificity and 2) modularity (also the two features that make Unix programs so powerful). What do we mean by this?

- [Specificity¹⁰⁷](#) means each Unix command or tool does one thing very well (or [DOTADIW¹⁰⁸](#))
- [Modularity¹⁰⁹](#) is the ability to mix and match these tools together with 'pipes,' a kind of grammatical glue that allows users to expand these tools in seemingly endless combinations

Having these skills have also made us more comfortable when we've had to interact with remote machines or different operating systems (Linux, per se). We will work through an example to demonstrate some of these features.

¹⁰³https://en.wikipedia.org/wiki/MacOS_version_history#Version_10.5:_%22Leopard%22

¹⁰⁴https://en.wikipedia.org/wiki/Single_UNIX_Specification

¹⁰⁵<https://www.windowcentral.com/how-install-bash-shell-command-line-windows-10>

¹⁰⁶https://en.wikipedia.org/wiki/Command-line_interface

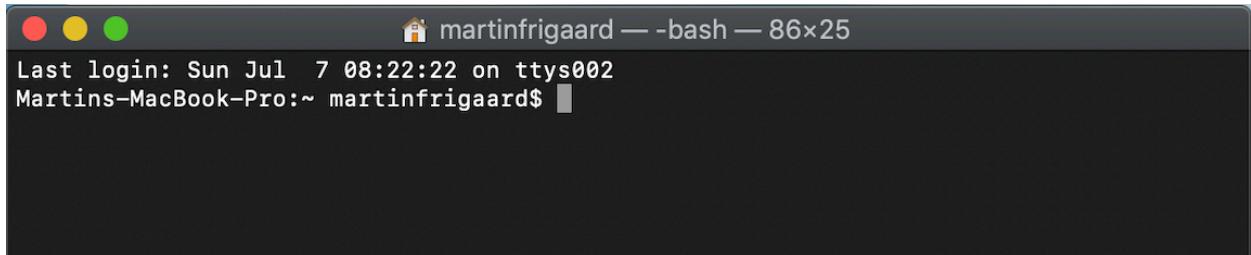
¹⁰⁷<https://www.dictionary.com/browse/specific>

¹⁰⁸https://en.wikipedia.org/wiki/Unix_philosophy#Do_One_Thing_and_Do_It_Well

¹⁰⁹https://en.wikipedia.org/wiki/Modularity#Table_1:_The_use_of_modularity_by_discipline[34]

The Terminal (macOS)

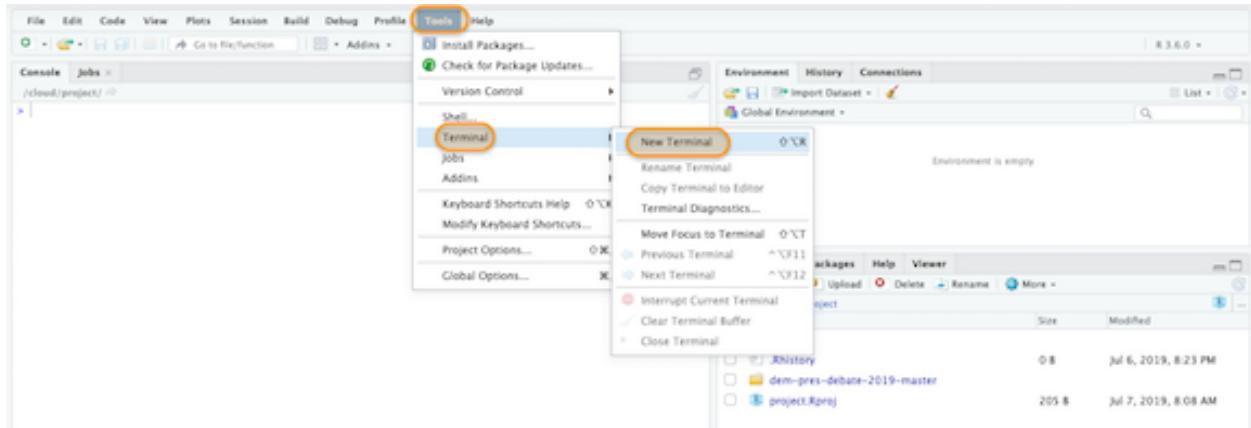
Below is an image of what the terminal application looks like on macOS. On Macs, the Terminal application runs a [bash shell¹¹⁰](#), which is why you can see the `bash -- 86x25` on the top of the window. Bash is a commonly-used shell, but there are other options too (see [Zsh¹¹¹](#), [tcsh¹¹²](#), and [sh¹¹³](#)). *Fun fact: bash is a pun for the sh shell: bourne-again shell.*



The Terminal is an emulator application for Mac users. Terminal is available as an application under the **Applications > Utilities > Terminal**.

The Terminal (RStudio)

The Terminal pane is also available in RStudio under **Tools > Terminal > New Terminal**.



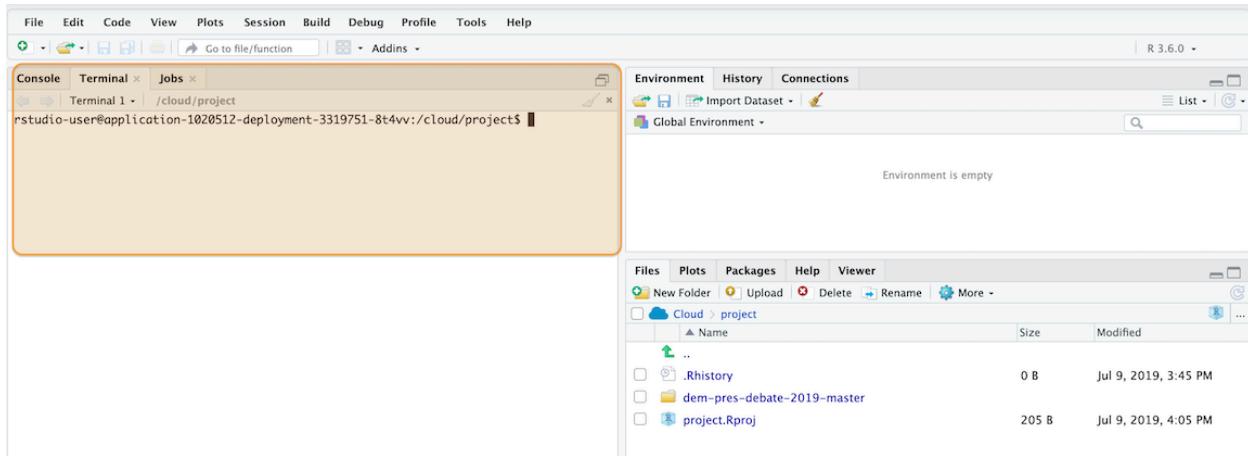
The Terminal pane will open in the same window as the Console pane.

¹¹⁰[https://en.wikipedia.org/wiki/Bash_\(Unix_shell\)](https://en.wikipedia.org/wiki/Bash_(Unix_shell))

¹¹¹<http://zsh.sourceforge.net/>

¹¹²<https://en.wikipedia.org/wiki/Tcsh>

¹¹³https://en.wikipedia.org/wiki/Bourne_shell



Now we will get some practice organizing our data science project using the command line.

Good enough command-line tools

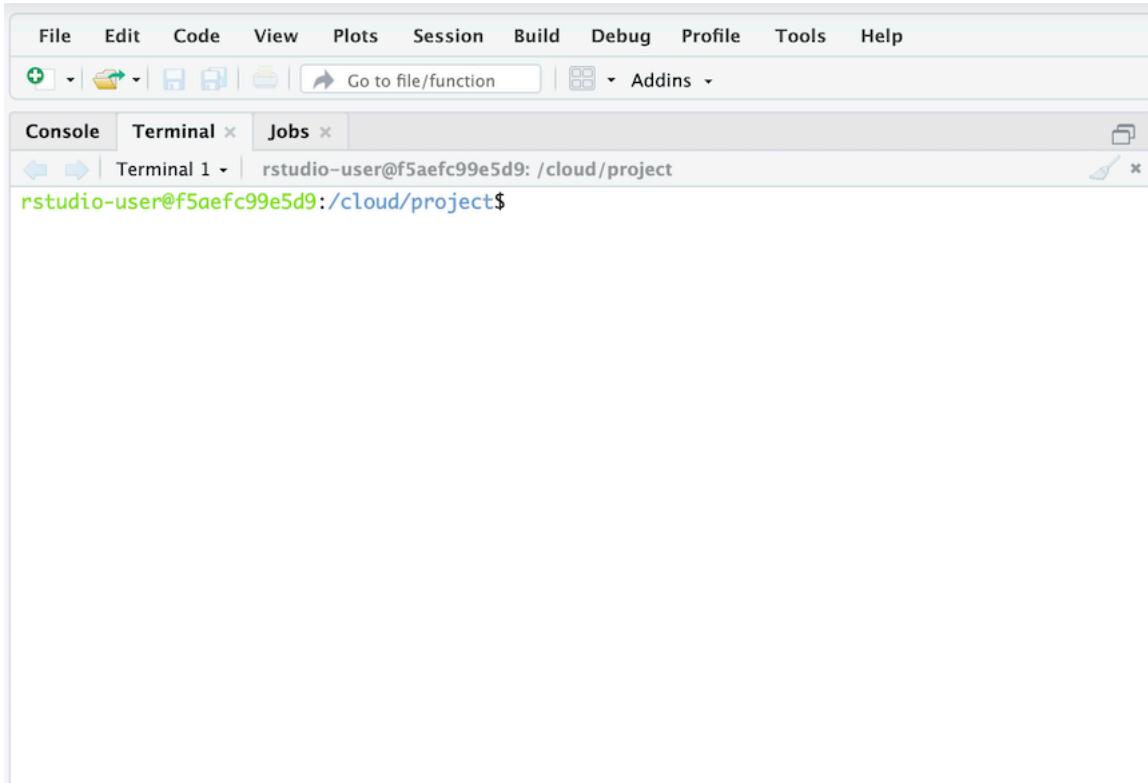
FAIR WARNING—command-line interfaces can be frustrating. Computers don't behave in ways that are easy to understand (that's why GUIs exist). Switching from a GUI to a CLI seems like a step backward at first, but the initial headaches pay off because of the gains we'll have in control, flexibility, automation, and reproducibility.

Here is a quick list of commonly used Terminal commands.

- **pwd** - print working directory
- **cd** - change directories
- **cp** - copy files from one directory to another
- **ls** - list all files
- **ls -la** - list all files, including hidden ones
- **mkdir** - make directory
- **rmdir** - remove a directory
- **cat** - display a text file in Terminal screen
- **echo** - outputs text as arguments, prints to Terminal screen, file, or in a pipeline
- **touch** - create a few files
- **grep** - “globally search a regular expression and print”
- **>>** and **>** - redirect output of program to a file (don't display on Terminal screen)
- **sudo** and **sudo -s** (**BE CAREFUL!!**) performing commands as **root** user can carry some heavy consequences.

Command-line skill #1: who is using what?

After downloading the files from Github, we've uploaded the zipped folder into the Cloud/project. In the RStudio.Cloud Terminal pane, we should see something like this:



The figure above might look like gobbledegook at first, but command-line interfaces have a recognizable pattern if we know what we're looking for:

- First, we can almost always expect some `user@machine` identifier to tell us who we're signed in as and on what machine
- Second, there's usually some way of displaying the `home` directory. In this case, it's the stuff between the colon (:) and the dollar sign \$ (`/cloud/project`)

Let's check a few things to help figure out what's going on.

```
1 rstudio-user@f5aefc99e5d9:/cloud/project$ whoami  
2 rstudio-user
```

We are the `rstudio-user` on this machine `f5aefc99e5d9`. The same information on a local MacBook laptop might look like this:

```
1 Martins-MacBook-Pro:~ martinfrigaard$ whoami  
2 martinfrigaard
```

In this case, the machine information would be `Martins-MacBook-Pro` and the location to be the `home` directory `~` (the top-level folder) for the user `martinfrigaard`.

Command-line skill #2: where am I?

In the RStudio.Cloud Terminal pane, enter the print working directory (`pwd`) command:

I've omitted everything preceding the prompt (\$) for easier printing

```
1 $ pwd  
2 /cloud/project
```

`pwd` tells us where we are, otherwise known as the current working directory. Imagine the current working directory as the spot we're standing, and file path `/cloud/project` as the way back to our `root` folder.

To get a sense of our surroundings lets list the files in `/cloud/project` using `ls`

```
1 $ ls  
2 dem-pres-debate-2019-master project.Rproj
```

We can see the folder (`dem-pres-debate-2019-master`) and the RStudio project file (`project.Rproj`). On a side note, it's always a good idea to pay attention to file extensions (`.Rproj`, `.R`, `.md`, etc.), because different files interact with the Terminal in different ways.

Command-line skill #3: moving around

Now that we know where we are, and what files and folders are in here with us, we can start to stretch our legs and move around. Let's start by changing directories `cd` to the `dem-pres-debate-2019-master` folder, then check with `pwd`.

```
1 $ cd dem-pres-debate-2019-master  
2 $ pwd  
3 /cloud/project/dem-pres-debate-2019-master
```

Now we can check the files in this new directory with `ls`

```

1 $ ls
2 README.Rmd README.md code data
3 dem-pres-debate-2019.Rproj figs project.Rproj

```

The output from `ls` shows me there are three sub-folders in the `dem-pres-debate-2019-master` folder (`code`, `data`, `figs`), three `.Rmd` files, one `README.md` file, and one `.Rproj` file.

Now that we've moved into this folder and looked around let's climb back out of it. We can always move up one folder by executing the `cd ..` command.

```

1 $ cd ..
2 $ pwd
3 cloud/project

```

Let's move back into `dem-pres-debate-2019-master` using `cd` again, but this time, we will move up one folder using `cd /cloud/project`.

```

1 $ cd /cloud/project
2 $ ls
3 dem-pres-debate-2019-master project.Rproj

```

We can also check the files in `dem-pres-debate-2019-master` using `ls` and the folder name.

```

1 $ ls dem-pres-debate-2019-master
2 README.Rmd README.md code data
3 dem-pres-debate-2019.Rproj figs project.Rproj

```

We can add the `-F` option to the end of the command to tell Terminal to list the files in the folder at the end of the file path.

```

1 ls dem-pres-debate-2019-master -F
2 README.Rmd README.md code/ data/
3 dem-pres-debate-2019.Rproj figs/ project.Rproj

```

Now we can see the folders have a `/` forward slash at the end of their name to separate them from the other files.

Absolute vs. relative file paths

An **absolute file path** starts at the root directory (`~` or `\`) and follows along the path, folder by folder, until it lands on the last folder or file.

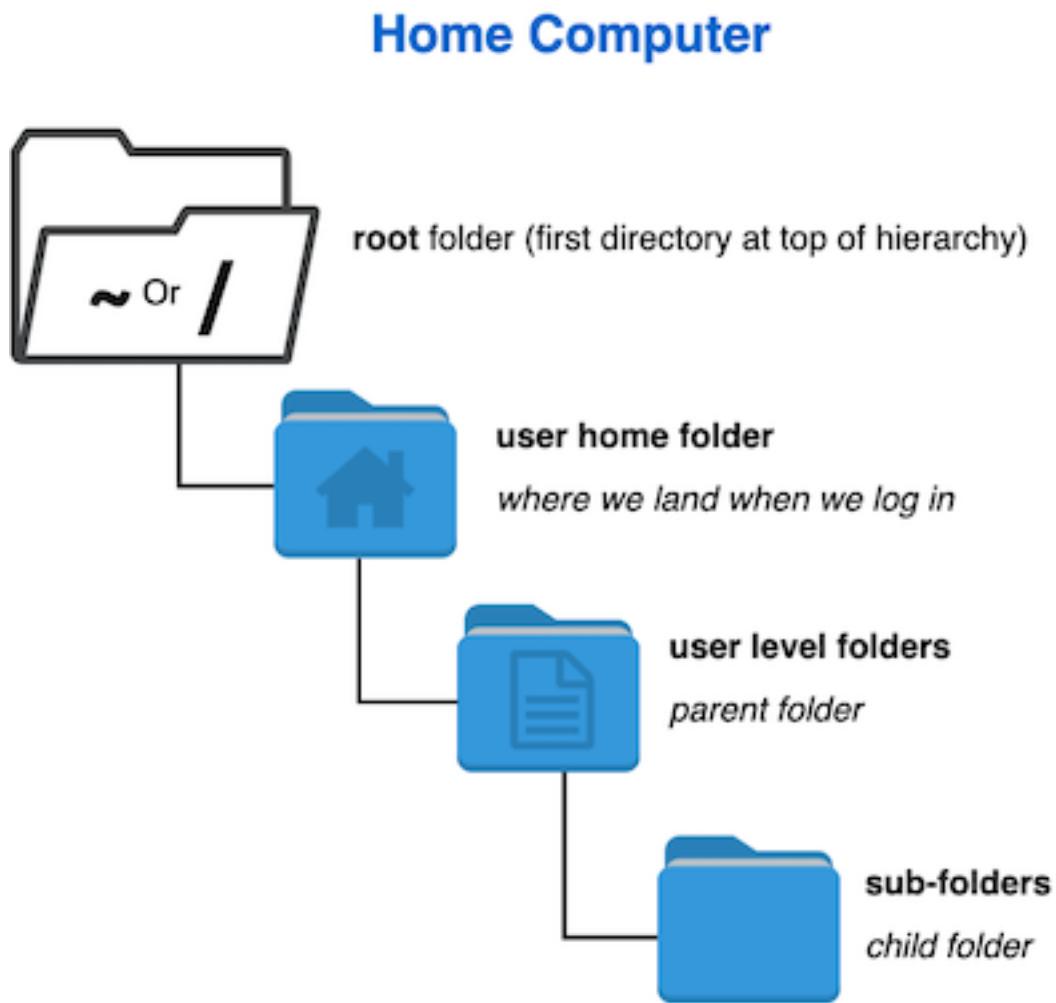
`/start/from/absolutely/where/i/tell/us`

A **relative file path** starts at a folder but leaves the rest 'relative' to wherever that folder is located.

`start/from/wherever/we/put/me`

Folder structure

Below is an example folder tree structure on a macOS. The structure won't be too different on a Windows operating system, but the folder paths will have a different syntax.



The **root** folder is the “uppermost” location of this machine’s folders and files. In macOS, root is represented with a tilde (~). In Windows, the root folder is located with the forward-slash (/). If we have the right privileges, we can log in as the root user, and the prompt will change from \$ to # (be careful here!)

When we log into a computer, we start in a **home** folder (usually with a shorter version of that user’s full name they used to set up their operating system). The **home** folder is the typical “starting point” for that user’s folders and files. If we are working on macOS, this is the folder with a little house on it.

Depending on the operating system, this location starts with some standard default folders (Desktop,

Documents, Downloads, and Applications)

Special considerations: Windows machines

On Windows machines, the file path to dem-pres-debate-2019-master might look like this:

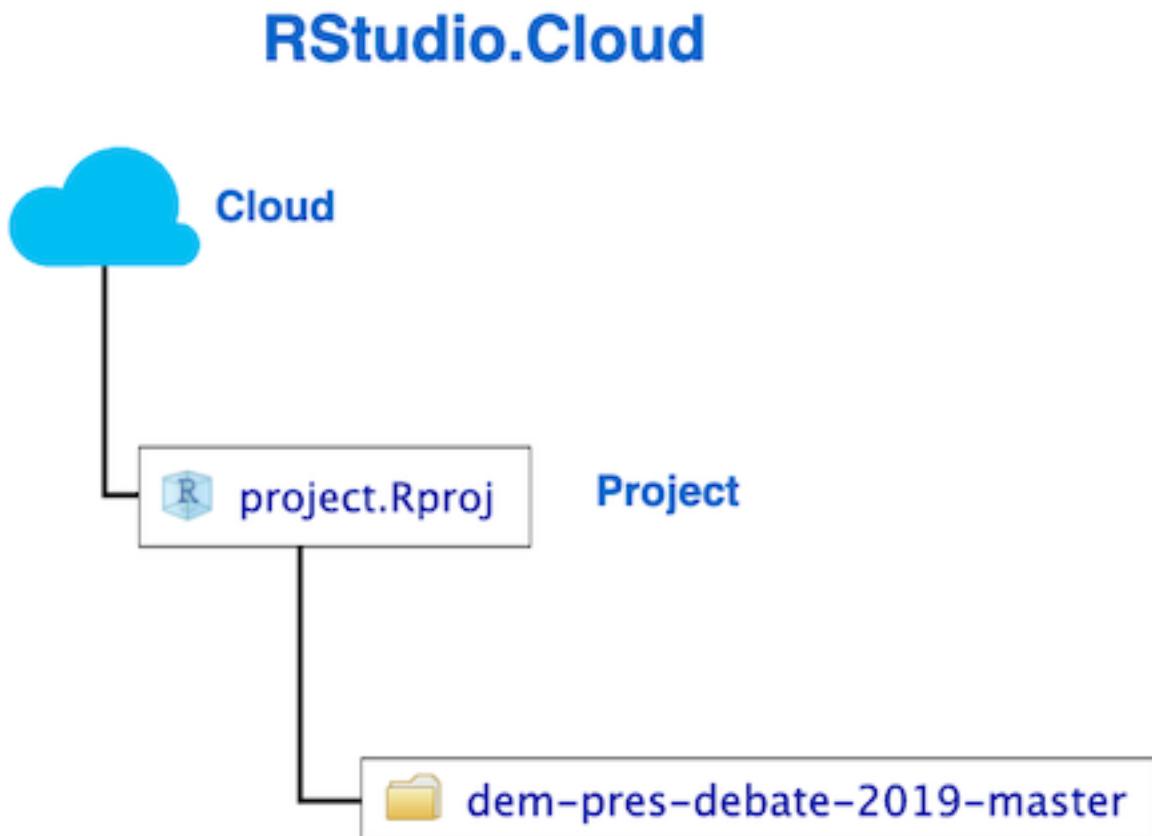
```
1 C:\Users\martinfrigaard\Documents\dem-pres-debate-2019-master
```

But we would need to write it like this:

```
1 C:\\\\Users\\\\martinfrigaard\\\\Documents\\\\dem-pres-debate-2019-master
```

This odd way of writing file paths is because, in R, the \ is called an escape character, so to navigate through folders we will have to use two backslashes \\.

Below is the folder tree on RStudio.Cloud:



Now, this image might be about as clear as mud, but it'll make more sense when we start moving things around.

Command-line skill #4: moving things around

We're working in RStudio.Cloud, but the GUI representation of our folder structure won't be much different if we were working on our local laptop.

Remember, we want to move the contents of `dem-pres-debate-2019-master` into `cloud/project`. The command for moving files from one place to another is `mv`, but we are going to add two options, `-v` and `*`. There are many other options for using `mv`, read about them [here¹¹⁴](#).

The sequence of commands we'll enter in the RStudio.Cloud Terminal are below:

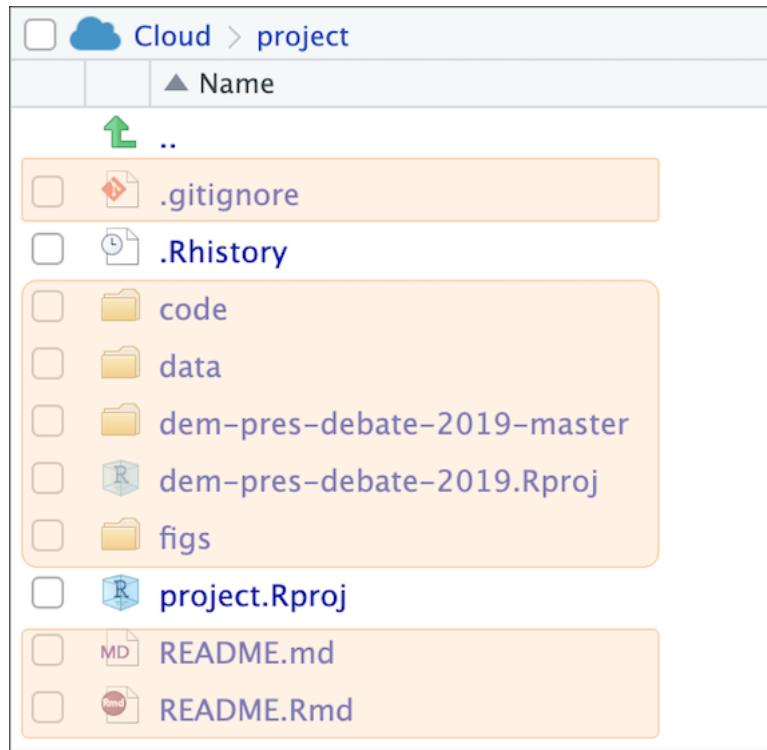
```
1 $ mv -v dem-pres-debate-2019-master/* /cloud/project
```

You will see the following changes in Terminal:

```
1 'dem-pres-debate-2019-master/README.Rmd' -> '/cloud/project/README.Rmd'
2 'dem-pres-debate-2019-master/README.md' -> '/cloud/project/README.md'
3 'dem-pres-debate-2019-master/code' -> '/cloud/project/code'
4 'dem-pres-debate-2019-master/data' -> '/cloud/project/data'
5 'dem-pres-debate-2019-master/dem-pres-debate-2019.Rproj' -> '/cloud/project/dem-pres\
6 -debate-2019.Rproj'
7 'dem-pres-debate-2019-master/figs' -> '/cloud/project/figs'
```

And the following changes in the Files pane:

¹¹⁴https://www.gnu.org/software/coreutils/manual/html_node/mv-invocation.html#mv-invocation



Now we know we've successfully moved all of the files. But we will want to get rid of the old folder, `dem-pres-debate-2019-master`.

Command-line skill #5: Delete things

To delete a folder, we can either use `rmdir` or `rm -R`.

```
1 $ rm dem-pres-debate-2019-master -Ri
2 rm: descend into directory 'dem-pres-debate-2019-master'?
```

This command is helpful because the `i` option tells `Terminal` to check with us before doing anything. Go ahead and enter `n` and try using `rmdir` to delete the `dem-pres-debate-2019-master` folder.

```
1 $ rmdir dem-pres-debate-2019-master
2 rmdir: failed to remove 'dem-pres-debate-2019-master': Directory not empty
```

`Terminal` does its best to save us from ourselves, but that's not always possible. As Doug Gwyn said,

"Unix was not designed to stop its users from doing stupid things, as that would also stop them from doing clever things."

Well, what does `rmdir` actually do then? We can figure this out with `rmdir --help`

```
1 $ rmdir --help
```

This command will print some useful information about the `rmdir` command:

```
1 $ rmdir --help
2 Usage: rmdir [OPTION]... DIRECTORY...
3 Remove the DIRECTORY(ies), if they are empty.
4 # else omitted...
```

Now we know this is not the right tool for the job (the folder isn't empty), so we will use `rm -R` `dem-pres-debate-2019-master`. Each folder and file will prompt a question that needs a response before **Terminal** can delete anything.

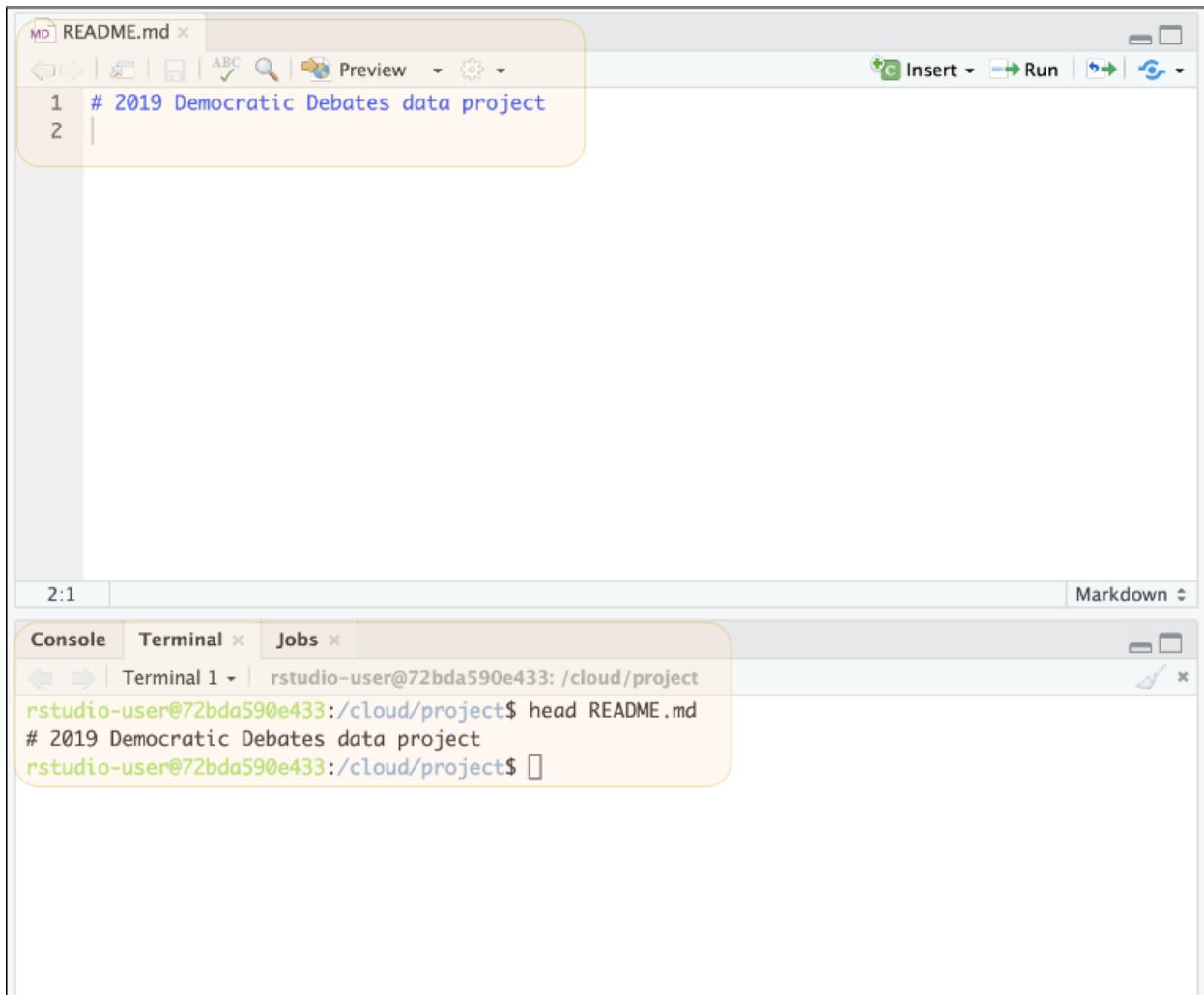
The **Terminal** pane should have the following contents when we're finished:

```
1 $ rm -R dem-pres-debate-2019-master
2 rm: descend into directory 'dem-pres-debate-2019-master'? y
3 rm: remove regular file 'dem-pres-debate-2019-master/.DS_Store'? y
4 rm: remove regular file 'dem-pres-debate-2019-master/.gitignore'? y
5 rm: remove directory 'dem-pres-debate-2019-master'? y
```

Command-line skill #6: Print things

Terminal works very well with plain text format. For example, I can use `head` and the name of a file I want to see.

```
1 $ head README.md
```



As we can see, this is the first few lines of the `README.md`. Markdown is a plain text format so that it will print clearly to the **Terminal** window. In addition to `head`, we can also use the `tail` command to view the bottom of the `README.md` file.

What if we want to see all the contents in `README.md`? Well, before printing all the contents, we want to see how big the file is, and we can do that using `wc` (which stands for “word count”).

```
1 $ wc README.md  
2 # 1 6 39 README.md
```

The three numbers above are the number of lines (1), the number of words (6), and the number of characters (39).

`wc` is telling us that `README.md` won’t be hard to read on the Terminal window. If it was, that’s where the `less` command comes in.

```
1 $ less README.md
```

less will display the contents of README.md in a way that allows us to scroll through the file using the arrow keys. After we're done viewing the file, we can exit less using q.

Another option to print is cat, but this will print the entire contents to the Terminal window, so use wc first to see if that's the best choice.

Command-line skill #7: Create things

Sometimes we might need to create a new file and add some text to it. This skill is handy if we don't have to open any new applications.

The touch command will create a new file (CHANGELOG.txt), and echo will put the "some thoughts" on this file (which we can verify with cat).

```
1 $ touch CHANGELOG.txt
2 $ echo "some thoughts" > CHANGELOG.txt
3 $ cat CHANGELOG.txt
4 some thoughts
```

The > symbol tells Terminal to send echo "some thoughts" to CHANGELOG.txt. When we use cat, we see these commands put "some thoughts" into the top lines of the new file, CHANGELOG.txt.

The CHANGELOG.txt file is for writing notes about changes to our project, but we should add a date to make sure they're listed chronologically. Unix has a date variable we can access using \$(date) (which 'attaches' the output from the command date with "some thoughts"), so we will repeat the process above, but include today's date with \$(date).

```
1 $ echo $(date) "some thoughts" > CHANGELOG.txt
2 $ cat CHANGELOG.txt
```

In Unix systems, we can always access today's date with the date or cal.

```
1 $ cal
2      July 2019
3 Su Mo Tu We Th Fr Sa
4   1  2  3  4  5  6
5   7  8  9 10 11 12 13
6  14 15 16 17 18 19 20
7  21 22 23 24 25 26 27
8  28 29 30 31
```

Command-line skill #8: Combine things

The commands above are great for creating new files and adding new text, but what if CHANGELOG.txt already exists and we wanted to add more thoughts to it? We can do this by changing the > symbol to >>.

```

1 $ echo $(date) "more thoughts" >> CHANGELOG.txt
2 $ cat CHANGELOG.txt
3 # Thu Jul 11 13:45:57 UTC 2019 some thoughts
4 # Thu Jul 11 13:49:42 UTC 2019 more thoughts

```

>> tells Terminal to append the output from echo to CHANGELOG.txt on a new line.

Another powerful tool in the Unix toolkit is the pipe (|). The pipe can be used to ‘direct’ outputs from one command to another. For example, if I wanted to see how many download .R script files are in the code folder, I could use the following:

```
1 $ ls code | grep "download" | less
```

The code above should display the following result:

```

1 00.2-download-538.R
2 00.3-download-google.R
3 00.4-download-tweets.R
4 00.5-download-wikipedia.R
5 (END)

```

We will leave the grep command for you to investigate with --help to figure out what’s happening here. Type q to leave this screen.

Other command line stuff: homebrew

The bash shell on macOS comes with a whole host of packages we can install with [homebrew¹¹⁵](#), the “The missing package manager for macOS (or Linux)”.

(You won’t be able to do this on RStudio Terminal, but there are other options we will list below)

After installing homebrew, we recommend installing the [tree¹¹⁶](#) package.

¹¹⁵<https://brew.sh/>

¹¹⁶<https://brewinstall.org/Install-tree-on-Mac-with-Brew/>

```
1 $ # install tree with homebrew
2 $ brew install tree
3 $ # get a folder tree for this project
4 $ tree
```

The tree command gives us output like the folder tree below.

```
1 └── README.Rmd
2 └── README.md
3 └── code
4     ├── 00.1-inst-packages.R
5     ├── 00.2-download-538.R
6     ├── 00.3-download-google.R
7     ├── 00.4-download-tweets.R
8     ├── 00.5-download-wikipedia.R
9     ├── 01-import.R
10    └── 02-wrangle.R
11   └── data/
12     ├── processed/
13     └── raw
14         └── 538/
15             └── 2019-07-06-Cand538Fav.csv
16         └── google-trends/
17             ├── 2019-07-10-Dems2020Night1Group1.rds
18             └── 2019-07-10-Dems2020Night1Group2.rds
19         └── twitter/
20             ├── 2019-07-06-Night01Tweets.rds
21             ├── 2019-07-06-Night01TweetsRaw.rds
22             ├── 2019-07-06-Night01TweetsUsers.rds
23             ├── 2019-07-06-Night02Tweets.rds
24             ├── 2019-07-06-Night02TweetsRaw.rds
25             └── 2019-07-06-Night02TweetsUsers.rds
26         └── wikipedia/
27             ├── 2019-07-10-WikiDemAirTime01Raw.csv
28             ├── 2019-07-10-WikiDemAirTime02Raw.csv
29             └── 2019-07-25-PollingCriterionRaw.csv
30   └── dem-pres-debate-2019.Rproj
31   └── figs/
32       └── 03-538-night-one-debates.png
33
34 9 directories, 23 files
```

Note that the CHANGELOG.txt file is not included in this tree (because we made these changes on a

local repository). Folder trees come in handy for documenting the project files (and any changes to them).

A note on Terminals vs. Shells: Sometimes you'll hear the term "shell" thrown around when researching command-line tools. Strictly speaking, the Terminal application is not a shell, but rather it *gives the user access to the shell*. Other terminal emulator options exist, depending on your operating system and age of your machine. **Terminal.app** is the default application installed on macOS, but you can download other options (see [iTerm2¹¹⁷](#)). For example, the **GNOME¹¹⁸** is a desktop environment based on Linux which also has a Terminal emulator, but this gives users access to the Unix shell.

Command line recap

We've covered eight command-line tools, and we hope you can see how these can be combined to create very efficient workflows and procedures. By tethering commands together, we can move inputs and outputs around with a lot of flexibility.

More on organizing project files

As we saw above, the `tree` output gave us a printout of the project folder in a hierarchy (i.e. a tree with branches).

The thing to notice is the separation of files into folders titled, `data`, `docs`, and `src` or `code`. We didn't choose these folder names at random—there is a way to organize a data science project. Start with the folder structure outlined by Greg Wilson we [mentioned earlier](#).¹¹⁹ If you already have an organization scheme, we still recommend reading at least [this section¹²⁰](#) of the paper—it's full of great information and links to other resources.

Building command-line skills

This section has been a concise introduction to command-line tools, but hopefully, we've demystified some of the terms for you. The reason these technologies still exist is that they're powerful. Probably, you're starting to see the differences between these tools and the standard GUI software installed on most machines. [Vince Buffalo¹²¹](#), sums up the difference very well,

"the Unix shell does not care if commands are mistyped or if they will destroy files; the Unix shell is not designed to prevent you from doing unsafe things."

The command line can seem intimidating because of its power and ability to destroy the world.

¹¹⁷<https://www.iterm2.com/>

¹¹⁸<https://en.wikipedia.org/wiki/GNOME>

¹¹⁹<https://swcarpentry.github.io/good-enough-practices-in-scientific-computing/#project-organization>

¹²⁰(<https://swcarpentry.github.io/good-enough-practices-in-scientific-computing/#project-organization>)

¹²¹<http://vincebuffalo.org/blog/>

Part 4: Keep track of changes with version control

In the previous sections, we've covered using Terminal in the RStudio. If you're unfamiliar with these topics, please start there. This section will include tracking changes with version control, specifically Git, Github, and RStudio.

Tracing our steps

'Sharing your' can take a few forms. You can finish a project, then share your work for people to see and use in what they're doing.

Another option is to share what we're currently working on in a way that allows other people can collaborate with us along the way.

To accomplish the second option, we need a means showing how our work has changed over time. For example, maybe we've used the 'Review' tools in Microsoft Word, or we've collaborated in a Google sheet document. Both are types of [version control](#)¹²² because they're a formal system of managing changes to information over time.

Consider a typical lifespan of a .docx file:

```
1 project/
2 └── paper/
3     ├── draft-1.docx
4     ├── draft-2.docx
5     ├── draft-3-mf.docx
6     ├── draft-3-mk.docx
7     ├── final.docx
8     ├── final-final.docx
9     └── no-really-final.docx
```

This is the inevitable progression of a document that requires collaboration *and* iteration. Revision features are great when we're using a single file, and that file should not be altered until changed have been approved/rejected. But what about all the supporting files that are needed to create that report or analysis in our Word documents? Every project we've ever worked on required more than a single Word document.

¹²²https://en.wikipedia.org/wiki/Version_control

2 Childhood overweight and obesity rates have tripled nationally over the past three decades^{1,2}.

3 The prevalence rates of overweight and obesity are distinctly higher in low-income and ethnic minority

4 children^{3,4}. The CDC considers children overweight or obese with a body mass index (BMI)-for-age equal

5 to or greater than the 85th and 95th percentile, respectively⁵. The recent reports from school health

6 assessments revealed higher rates of childhood overweight and obesity in California schools compared

7 to national rates⁶. The 2008 Pediatric Nutrition Surveillance Survey (PedNSS) , identified the combined

8 prevalence of overweight and obesity in Butte County as 47% for Latinos, African Americans, and Native

9 Americans low income children and adolescents aged 9-11⁷.

10 High BMI in youth is associated with a number of comorbidities, including increased risk for

11 hypertension^{8, 9, 10}. The diagnostic criteria for childhood hypertension is an average blood pressure (BP)

12 reading greater than or equal to the 95th percentile for sex, age and height on at least three separate

13 occasions. Children that have an average BP reading greater than or equal to the 90th percentile for sex,

14 age and height on at least three separate occasions or BP readings greater than or equal to an absolute

15 value of 120/80 mmHg are considered pre-hypertensive¹¹. Elevated BP (EBP) and pre-elevated BP (pre-

16 EBP) consist of any measurements that have been taken, but do not meet the necessary diagnostic

17 criteria.

We can see the proposed changes, what they are, and who made them.

Frigaard, Martin
Deleted: the mandatory adolescent

Frigaard, Martin
Maybe remove this sentence and add the reference to the previous statement?

Frigaard, Martin
Deleted: an instrument that specifically tracks the health of low-income U.S. children and adolescents

Frigaard, Martin
Deleted: Childhood hypertension is defined as having

Frigaard, Martin
Deleted: Pre-hypertension in

Frigaard, Martin
Deleted: c

Frigaard, Martin
Deleted: is defined as

Frigaard, Martin

Revision history can track changes in a specific document, but what about all the other files in this project?

The file displayed above is an early version of a peer-reviewed manuscript. In this situation, a coauthor has suggested some changes to our results section. While making changes, the document is essentially ‘locked’ until the next coauthor can review. This is a serious limitation when a project involves many files, and those files require frequent updates, or *versions*.

Enter version control

Sound version control systems let us see four important aspects of change:

1. What the differences are (*i.e. the actual changes*),
2. Who made the changes (*this allows us to separate our changes from everyone else's*),
3. The time/date of the proposed changes (*without this, it's impossible to know what order the changes took place in*), and
4. Any comments about the change (*the 'how and why' for the changes*)

Unfortunately, tracked changes in Microsoft Word only applies to a single document at a time. When we’re working collaboratively on a project (which is quite often), we know asking someone to change a single sentence can result in changes to dozens of files. That’s why we need a way to track changes across multiple files in a project.

Git

Git is a [version control system¹²³](#) (VCS), which is somewhat like the *Tracked Changes* in Microsoft Word or the *Version History* in Google Docs, but extended to every file in a project. Git will help us keep track of our documents, datasets, code, images, and anything else we tell it to keep an eye on.

Why use Git?

You will eventually ask yourself, *why am I subjecting myself to this—is there another way?*

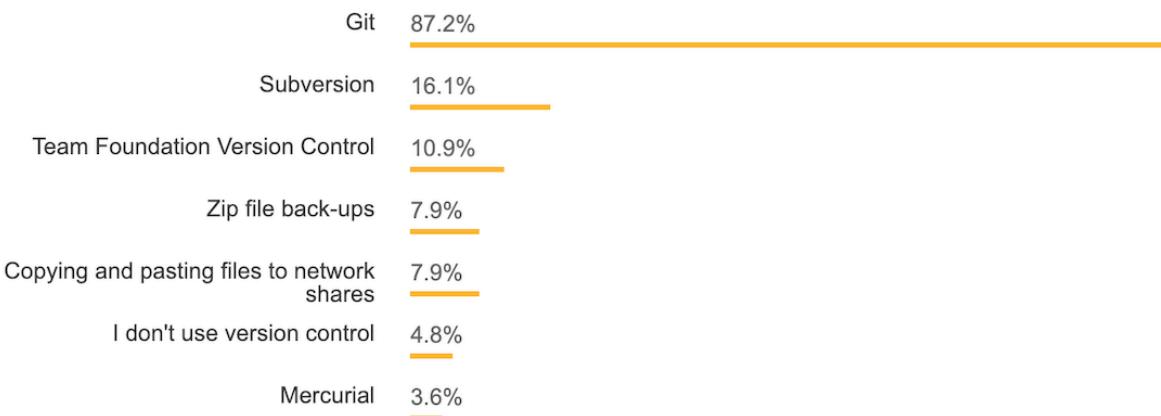
We've included these sections to remind you that you're making a sound choice.

Plain text + Git

Software developers keep track of their code with Git. We learned in the last chapter that most code files get kept in plain text, so Git plays well with plain text.

Everyone else is doing it

Git has become the most common version control system used by [programmers¹²⁴](#).



74,298 responses; select all that apply

Git is the dominant choice for version control for developers today, with almost 90% of developers checking in their code via Git.

source: [StackOverflow Developer Survey Results¹²⁵](#)

¹²³https://en.wikipedia.org/wiki/Version_control

¹²⁴<https://insights.stackoverflow.com/survey/2018#work-version-control>

¹²⁵<https://insights.stackoverflow.com/survey/2018#overview>

Git is a useful way to think about making changes

Git is also a helpful way of thinking about the file changes in our project. The terminology of Git is strange at first, but after we use Git long enough, we'll start thinking about our code in terms of 'adds', 'commits', 'pushes', 'pulls' and 'repos.' We'll go over these terms in-depth in the next chapter.

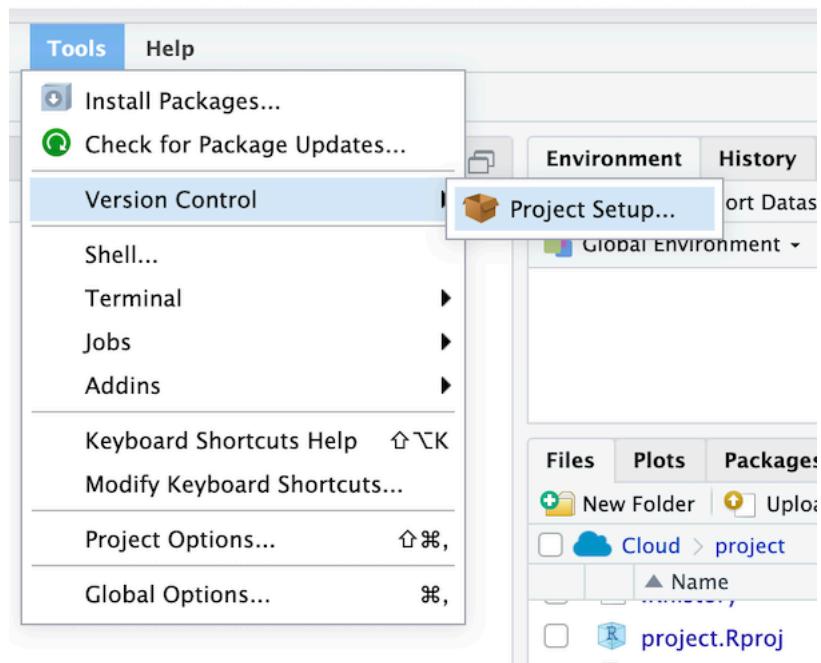
As someone who analyzes data regularly, we can start to think about how to quantify these concepts, too. For example, we can count changes to files, or lines of code, or even measure how changes to code in one file can alter what other files do. All of this is exciting because it means we can start to quantify the changes we make, and we begin to think about our work in exciting new ways.

Setting up Git

If you'd like to install Git on your local machine, you can do so following these two links:

1. Download and install [Git](#).¹²⁶
2. Create a [Github](#)¹²⁷ account.

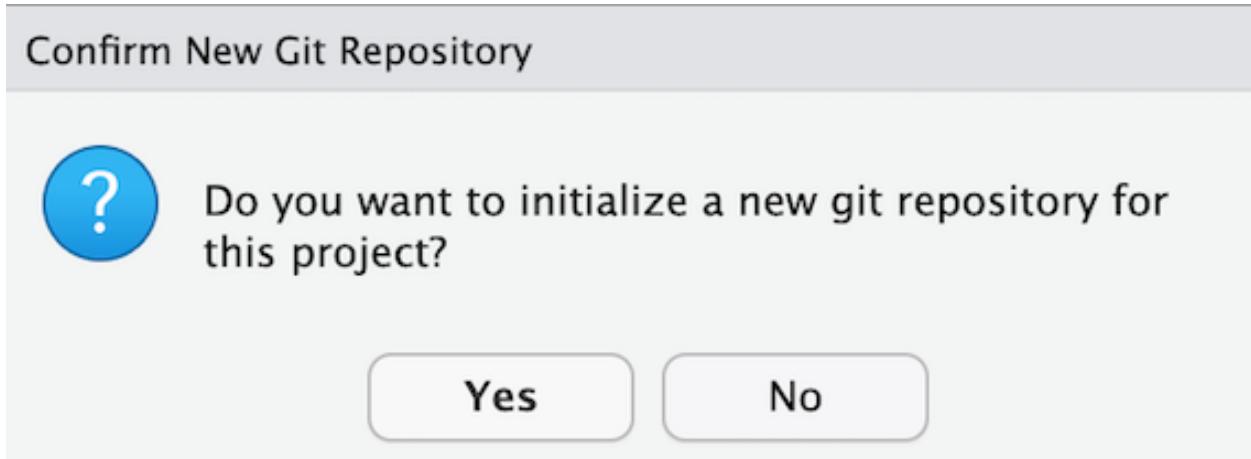
In RStudio.Cloud, we want to add version control to this project from *Tools > Version Control > Project Setup*



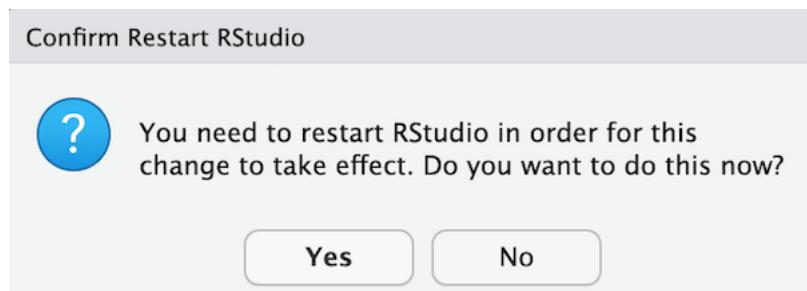
¹²⁶<http://git-scm.org/>

¹²⁷<https://github.com/>

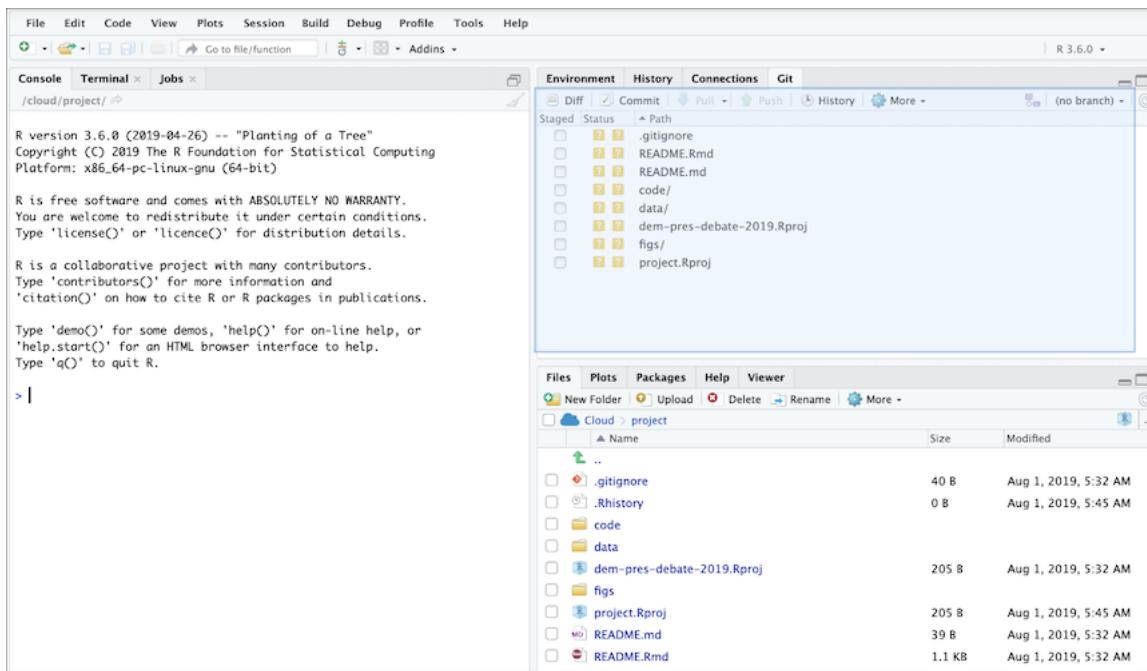
From here, we will see the *Git/SVN* option on the sidebar, where we will select *Git* from the drop-down list next to *Version control system*. After this, RStudio.Cloud will ask if we want to *initialize a new git repo*, which we do.



Then we will be asked if we are OK to restart RStudio.Cloud (and we are).



After restarting the RStudio IDE, we should see the **Git** tab in one of the panes.



Configuring Git

Git needs a little configuration before we can start using it and linking it to Github. There are three levels of configuration within Git, system, user, and project.

1) For **system** level configuration use:

```
git config --system
```

2) For **user** level configuration, use:

```
git config --global
```

3) For **Project** level configuration use:

```
git config
```

We'll set our Git `user.name` and `user.email` with `git config --global` so these are configured for all projects.

```
1 $ git config --global user.name "Martin Frigaard"
2 $ git config --global user.email "mjfrigaard@gmail.com"
```

We can check what we've configured with `git config --list`.

```
1 $ git config --list
```

At the bottom of the output, we can see the changes.

```
1 user.name=Martin Frigaard
2 user.email=mjfrigaard@gmail.com
3 core.repositoryformatversion=0
4 core.filemode=true
5 core.bare=false
6 core.logallrefupdates=true
```

On our local machine, the `user.name` and `user.email` are in the `.gitconfig` file.

We can view this using:

```
1 $ cat .gitconfig
2 [user]
3   name = Martin Frigaard
4   email = mjfrigaard@gmail.com
```

Synchronizing RStudio and Git/Github

Jenny Bryan¹²⁸ has created the online resource Happy Git and GitHub for the useR¹²⁹ has all the information anyone would need for connecting RStudio and Git/Github. We echo a lot of this information below (with copious screenshots).

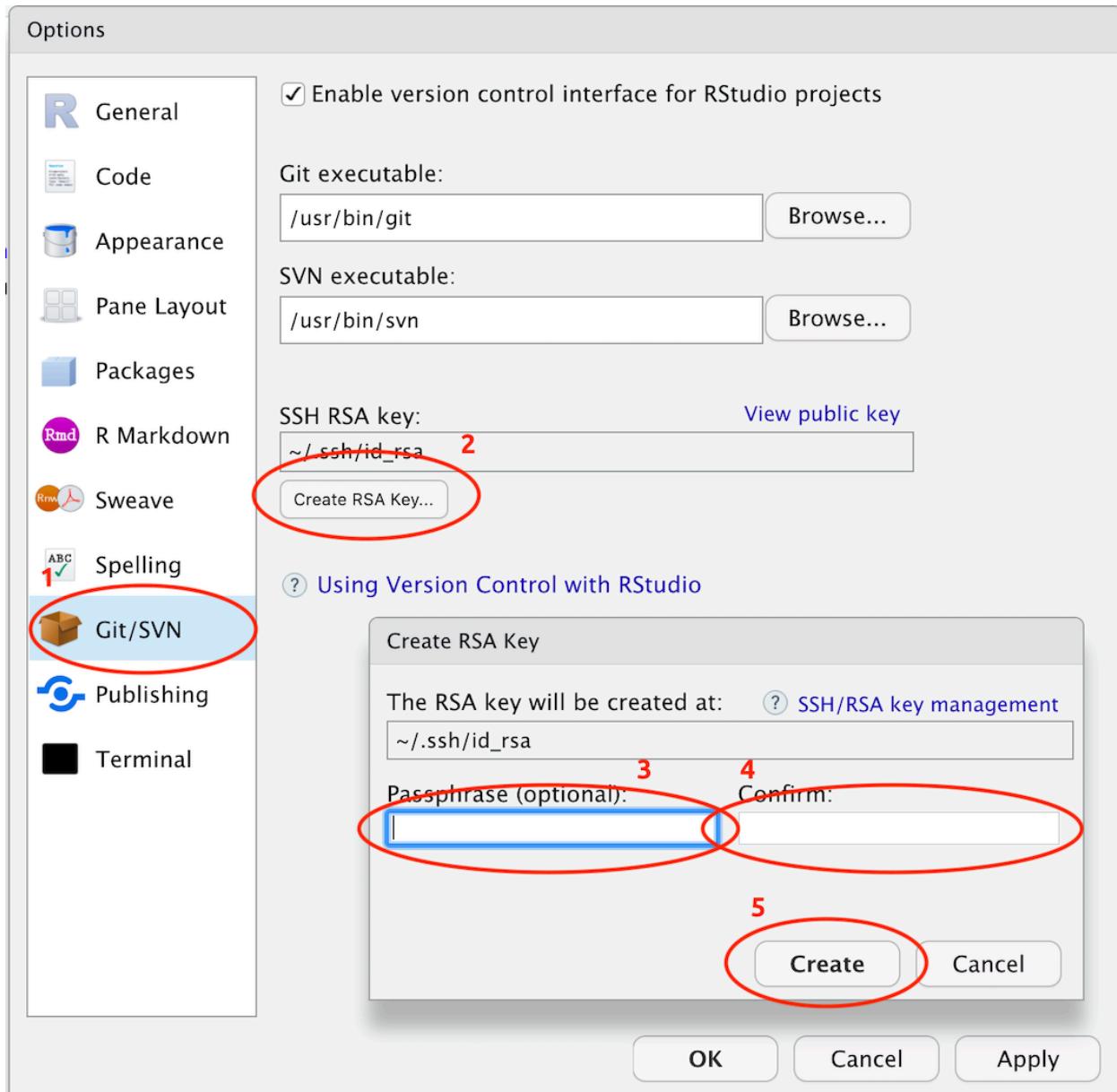
The first step is setting up the RSA Key and passphrase.

Go to *Tools > Global Options > ...*

- 1. Click on *Git/SVN*
- 1. Then *Create RSA Key...*
- 3, 4, and 5. In the dialog box, enter a passphrase (and store it in a safe place), then click *Create*.

¹²⁸<https://jennybryan.org/>

¹²⁹<http://happygitwithr.com/>



The result should look something like this:

Generating public/private rsa key pair.

Your identification has been saved in /home/rstudio-user/.ssh/id_rsa.

Your public key has been saved in /home/rstudio-user/.ssh/id_rsa.pub.

The key fingerprint is:

The key's randomart image is:

-----[RSA 2048]-----

$\text{---} = 0 = \dots 0000$

| Page

8.8.8.8.8.8.8.8.8.

Xoo Fo I

• 800-234-2345

—

—
—

1

—
—

ESU12E67

-----[SHA256]-----

```
2 Generating public/private rsa key pair.
3 Enter file in which to save the key (/Users/username/.ssh/id_rsa):
4 Enter passphrase (empty for no passphrase):
5 Enter same passphrase again:
6 Your identification has been saved in /Users/username/.ssh/id_rsa.
7 Your public key has been saved in /Users/username/.ssh/id_rsa.pub.
8 The key fingerprint is:
9 SHA483:g!bB3r!sHg!bB3r!sHg!bB3r!sH USEFUL-COMMENT
10 The key's randomart image is:
11 +---[RSA 2891]---+
12 |      o+ . .
13 |      .=.o . +
14 |      ..= + +
15 |      .+* E
16 |      .= So =
17 |      . +. = +
18 |      o... = ..* .
19 |      o ++=.o =o.
20 |      ..o.++o.=+.
21 +---[SHA483]---+
```

Great! We need to go back to Terminal and store this SSH from RStudio.

Adding a key SSH in Terminal

In the Terminal, we enter the following commands.

```
1 $ eval "$(ssh-agent -s)"
```

The response tells us we're an Agent.

```
1 Agent pid 007
```

Now we want to add the *SSH RSA* to the keychain. There are three elements in this command: the `ssh-add`, the `-K`, and `~/.ssh/id_rsa`.

- The `ssh-add` is the command to add the *SSH RSA*
- The `-K` stores the passphrase we generated, and
- `~/.ssh/id_rsa` is the location of the *SSH RSA*.

```
1 $ ssh-add -k /home/rstudio-user/.ssh/id_rsa
```

Enter the passphrase, and then git should tell us the identity has been added.

```
1 Enter passphrase for /home/rstudio-user/.ssh/id_rsa:  
2 Identity added: /home/rstudio-user/.ssh/id_rsa (/home/rstudio-user/.ssh/id_rsa)
```

Create the `.ssh/config` file

Most operating systems require a `config` file. We can do this using the Terminal commands above.

First, I move into the `.ssh/` directory.

```
1 $ cd /home/rstudio-user/.ssh
```

Then we create this `config` file with `touch`

```
1 $ touch config  
2 # verify  
3 $ ls  
4 config id_rsa id_rsa.pub
```

We use `echo` to add the following text to the `config` file.

```

1 Host *
2 AddKeysToAgent yes
3 UseKeychain yes

```

Recall the `>>` will send the text to the config file.

```

1 # add text
2 $ echo "Host *
3 > AddKeysToAgent yes
4 > UseKeychain yes" >> config

```

Finally, we can check config with cat

```

1 # verify
2 $ cat config
3 Host *
4 AddKeysToAgent yes
5 UseKeychain yes

```

Great! Now I am all set up to use Git with RStudio. In the next section, we'll extend our Github skills by moving the contents of a local folder to Github.

A quick git terminology overview

Below are some commonly used terms/commands associated with Git and Github.

init - the command `git init` is used to initialize a new git repository (it tells Git to start tracking changes in this directory).

status - `git status` will tell you what you've done and what is happening. You can check the status of a git repository with `git status` (use this liberally).

add - For Git to keep track of the changes we make to files, we have to tell Git which files to pay attention to. We can do this using `git add`. The `git add -A` tells git to stage *ALL* the files that are in an initialized repo.

commits - commits are the staple in Git/Github the workflow. Commits are what Git uses to track the changes you've made to files or folders. Commits are confusing because they can be nouns ("I'm creating a commit with these changes") or verbs (I am going to commit these changes to my project").

To quote David Demaree,

- “Semantically, each commit represents a complete snapshot of the state of your project at a given moment in time; its unique identifier serves to distinguish that state from the way the files in your project looked at any other moment in time.”*

repository - repos are the files and folders in your project and all the changes you make while working on them. On your local computer, a repository can exist in a folder you initialize a repository in (see below). On Github, a repo has the following structure:

`https://github.com/<username>/<repository_name>.`

clone - this command copies all files and changes into a new working directory from a remote, initializes (`init`) a new Git repository, and it adds a remote called `origin`.

diff - this is how Git shows differences between files. Read more about how changes are formatted/displayed [here](#).¹³⁰

More on Git and Github and data organization

Fortunately, many articles have come out in the last few years with excellent, practical advice on organizing data analysis projects. I recommend reading these before getting started (you'd be surprised at the cacophony of files a single project can produce). We've listed a few 'must-reads' below:

- [the importance of using version control](#)¹³¹
- [sharing data with collaborators](#)¹³²
- [how to name your files](#)¹³³

¹³⁰<https://www.git-tower.com/learn/git/ebook/en/command-line/advanced-topics/diffs>

¹³¹<https://www.nature.com/news/democratic-databases-science-on-github-1.20719>

¹³²<https://www.tandfonline.com/doi/full/10.1080/00031305.2017.1375987>

¹³³<https://speakerdeck.com/jennybc/how-to-name-files>

Part 5: RStudio.Cloud

In this chapter, we're going to introduce R packages (the source of R's power), how to quickly navigate RStudio's panes, and more benefits of working in Rmarkdown/plain text.

What have we done so far?

The last few chapters covered a lot, so let's do a quick recap:

- 1) *We created RStudio.Cloud and Github accounts*
- 2) *We downloaded a repository from Github onto our local computer.*
- 3) *We uploaded the repository files into RStudio.Cloud, then moved the files into the project folder using the Terminal.*
- 4) *We added and committed these changes to our local Git repository using RStudio's Git integration.*

Now we have a project folder in RStudio with the following structure:

```
1 project/
2   └── CHANGELOG.txt
3   └── README.Rmd
4   └── README.md
5   └── code/
6   └── data/
7   └── dem-pres-debate-2019.Rproj
8   └── figs/
9   └── project.Rproj
```

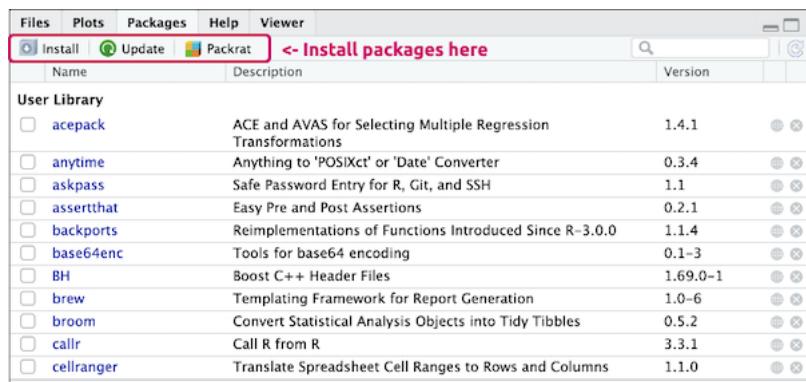
There are three folders in this project: `code/`, `data/`, and `figs/`. We will go over the contents of these folders and files in more depth in the following sections. First, we will cover some tips for navigating RStudio's various panes.

We covered the **Terminal** and **Git** panes in chapters three and four. Now that we've run some code, we will go over RStudio's other panes.

Packages pane

Packages are a vital part of the R ecosystem. R packages are collections of functions and objects collected together with a specific purpose. When we started our RStudio.Cloud session, a few packages get loaded automatically. R users typically refer to these default packages and functions as “base R.” Base R packages cover a wide range of statistical procedures and visualizations. Unfortunately, base R can also be challenging to learn because of its inconsistent syntax and style conventions.

A list of the packages and their descriptions are available in the Packages pane. We’ve installed the packages we will use with the `00.1-inst-packages.R` file. If we wanted to, we could also install more packages using *Install* icon, and then enter its name. We don’t recommend this, though, because you’ll want a record of the packages you used in the project.



The screenshot shows the RStudio Cloud interface with the 'Packages' tab selected. At the top, there are three buttons: 'Install' (highlighted in red), 'Update', and 'Packrat'. Below them, a search bar contains the placeholder text '<- Install packages here'. The main area is a table with columns for 'Name', 'Description', and 'Version'. The table lists various R packages from the 'User Library', each with a checkbox to the left of its name. The packages listed are: acepack, anytime, askpass, assertthat, backports, base64enc, BH, brew, broom, callr, and cellranger. Each package has a brief description and its version number.

Name	Description	Version
User Library		
<input type="checkbox"/> acepack	ACE and AVAS for Selecting Multiple Regression Transformations	1.4.1
<input type="checkbox"/> anytime	Anything to 'POSIXct' or 'Date' Converter	0.3.4
<input type="checkbox"/> askpass	Safe Password Entry for R, Git, and SSH	1.1
<input type="checkbox"/> assertthat	Easy Pre and Post Assertions	0.2.1
<input type="checkbox"/> backports	Reimplementations of Functions Introduced Since R-3.0.0	1.1.4
<input type="checkbox"/> base64enc	Tools for base64 encoding	0.1-3
<input type="checkbox"/> BH	Boost C++ Header Files	1.69.0-1
<input type="checkbox"/> brew	Templating Framework for Report Generation	1.0-6
<input type="checkbox"/> broom	Convert Statistical Analysis Objects into Tidy Tibbles	0.5.2
<input type="checkbox"/> callr	Call R from R	3.3.1
<input type="checkbox"/> cellranger	Translate Spreadsheet Cell Ranges to Rows and Columns	1.1.0

You can click on the names of each package to learn more about them and load them into the RStudio session. One of the great things about R is the many user-created packages that greatly expand the number of functions. At the time of this writing, R users have contributed [14638¹³⁴](#) packages available for us to download and use.

tidyverse packages

R is open-source software, so users can write packages to expand and enhance its functionality. The [tidyverse¹³⁵](#) is a collection of R packages lead by RStudio’s Chief Scientist, [Hadley Wickham¹³⁶](#). These package have been developed for doing data science. All tidyverse packages share a few similar underlying principles that allow them to work well together.

Unlike base R, the tidyverse also a consistent grammar and syntax, which makes it easier to read and write. You can learn more about this syntax in the [R for Data Science¹³⁷](#) text or on the [tidyverse webpage¹³⁸](#).

¹³⁴<https://cran.r-project.org/web/packages/>

¹³⁵<https://cran.r-project.org/web/packages/tidyverse/vignettes manifesto.html>

¹³⁶<http://hadley.nz/>

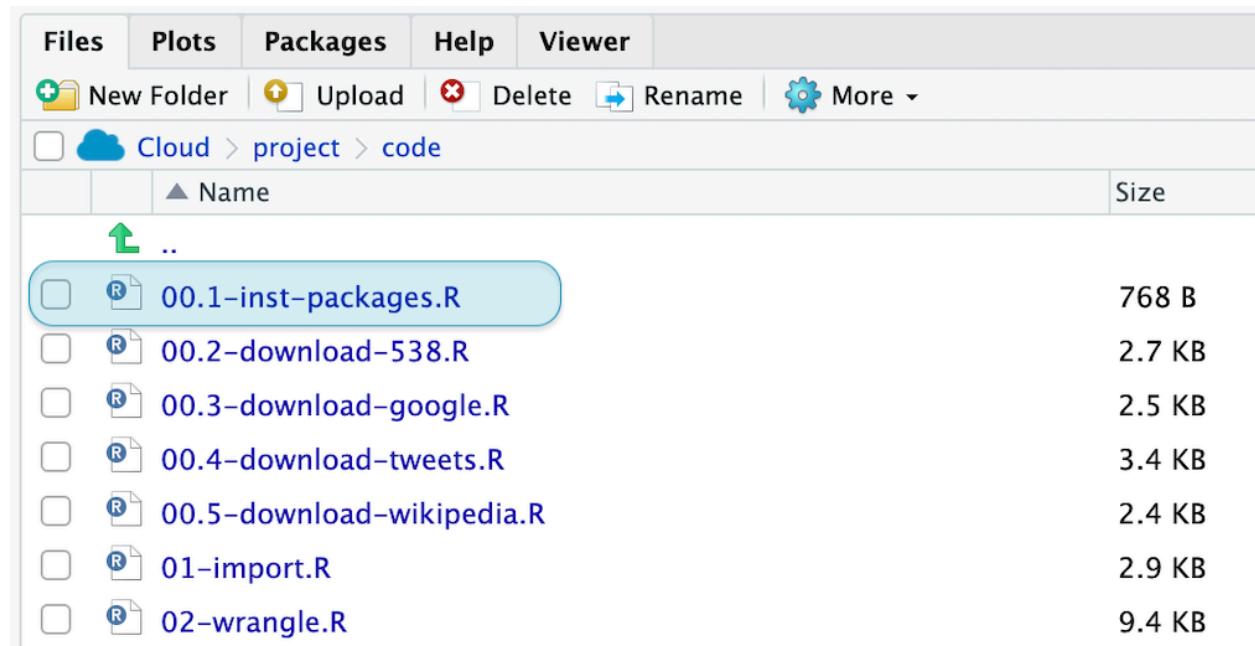
¹³⁷<https://r4ds.had.co.nz/>

¹³⁸<https://www.tidyverse.org/>

So far, the code we've run comes from base R. Going forward; we're going to use various packages from the `tidyverse`.

Files pane

The **Files** pane displays the files and folders in this project. The file path is visible in the top portion of the pane, beneath the options for *New Folder*, *Upload*, *Delete*, *Rename*, and *More*. An image of this folder's contents is below:



The screenshot shows the RStudio.Cloud interface with the 'Files' tab selected. The file path is shown as 'Cloud > project > code'. The list of files includes:

	Name	Size
<input type="checkbox"/>	..	
<input type="checkbox"/>	00.1-inst-packages.R	768 B
<input type="checkbox"/>	00.2-download-538.R	2.7 KB
<input type="checkbox"/>	00.3-download-google.R	2.5 KB
<input type="checkbox"/>	00.4-download-tweets.R	3.4 KB
<input type="checkbox"/>	00.5-download-wikipedia.R	2.4 KB
<input type="checkbox"/>	01-import.R	2.9 KB
<input type="checkbox"/>	02-wrangle.R	9.4 KB

The script to install the packages for this project (`00.1-inst-packages.R`) is in the code folder. Navigate to this folder either by clicking on the file path (Cloud/project/code).

We can see `00.1-inst-packages.R` and other scripts are in the **Files** pane. We now know that because the `.R` files are plain text files, so they have code for the computer to execute, and comments for a human to tread.

In R, we can create comments with a preceding `#` on any line.

```
1 # this is a comment
```

When we click on the `00.1-inst-packages.R` file, we can see it open in the **Source** pane. RStudio.Cloud is giving us a few hints about the packages referenced in the `.R` file. We're being told we need to install a package before moving forward:

Package devtools required but is not installed.

RStudio gives us a choice to **Install** or * **Don't Show Again**, and we'll click on the ***Install** option. Packages vary in the length they take to install, but we'll wait patiently for devtools package finish downloading. After the install has finished, you should see the > prompt in the **Console** pane.

Source pane

When working in RStudio, sometimes we'll write some code, and want to run that code and see if it works. To do this, we can highlight the rest of the script with the mouse cursor (lines 14-22), then hold down **ctrl** or **cmd** and hit **enter** or **return**.

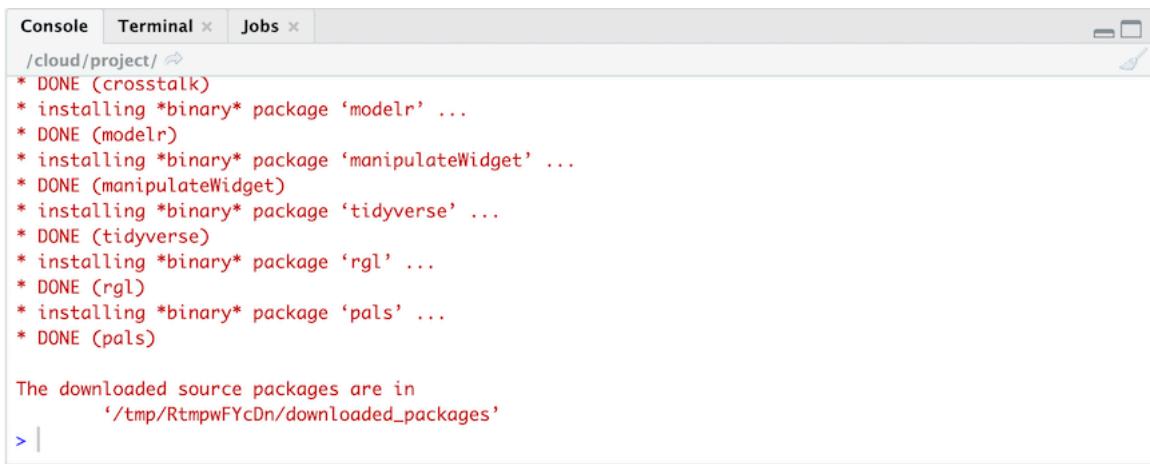
```

1 #=====
2 # This is code to create: 00-packages.R
3 # Authored by and feedback to mjfrigaard@gmail.com
4 # MIT License
5 # Version: 1.1
6 #
7
8 # first install devtools
9 install.packages("devtools")
10
11 # now get the gtrends package
12 devtools::install_github('PMassicotte/gtrendsR')
13
14 # install the other packages for mapping
15 install.packages(c("gtrendsR", "maps",
16                   "lettercase", "viridis",
17                   "pals", "tidyverse",
18                   "ggrepel", "scico",
19                   "ggthemes", "skimr",
20                   "Hmisc", "janitor",
21                   "DT", "inspectDF",
22                   "rtweet", "visdat"))
23
24
25

```



After the installation has completed for all the packages, we'll see the following in the **Console** pane.



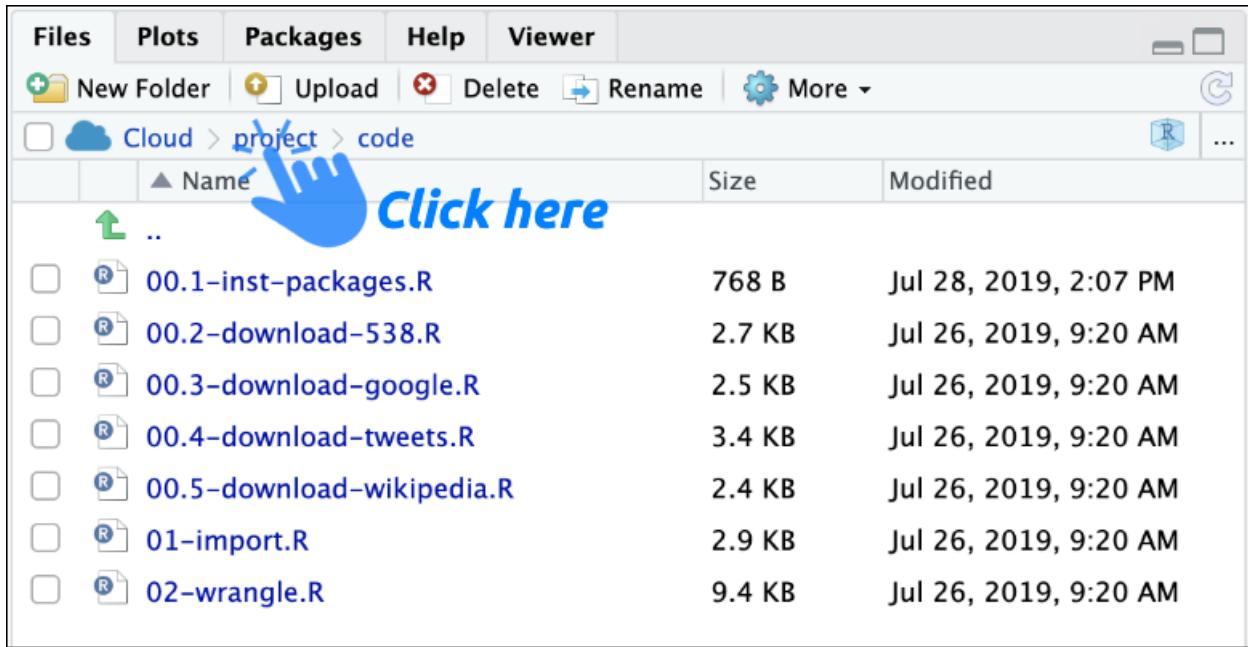
```

Console Terminal x Jobs x
/cloud/project/ ↗
* DONE (crosstalk)
* installing *binary* package 'modelr' ...
* DONE (modelr)
* installing *binary* package 'manipulateWidget' ...
* DONE (manipulateWidget)
* installing *binary* package 'tidyverse' ...
* DONE (tidyverse)
* installing *binary* package 'rgl' ...
* DONE (rgl)
* installing *binary* package 'pals' ...
* DONE (pals)

The downloaded source packages are in
  '/tmp/RtmpwFYcDn/downloaded_packages'
> 

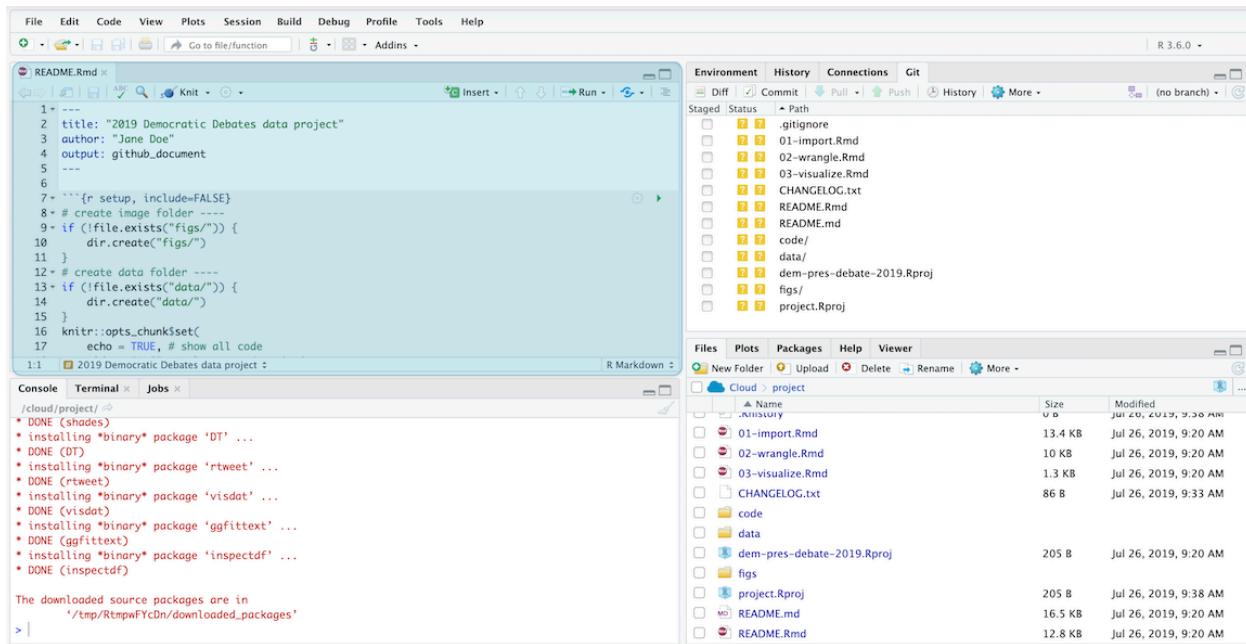
```

Back in the **Files** pane, we can navigate back to the project/ folder by clicking on the path above the files.



	Name	Size	Modified
	..		
<input type="checkbox"/>	00.1-inst-packages.R	768 B	Jul 28, 2019, 2:07 PM
<input type="checkbox"/>	00.2-download-538.R	2.7 KB	Jul 26, 2019, 9:20 AM
<input type="checkbox"/>	00.3-download-google.R	2.5 KB	Jul 26, 2019, 9:20 AM
<input type="checkbox"/>	00.4-download-tweets.R	3.4 KB	Jul 26, 2019, 9:20 AM
<input type="checkbox"/>	00.5-download-wikipedia.R	2.4 KB	Jul 26, 2019, 9:20 AM
<input type="checkbox"/>	01-import.R	2.9 KB	Jul 26, 2019, 9:20 AM
<input type="checkbox"/>	02-wrangle.R	9.4 KB	Jul 26, 2019, 9:20 AM

We can see the `README.Rmd` file in the **Files** pane in the lower right corner. We'll click on it, so the file opens in the **Source** pane (see image below).

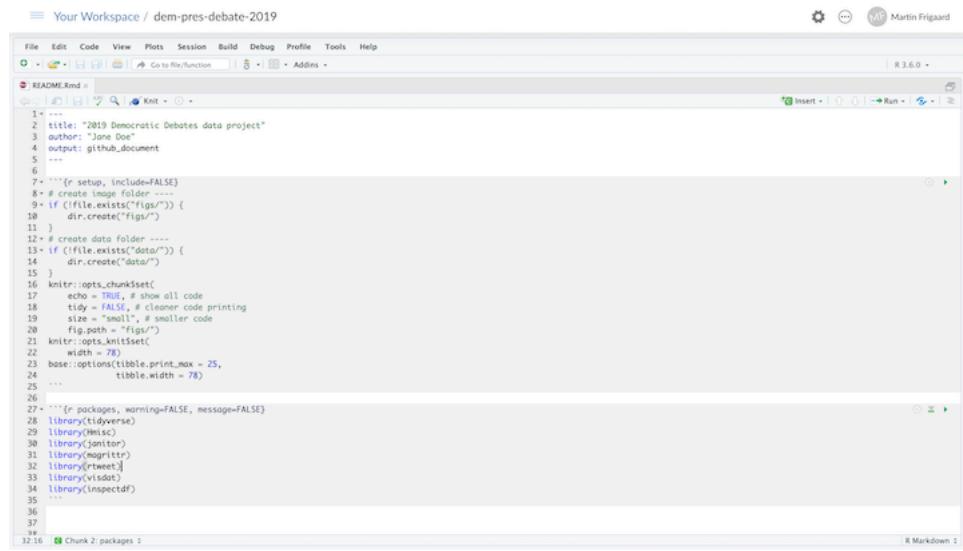


We've highlighted the pane of RStudio that's displaying the `README.Rmd` file. If you see another prompt to install additional packages, click on *Install* and wait for the installation to finish. As you can see, this is a relatively small area in the IDE. Working in a tiny corner of the screen can be hard, but fortunately, RStudio gives us the ability to expand any pane into a fullscreen view. In this case, we'll zoom in the `README.Rmd` file.

If we want to focus on the **Source** pane, and can zoom in using `shift+control+1`.



As soon as you click 1, your screen should expand to look like the image below:



The screenshot shows the RStudio Cloud interface with the Source pane active. The code in the pane is as follows:

```

1: ---
2: title: "2019 Democratic Debates data project"
3: author: "Jane Doe"
4: output: github_document
5: ---
6:
7: ```{r setup, include=FALSE}
8: # create image folder ----
9: if (!file.exists("figs/")) {
10:   dir.create("figs/")
11: }
12: # create data folder ----
13: if (!file.exists("data/")) {
14:   dir.create("data/")
15: }
16: knitr::opts_chunk$set(
17:   echo = TRUE, # show all code
18:   tidy = FALSE, # cleaner code printing
19:   size = "small", # smaller code
20:   fig.path = "figs/"
21:   knitr::opts_chunk$set(
22:     width = 70)
23: base::options(tibble.print_max = 25,
24:               tibble.width = 78)
25: ...
26: ...
27: ```{r packages, warning=FALSE, message=FALSE}
28: library(tidyverse)
29: library(Hmisc)
30: library(janitor)
31: library(magrittr)
32: library(sweave)
33: library(cldst)
34: library(inspectdf)
35: ...
36: ...
37: ...
38: ...
39: ...
40: ...
41: ...
42: ...
43: ...
44: ...
45: ...
46: ...
47: ...
48: ...
49: ...
50: ...
51: ...
52: ...
53: ...
54: ...
55: ...
56: ...
57: ...
58: ...
59: ...
60: ...
61: ...
62: ...
63: ...
64: ...
65: ...
66: ...
67: ...
68: ...
69: ...
70: ...
71: ...
72: ...
73: ...
74: ...
75: ...
76: ...
77: ...
78: ...
79: ...
80: ...
81: ...
82: ...
83: ...
84: ...
85: ...
86: ...
87: ...
88: ...
89: ...
90: ...
91: ...
92: ...
93: ...
94: ...
95: ...
96: ...
97: ...
98: ...
99: ...
```

```

Since the **Source** pane is where we'll write most of our code, it's also the pane we are spending most of our time. We've seen that's where our files open, so being able to focus on the area quickly is helpful.

## Console

The **Console** is probably the second most used area in RStudio (it displays most of the output), so we also should know how to zoom in on this pane. To focus on this pane, we'll use **shift+control+2**.

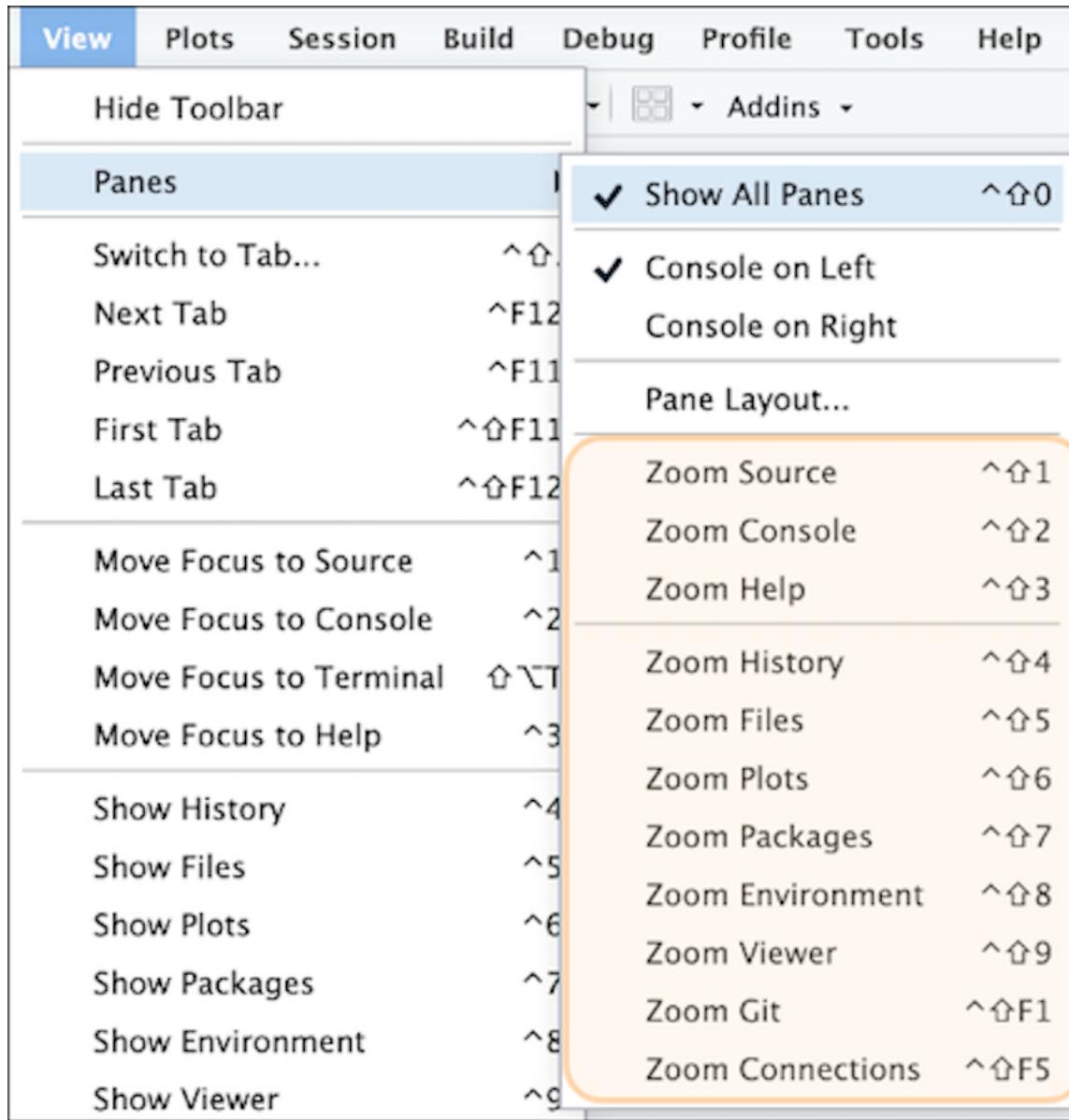


What if we accidentally zoom into the wrong pane? We can easily resize the IDE to its original position by holding down the **shift+control** buttons, then clicking either number again. Combining those keys should return the IDE to its default arrangement.

## Working in one pane

We can move through every pane in RStudio by holding down shift+control and clicking on numbers 1-9. Go ahead and do that now.

If you ever forget which number corresponds to which pane, you can always find them under *View > Panes* (see image below)

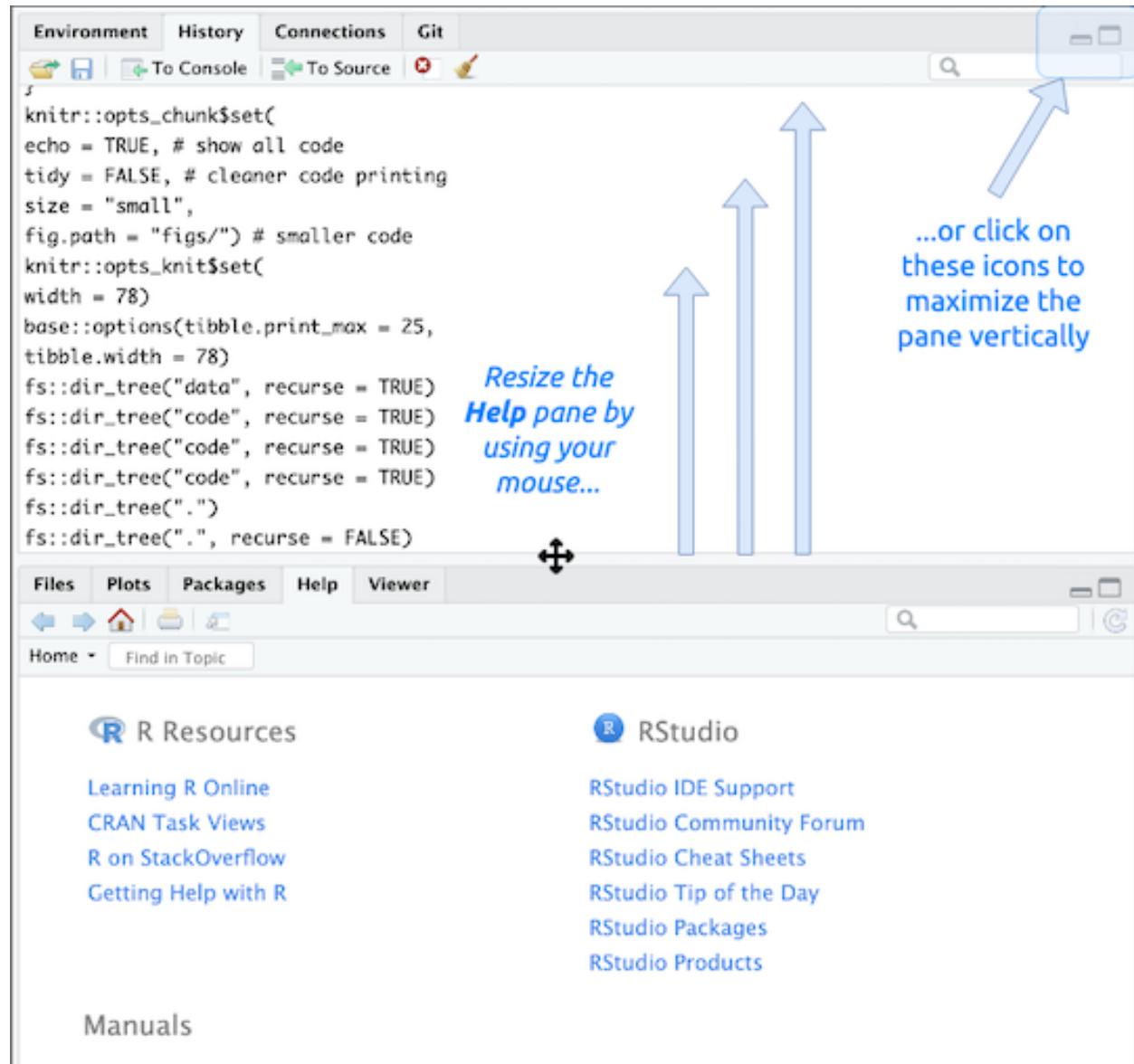


## Help

Inevitably, we'll write something that doesn't work. When things aren't working (try to remain calm), you're going to want a place to start looking for solutions. The **Help** pane is accessible in the

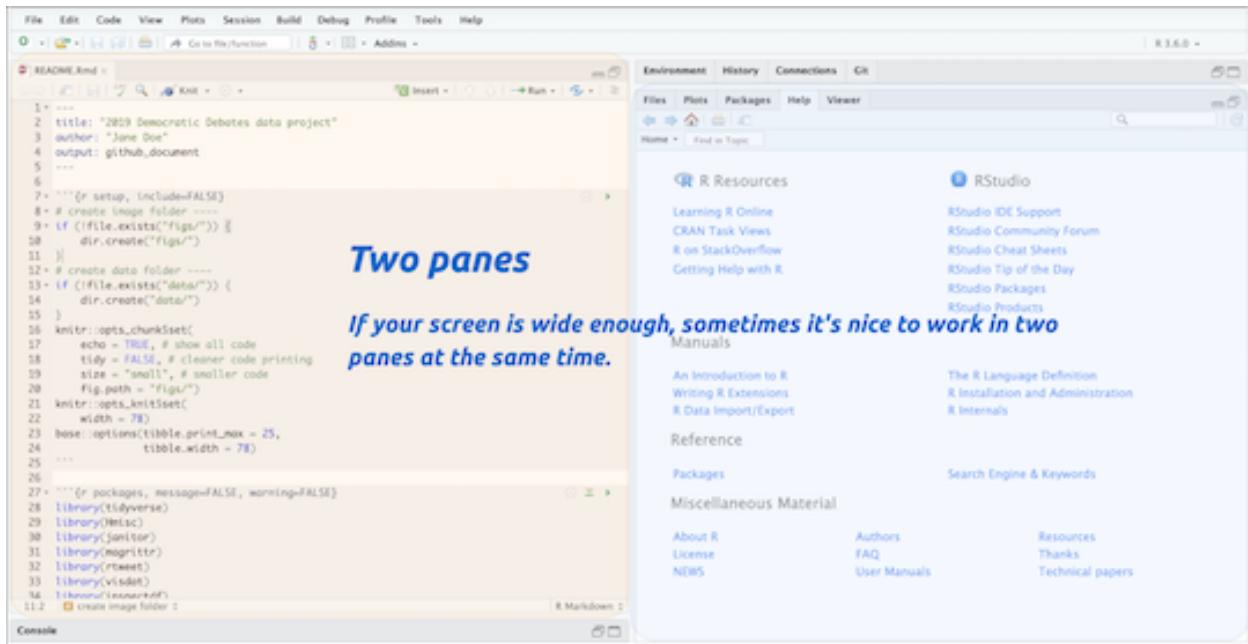
lower right corner of the IDE.

If it's hard to read the **Help** pane, you can adjust the size by following the instructions in the image below:



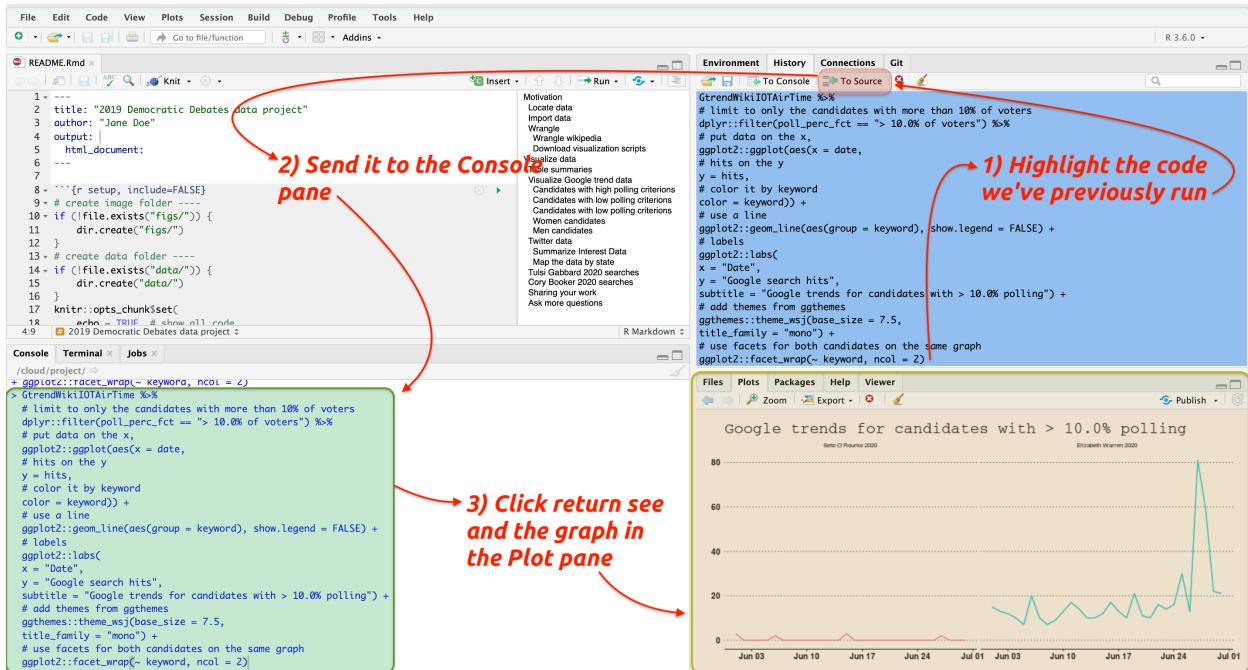
## Working in two panes

Sometimes we'll want to do our work in more than one pane at a time. For example, what if we're working in the **Source** pane, but have a question about a function? We can get answers to a lot of problems using RStudio's internal **Help** pane. After resizing the **Help** pane, we should see the following layout in our IDE.



## History and Plots

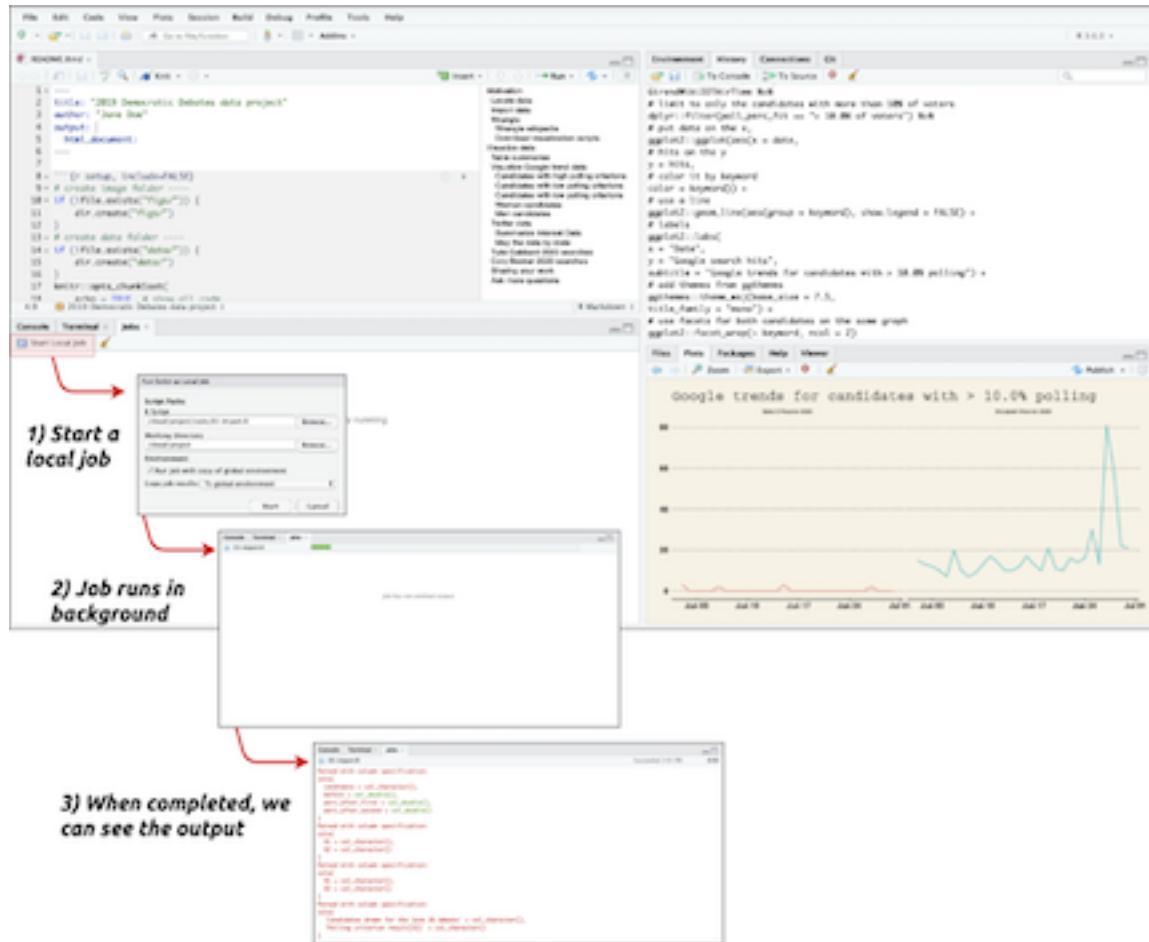
The History panel has the code we've previously run and can be sent directly to the Console or Source panes.



The History comes in handy when we're looking to reproduce a plot without having to run an entire script.

## Jobs and Connections

The **Jobs** pane allows us to run multiple jobs at the same time. For example, if we wanted to run the `01-import.R` and then `02-wrangle.R`, we could do this using the **Jobs** pane. See the image below:



The **Connections** pane tells us what databases we're connected to, and their connection status. This pane require us to cover some advanced topics, but nothing are too complicated if you've been using RStudio and working with data for awhile. Check out more information [here](#)<sup>139</sup>.

## Recap on RStudio panes

As we pointed out earlier, RStudio is like a workbench built around different panes. Each pane serves a specific purpose and being able to move between them allow us to work quickly and efficiently.

**Keyboard shortcuts are time-savers.** Not having to switch from keyboard to mouse over and over again keeps us focused on the task at hand.

<sup>139</sup><https://db.rstudio.com/rstudio/connections/>

## Where should we write code?

You might be wondering why we have provided you with both .R scripts and .Rmd files. Well, we want to show you a few options for documenting your work in RStudio. The two sections below outline two standard options (but these are by no means the only way to work with these tools!)

### Option 1) everything goes in scripts

rmarkdown is a relatively new package, so early R users didn't have an option to put everything in .Rmd files (we fit into this group). We have found this approach is excellent when you're reasonably confident about the project. If you're familiar with the data files, wrangling techniques, visualizations, and products, you can quickly outline the scripts and predetermine their names and contents.

For example, the folder tree of our code folder tells us what we should expect from this project.

```
1 code
2 └── 00.1-inst-packages.R
3 └── 00.2-download-538.R
4 └── 00.3-download-google.R
5 └── 00.4-download-tweets.R
6 └── 00.5-download-wikipedia.R
7 └── 01-import.R
8 └── 02-wrangle.R
```

We can see there are seven script files, each one with a numerical prefix. These prefixes tell us the order to run the code files, too. The file names also tell us what each script does (`download`, `import`, `wrangle`, etc.).

The layout above is an example of how to organize a set of .R scripts that follow the guidelines mentioned in previous chapters. These naming conventions make it easier for newcomers to the project to get oriented to its contents.

### Option 2) everything goes in the Rmarkdown

The second option is well illustrated in the [tweet<sup>140</sup>](#) below by Andrew MacDonald:

---

<sup>140</sup><https://twitter.com/polesasunder/status/953624238266646529?s=20>



**Andr(é) MacDonald**  
@polesasunder



ok now listen, the harsh truth is you're better off writing one thick, messy .Rmd where you keep all your garbage models and weird musings then going off on some precious folder structure and artfully-named .R files where you can't find a damned thing ever. [#rstats #oldman](#)

5:45 AM · Jan 17, 2018 · [TweetDeck](#)

In our experience, this accurately captures the reality of data projects (especially in the beginning stages). Throwing everything into the Rmarkdown file gives us a lot more flexibility by allowing us to add multiple types of content.

All that being said, it's always good practice to come back and revise the README.Rmd document, so its contents outline all steps in the analysis thoroughly. We consider documentation to be a form of communication with our future selves, and this typically involves clearly describing each step. During the revision, we might also break down each of the chunks into individual scripts. We can also add more details to the text portions, with links to external content, etc.

We've found each project usually starts at a bit of a sprint, so we try to capture as much code and content in .Rmd files early. We can focus on trimming it down later, but the flexibility of markdown combined with functional code is better than limiting ourselves to code and comments .R script files.

---

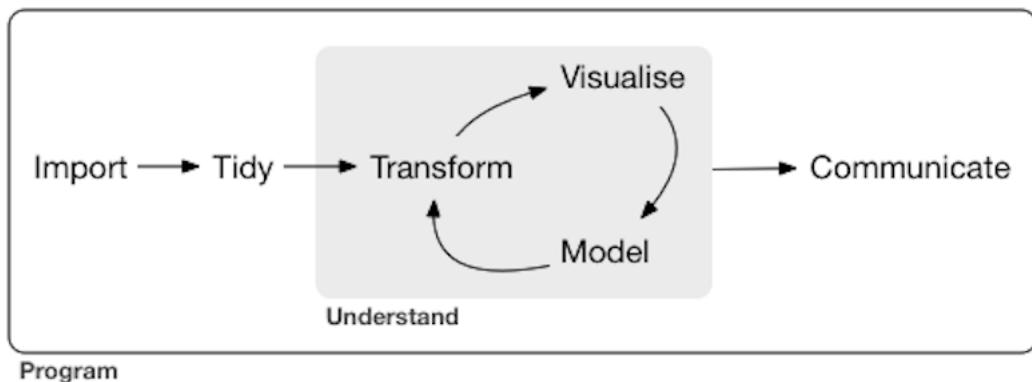
## Rmarkdown

We're going to move forward assuming we're documenting everything in a README.Rmd file. In the next few sections, we'll add various components to the README.Rmd we've placed in the project folder.

If you're ever wondering how to outline your README.Rmd file, the figure from [R for Data Science](#)<sup>141</sup> isn't a bad place to start.

---

<sup>141</sup><https://r4ds.had.co.nz/introduction.html>



## Rmarkdown step 1: create a YAML header

At the top of the README.Rmd document, the first thing we see is what's called the "YAML header", and it's going to tell RStudio.Cloud the file's title, the author, and what the output file will be.

```

1 ---
2 title: "2019 Democratic Debates data project"
3 author: "Jane Doe"
4 output: github_document
5 ---

```

The YAML header always goes at the top of the README.Rmd file, between two sets of three dashes:

```

1 ---
2 title: "2019 Democratic Debates data project"
3 author: "Jane Doe"
4 output: github_document
5 ---

```

People typically use YAML in configuration files, which makes it perfect for setting some default options in our README.Rmd document. Read more about YAML headers in the [RMarkdown book<sup>142</sup>](#) and on the [YAML website<sup>143</sup>](#). YAML actually stands for "YAML Ain't Markup Language<sup>144</sup>".

We can change the output argument to `html_document`<sup>145</sup>, `word_document`<sup>146</sup>, or `pdf_document`<sup>147</sup> and create a different file from the plain text we are going to be working in. For now, we are going to

<sup>142</sup><https://bookdown.org/yihui/rmarkdown/>

<sup>143</sup><https://yaml.org/>

<sup>144</sup><https://en.wikipedia.org/wiki/YAML>

<sup>145</sup><https://bookdown.org/yihui/rmarkdown/html-document.html>

<sup>146</sup><https://bookdown.org/yihui/rmarkdown/word-document.html>

<sup>147</sup><https://bookdown.org/yihui/rmarkdown/pdf-document.html>

focus on the `github_document`<sup>148</sup> output.

## Rmarkdown step 2: know the Knit output options

Another benefit of working in the `README.Rmd` document is that when we combine it with the powerful `knitr` package<sup>149</sup>, we drastically extend the kind of files we can produce from our analysis. `Knitr` follows a principle of `literate programming` put forth by Donald E. Knuth<sup>150</sup>.

“Let us change our traditional attitude to the construction of programs: Instead of imagining that our main task is to instruct a computer what to do, let us concentrate rather on explaining to humans what we want the computer to do.” - Donald E. Knuth, *Literate Programming*, 1984

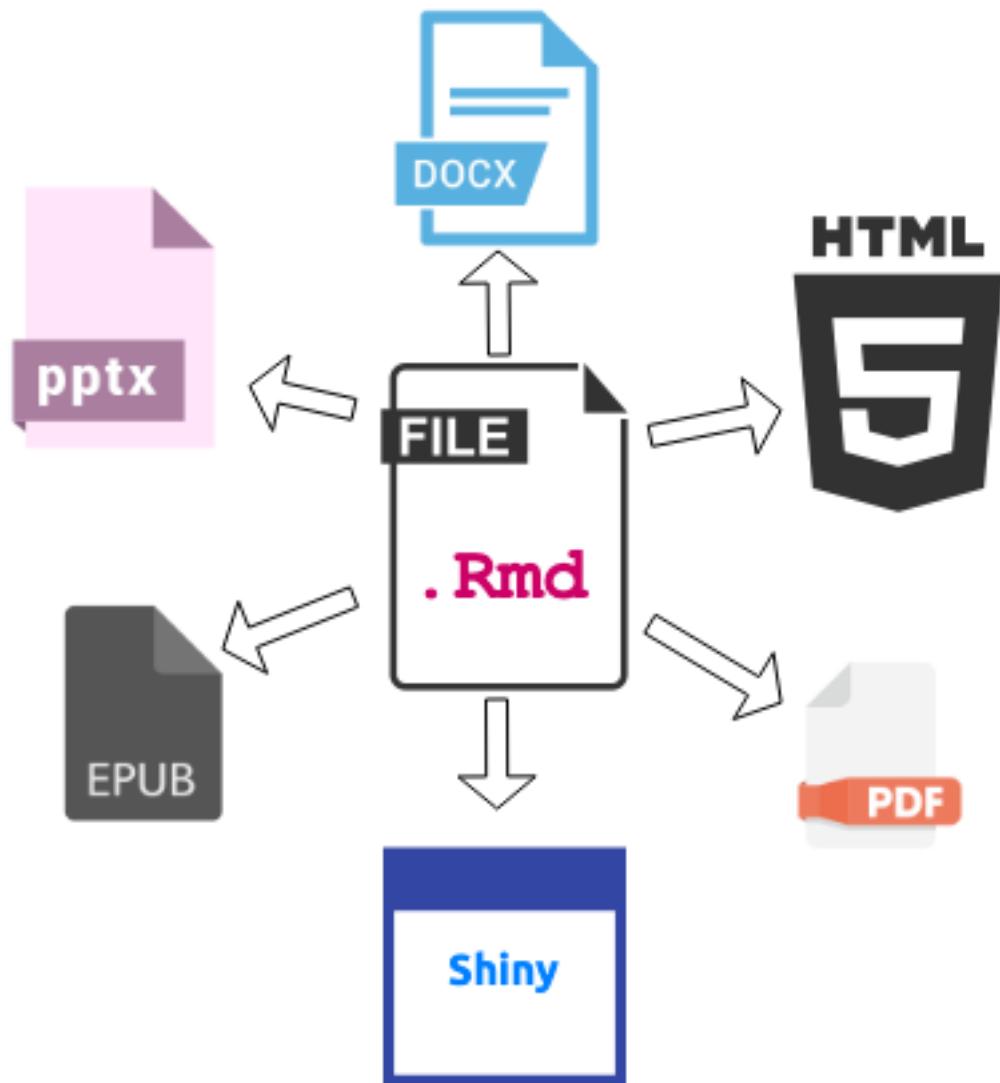
The **Knit** button will render the plain text `.Rmd` document into a variety of different output options.

---

<sup>148</sup>[https://rmarkdown.rstudio.com/github\\_document\\_format.html](https://rmarkdown.rstudio.com/github_document_format.html)

<sup>149</sup><https://yihui.name/knitr/>

<sup>150</sup><https://www-cs-faculty.stanford.edu/~knuth/lp.html>



The .Rmd files also give us the ability to document our intentions, write and execute code, and interpret and explain the results. After we've outlined (and revised) our thought process, we can go about organizing the code in more efficient ways to carry out our intentions.

### Rmarkdown step 3: Compose functional code chunks

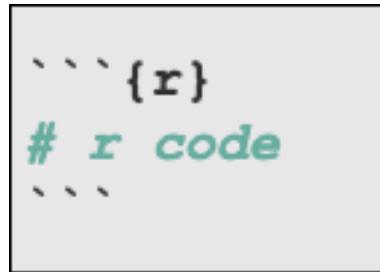
We like to think of Rmarkdown documents like a stack of two sheets of paper, and each piece representing a different file type (.R and .md). When you lay the markdown file on top of the .R script, you get the Rmarkdown file.



R code chunk

Combining markdown R script files gives us syntax for formatting our text (stuff like *italic*, **bold**, code, etc.), and a functional coding script we can get access to by inserting code chunks.

For example, in the README.Rmd file we can type directly onto the paper using the markdown syntax. But if we want to run some R code, imagine tearing a little hole in the markdown paper, and revealing an R script underneath.



R code chunk

These code chunks allow us to run R code in-between the markdown text.

## R Code chunk options & labels

Code chunks come with a long list of options (feel free to experiment with all possible combinations found [here<sup>151</sup>](#) and [here<sup>152</sup>](#)). The most common are echo, eval, and include.

- echo = *show the code in the output?*
- eval = *run the code chunk?*
- include = *put the results from the code in the output?*

We will show a few examples of these with labels below:

<sup>151</sup><https://yihui.name/knitr/options/>

<sup>152</sup><https://rmarkdown.rstudio.com/lesson-3.html>

```
```{r run-show-nothing, include=FALSE}  
# runs and shows nothing  
```
```

```
```{r run-show-code, eval=TRUE, echo=TRUE}  
# runs code and shows results  
```
```

```
```{r run-hide-code, eval=TRUE, echo=FALSE}  
# runs code, shows results, but hides code  
```
```

```
```{r no-run-show-code, eval=FALSE, echo=TRUE}  
# shows code, but doesn't run code  
```
```

The code chunks above are labeled, and we recommend *labeling all code chunks*. We've found it forces us to think through the analysis as a series of operations, and reduces random calculations spread throughout the document.

## The setup code chunk

After the YAML header, our setup chunk tells us what we're going to be doing with the code, text, figures, and output in this README.Rmd file.

Our setup chunk does the following:

**Project folders:** We check to see if the project has image or data folders, and if not, it creates them.

```

1 # create image folder ----
2 if (!file.exists("figs/")) {
3 dir.create("figs/")
4 }
5 # create data folder ----
6 if (!file.exists("data/")) {
7 dir.create("data/")
8 }
```

**Chunk options:** We've included `echo=TRUE`, which means the code will print in the output file. The `tidy=FALSE` makes sure the code doesn't get reformatted when the document gets rendered. We set both options to their default values, but we've included them below so that you can see more examples.

The `size = "small"` and `fig.path = "figs/"` are used to change the size of the printed code, and the location of any output visualizations.

```

1 knitr::opts_chunk$set(
2 echo = TRUE, # show all code
3 tidy = FALSE, # cleaner code printing
4 size = "small", # smaller code
5 fig.path = "figs/") # where the figures will end up
```

**Knit options:** The `knitr::opts_knit$set()` function gives us the ability to stipulate options for what happens when we render the document (knitted and rendered are somewhat synonymous).

```

1 knitr::opts_knit$set(
2 width = 78)
```

**Base options:** the final setting tells RStudio.Cloud how we want the tables to print ("25' rows, and 78 columns).

```

1 base::options(tibble.print_max = 25,
2 tibble.width = 78)
```

Why 78? The standard code file is 80 columns across, which is the length of a punch card (read more [here<sup>153</sup>](#)). We go two columns less than 80 ( to give a little wiggle room).

---

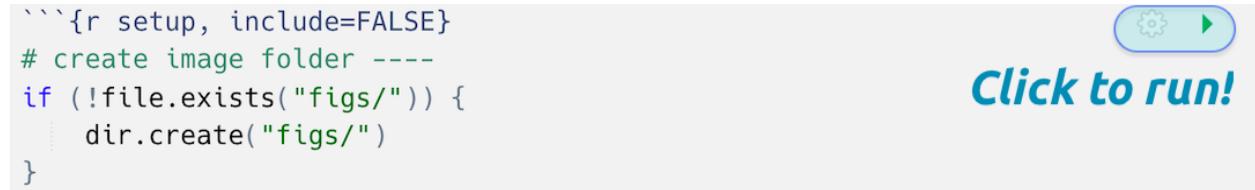
<sup>153</sup><http://bit.ly/code-80-cols>

## Running code chunks

To run the code in your document, we have a few options. First, we can use the same keyboard shortcuts we use for executing code in the .R scripts (ctrl+enter or cmd+return).

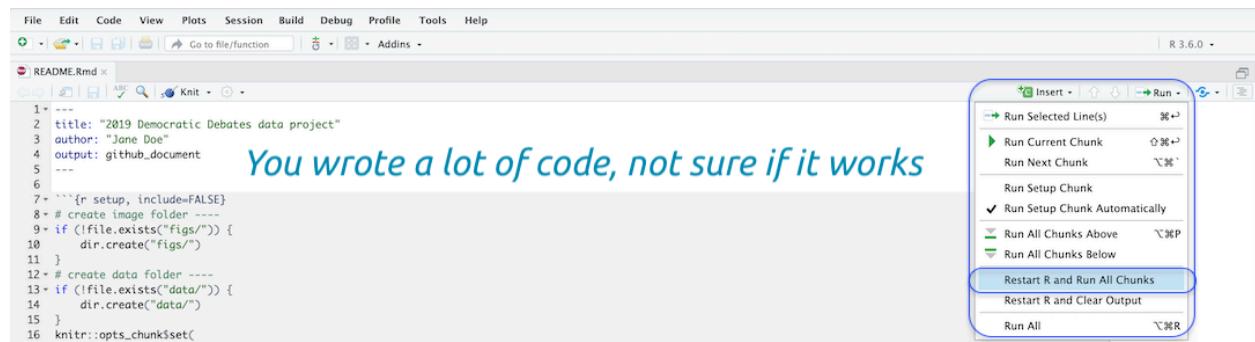
The second option is the small green play icon on the far right-hand side of the code chunk. Running code one chunk at a time is a great way to test what the output of a graph looks like, or to see the data table from a join.

```
```{r setup, include=FALSE}
# create image folder ----
if (!file.exists("figs/")) {
  dir.create("figs/")
}
```



Click to run!

The third option is we run a series of chunks (or the entire document) using the drop-downs at the top of the .Rmd file. We use this option if we've added a few chunks, and we want to make sure everything works together.



Don't worry if all of this is confusing! We will be creating our very own README.Rmd for this project, so you'll get some practice!

Documenting our projects

The README.Rmd file we've provided for this project is in the Source pane, and we can start by looking at the contents of the file (look at the YAML header, setup chunk, and packages).

As we've stated before, we won't be going through the code collects and downloads the data for this project. We'll provide links to those resources, but we're going to focus on the *workflow for using data to create tables, images, and products*.

Hit the enter/return key until we a few lines of padding beneath the packages code chunk. Our cursor should be around line 37.

We will begin this document with the following text under our first line header (#).

```
1 # Motivation  
2  
3 We read an interesting article on `fivethirtyeight` about the democratic debates in \  
4 June of 2019. In particular, the image below displays how voters had changed their m\  
5 inds after watching the candidates.
```

As the text above states, we saw something cool, and we're going to see if we can recreate or verify some of the information. The source of this inspiration comes from [this fivethirtyeight article¹⁵⁴](#).

Including an image in the .Rmd file: If you want to include an image in your README.Rmd file, read about the `knitr::include_graphics()` function by entering `?include_graphics` in the **Console** pane.

Import

Now that we have an idea why we're here, we want to start exploring the data. We're going to document all of our steps in the README.Rmd file.

We will start by importing all of the data for this project. We've put the raw data files in the `data/raw/` folder. **Always leave raw data in its own folder, and never alter the raw data files (read more [here¹⁵⁵](#)).**

In our README.Rmd file, we'll create a new level-two header (##) called 'Data files' and a code chunk labeled "locate-data" (it should look like the image below).

```
41 ## Data files  
42  
43 Below are the data files sorted by size:  
44  
45 ```{r locate-data}  
46  
47 ...  
48
```

¹⁵⁴<https://projects.fivethirtyeight.com/democratic-debate-poll/>

¹⁵⁵<https://amstat.tandfonline.com/doi/abs/10.1080/00031305.2017.1375987>

File management

To locate the files in `data/raw/` folder, we like to use the `fs` package. `fs` stands for ‘file system’, and this package gives us the ability to navigate our project folders and files.

Check out the package website [here¹⁵⁶](https://fs.r-lib.org/) for a full description and examples. `fs` is loaded as part of the `tidyverse`.

Start with a folder tree of the data folder.

```

1 # where are the data?
2 fs::dir_tree("data")
3 # data
4 #   └── processed
5 #     └── raw
6 #       ├── 538
7 #         └── 2019-07-06-Cand538Fav.csv
8 #       ├── google-trends
9 #         ├── 2019-07-10-Dems2020Night1Group1.rds
10 #        └── 2019-07-10-Dems2020Night1Group2.rds
11 #       ├── twitter
12 #         ├── 2019-07-06-Night01Tweets.rds
13 #         ├── 2019-07-06-Night01TweetsRaw.rds
14 #         ├── 2019-07-06-Night01TweetsUsers.rds
15 #         ├── 2019-07-06-Night02Tweets.rds
16 #         ├── 2019-07-06-Night02TweetsRaw.rds
17 #         └── 2019-07-06-Night02TweetsUsers.rds
18 #       └── wikipedia
19 #         ├── 2019-07-10-WikiDemAirTime01Raw.csv
20 #         ├── 2019-07-10-WikiDemAirTime02Raw.csv
21 #         └── 2019-07-25-PollingCriterionRaw.csv

```

The raw data are in four separate folders, each representing the data source (538, `google-trends`, `twitter`, `Wikipedia`).

Size can be a significant impediment to getting your work done quickly, so it’s best to determine the size of the raw data files before importing.

The code chunk below tells us how big each file is and the folder in which its located.

¹⁵⁶<https://fs.r-lib.org/>

```

1 # how big are each of these files?
2 fs::dir_info(path = "data", recurse = TRUE) %>%
3   # only files
4   filter(type == "file") %>%
5   group_by(path) %>%
6   # sort by size
7   tally(wt = size, sort = TRUE)
8 # A tibble: 12 x 2
9 #   path                               n
10 #   <fs::path>                         <fs::bytes>
11 # 1 data/raw/twitter/2019-07-06-Night02Tweets.rds      7.99M
12 # 2 data/raw/twitter/2019-07-06-Night02TweetsRaw.rds    7.87M
13 # 3 data/raw/twitter/2019-07-06-Night01Tweets.rds      6.92M
14 # 4 data/raw/twitter/2019-07-06-Night01TweetsRaw.rds    6.83M
15 # 5 data/raw/twitter/2019-07-06-Night02TweetsUsers.rds   1.9M
16 # 6 data/raw/twitter/2019-07-06-Night01TweetsUsers.rds   1.65M
17 # 7 data/raw/google-trends/2019-07-10-Dems2020Night1Group2.rds 111.47K
18 # 8 data/raw/google-trends/2019-07-10-Dems2020Night1Group1.rds 109.52K
19 # 9 data/raw/wikipedia/2019-07-25-PollingCriterionRaw.csv 722
20 # 10 data/raw/538/2019-07-06-Cand538Fav.csv            581
21 # 11 data/raw/wikipedia/2019-07-10-WikiDemAirTime02Raw.csv 202
22 # 12 data/raw/wikipedia/2019-07-10-WikiDemAirTime01Raw.csv 191

```

The output tells us the twitter data files are the largest (7.99M). All of these files are small enough for RStudio.Cloud to handle, though.

We can import the data using the `01-import.R` file in the `code` folder. Feel free to open this file and examine its contents in the **Source** pane. We will use the `base::source()` function to run all of the code in the `01-import.R` file.

```

1 # import data using 01-import.R file
2 source("code/01-import.R")

```

The `01-import.R` file loads all of the data files into the RStudio.Cloud session. We can verify this by examining the contents of the environment with `base::ls()`.

```
1 base::ls()
```

The output tells us the following objects are in our working environment.

```

1 # [1] "google_data_files"
2 # [2] "google_data_path"
3 # [3] "GoogleData"
4 # [4] "GSheetCand538Fav"
5 # [5] "GTrendDems2020Night1G1"
6 # [6] "GTrendDems2020Night1G2"
7 # [7] "twitter_data_files"
8 # [8] "twitter_data_path"
9 # [9] "twitter_users_data_files"
10 # [10] "TwitterData"
11 # [11] "TwitterUsersData"
12 # [12] "wiki_data_files"
13 # [13] "wiki_data_path"
14 # [14] "WikiData"
15 # [15] "WikiDemAirTime01Raw"
16 # [16] "WikiDemAirTime02Raw"
17 # [17] "WikiPollCriterionRaw"

```

We can also see these objects in the **Environment** pane.

Data = tibbles & Lists

Values = vectors

The *Global Environment* holds all of our imported data.

Document changes with Git

So, we've imported some data into the RStudio.Cloud session, but we need to make sure we're keeping track of the changes to our files.

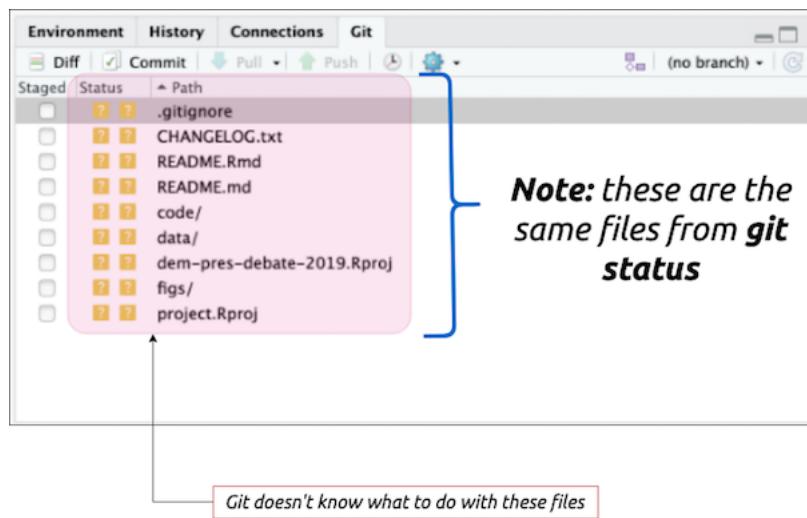
We can do this by checking our `git status` in the **Terminal** pane.

```

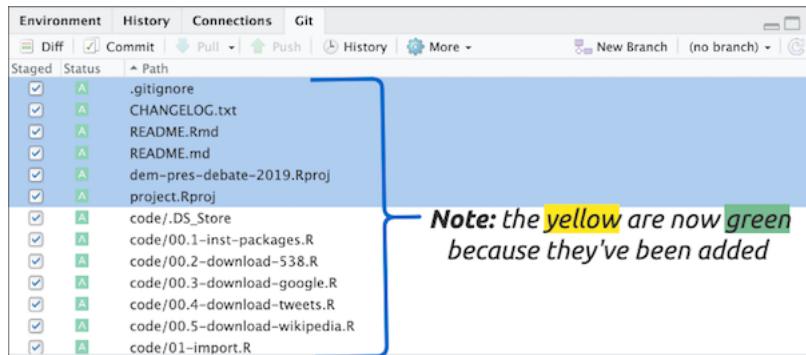
1 $ git status
2 On branch master
3
4 Initial commit
5
6 Untracked files:
7   (use "git add <file>..." to include in what will be committed)
8
9     .gitignore
10    CHANGELOG.txt
11    README.Rmd
12    README.md
13    code/
14    data/
15    dem-pres-debate-2019.Rproj
16    figs/
17    project.Rproj
18
19 nothing added to commit but untracked files present (use "git add" to track)

```

Git is telling us we have quite a few new files in our project (which we expected), but that to Git to pay attention to them, we need to add them. We can do this in the **Git** pane RStudio provides (it's by the **Environment** pane).



We will click on each checkbox next to the files in the **Git** pane. Don't be alarmed if this list expands to include all the subfiles and folders.



Now we can re-check git status to see what just happened.

```

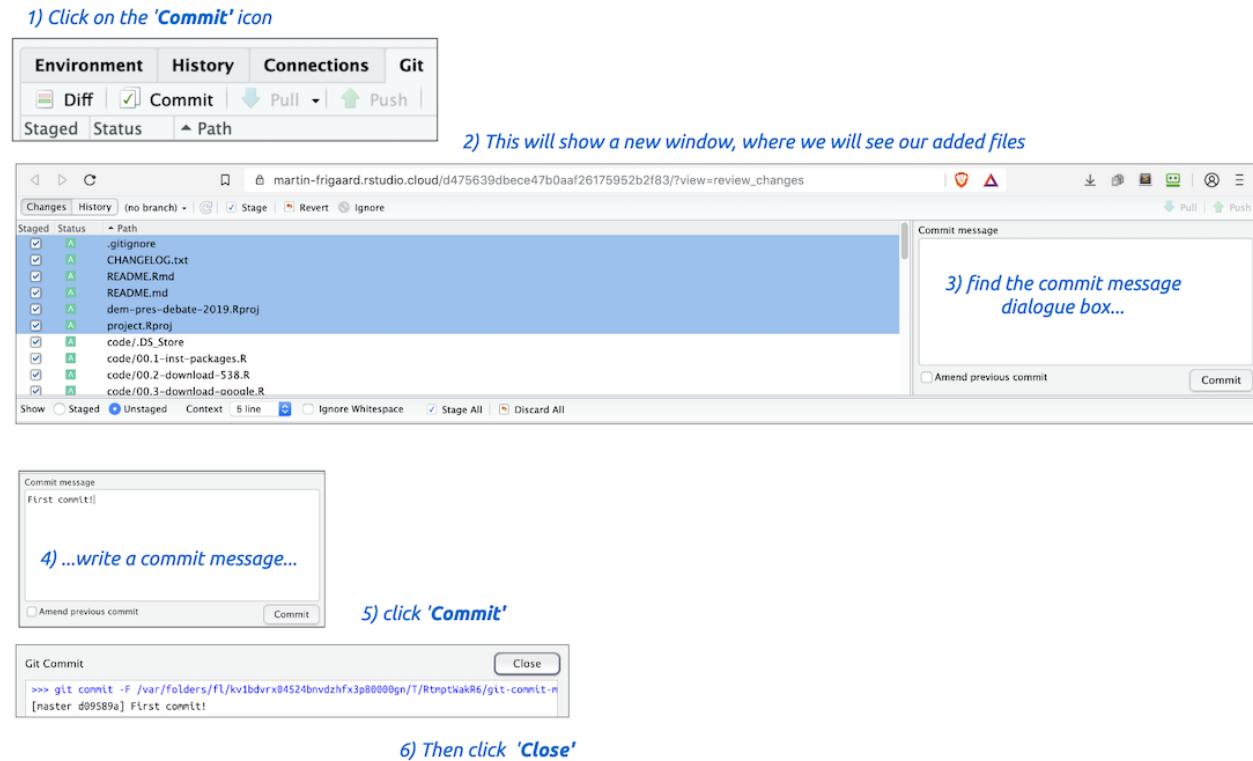
1 $ git status
2
3 On branch master
4
5 Initial commit
6
7 Changes to be committed:
8   (use "git rm --cached <file>..." to unstage)
9
10    new file:  .gitignore
11    new file:  CHANGELOG.txt
12    new file:  README.Rmd
13    new file:  README.md
14    new file:  code/.DS_Store
15    new file:  code/00.1-inst-packages.R
16    new file:  code/00.2-download-538.R
17    new file:  code/00.3-download-google.R
18    new file:  code/00.4-download-tweets.R
19    new file:  code/00.5-download-wikipedia.R
20    new file:  code/01-import.R
21    new file:  code/02-wrangle.R
22    new file:  data/.DS_Store
23    new file:  data/processed/.DS_Store
24    new file:  data/raw/.DS_Store
25    new file:  data/raw/538/2019-07-06-Cand538Fav.csv
26  # omitted

```

This long list of files shows all the contents have been added and are ready to be committed.

Committing the changes

If we click on the *Commit* icon in the Git pane, this will open a new window in RStudio.Cloud. In this window, we will enter a commit message (“First commit!”), and click the *Commit* button. We’ve outlined this entire process in the schematic below:



A quick lesson in compassionate programming

Your code will always be communicating to at least two audiences: your computer, and your future self. Be friendly to both of them!

R code can be difficult to read, but pipes (%>%) can help with clarity. The pipe is part of the `magrittr` package¹⁵⁷, and it takes code written like this:

```
1 outer_function(inner_function(Data_X), Data_Y)
```

And makes it look like this:

¹⁵⁷<https://cran.r-project.org/web/packages/magrittr/vignettes/magrittr.html>

```

1 Data_X %>% # do this
2   inner_function() %>% # then do this
3   outer_function(Data_Y)

```

`%>%` is a form of [syntactic sugar¹⁵⁸](#), which is a fancy way of saying “*something that helps us communicate better.*” You’ll see the pipe throughout the project’s R code files, and you can always read it as, “*do this, then do this.*”

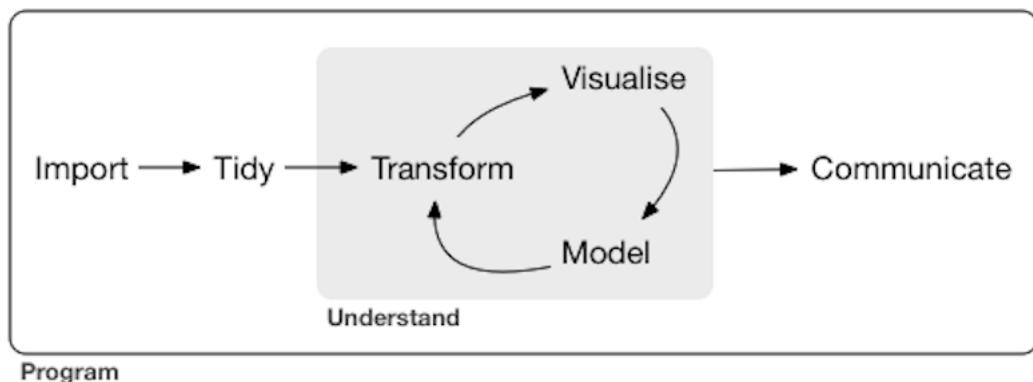
Wrangle

In the last section, we used the `README.Rmd` file to import our data into RStudio. We’ll continue using the `README.Rmd` file to add the contents from one additional .R script, `02-wrangle.R`. We will then create a new section and code file for visualizing the data (`03-visualize.R`).

We’ve provided a lot of code and comments in the `02-wrangle.R` script for you to explore, revise, and adapt to your liking. In the next few sections, we are going to move the code from the `02-wrangle.R` script into a new part of the `README.Rmd` (you probably guessed it “**Wrangle**”).

Data wrangling is whatever steps we needed to take the raw data into something we can use to create a table, visualization, model, etc. You’ll sometimes see ‘wrangling’ referred to as cleaning or munging.

The `02-wrangle.R` script prepares the data from the `01-import.R` for the visualizations. If you think back to the process outlined in the figure in [R for Data Science¹⁵⁹](#), you will notice that wrangle isn’t listed explicitly. Wrangling includes both ‘Tidy’ and ‘Transform’ steps (both of these need to happen before any visualizations or models can be properly run).



The simplest way to include the wrangling script in the `README` file is to create a code chunk, insert the `base::source()` function, and enter the path to the `02-wrangle.R` file.

However, we want to be helpful to our future selves, so we will include some language that describes what the functions are doing above each code chunk.

¹⁵⁸https://en.wikipedia.org/wiki/Syntactic_sugar

¹⁵⁹<https://r4ds.had.co.nz/introduction.html>

The first data that we need to wrangle is the Wikipedia tables (seen in the `02-wrangle.R` file on the section below).

```

21 # import -----
22 source("code/01-import.R")
23
24 # wrangle wikipedia data night 1 -----
25 WikiDemAirTime01 <- WikiDemAirTime01Raw %>%
26   magrittr::set_colnames(value = c("candidate", "airtime_night1"))
27
28 WikiDemAirTime01 <- WikiDemAirTime01 %>%
29   dplyr::filter(candidate %in% c("Night one airtime", "Candidate") &
30                 airtime_night1 %in% c("Night one airtime", "Airtime (min.)[58]")) %>%
31   dplyr::mutate(airtime_night1 = as.numeric(airtime_night1))
32
33 # wrangle wikipedia data night 2 -----
34 WikiDemAirTime02 <- WikiDemAirTime02Raw %>%
35   magrittr::set_colnames(value = c("candidate", "airtime_night2"))
36 WikiDemAirTime02 <- WikiDemAirTime02 %>%
37   dplyr::filter(candidate %in% c("Night two airtime", "Candidate") &
38                 airtime_night2 %in% c("Night one airtime", "Airtime (min.)[58]")) %>%
39   dplyr::mutate(airtime_night2 = as.numeric(airtime_night2))

```

packages
import
wrangle wikipedia data night 1
wrangle wikipedia data night 2
wrangle polling criterion data
export wiki tables
wrangle Google trend data
bind
gender
join Gtrend with wikipedia data
poll_perc_cat
mapping data (by region)
create processed data folder
export GtrendDems2020Intere...
export GtrendWikiOTAirTime

Both data sets had initially been Wikipedia (.html) tables, so they aren't immediately ready to go. We will create some new column names, remove some rows that used to be column headers, and make the airtime variable numeric.

The polling criterion Wikipedia data starting at the section titled, `wrangle polling criterion data`. This section creates a list of candidates (in `cand_names_wiki`) and uses it to filter out the observations we want. Check out [this webinar¹⁶⁰](#) to get an understanding of how `dplyr`¹⁶¹'s verbs work.

```

1 # wrangle polling criterion data -----
2 # create list from names using dput()
3 # dput(WikiPollCriterionRaw[ 1:11, 1])
4 cand_names_wiki <- c("Warren[note 2]",
5                      "O'Rourke[note 2]",
6                      "Booker[note 2]",
7                      "Klobuchar[note 2]",
8                      "Castro[note 2]",
9                      "Gabbard",
10                     "Ryan",
11                     "Inslee",
12                     "de Blasio",
13                     "Delaney")
14 # subset WikiPollCriterionRaw with list from above
15 WikiPollCriterion <- WikiPollCriterionRaw %>%

```

¹⁶⁰<https://www.rstudio.com/resources/webinars/data-wrangling-with-r-and-rstudio/>

¹⁶¹<https://dplyr.tidyverse.org/>

```

16 # this will remove all candidates not listed above
17 dplyr::filter(`Candidates drawn for the June 26 debate` %in% cand_names_wiki)

```

After we have wrangled the Wikipedia tables, the script exports these files to a new processed/ folder. Exporting them into a separate folder helps ensure they won't be accidentally altered or mistaken for the data files in the raw/ data folder.

The export section also timestamps each file, so we know the last time it was created. Read more about importing and exporting data in [this RStudio cheatsheet¹⁶²](#).

The Google trend data are a little more complicated because they come into RStudio.Cloud as a list, which is a data container in R that [doesn't have to be rectangular¹⁶³](#).

The image below outlines what each portion of code is doing. These are relatively everyday wrangling tasks, so we recommend going back or bookmarking these files as a reference.

The screenshot shows an RStudio interface with two tabs open: 'README.Rmd' and '02-wrangle.R'. The '02-wrangle.R' tab contains R code for wrangling Google trend data. The code is annotated with several callouts:

- Get data out of an R object into a rectangle we can manipulate easily**: Points to the first few lines of code that convert lists into tibbles.
- Stick two data sets (rectangles) together**: Points to the 'bind_rows' function used to combine two tibbles.
- Create a new variable (column) on multiple conditions**: Points to the 'case_when' function used to create a gender column based on keyword presence.
- A sidebar on the right lists various wrangling operations:

 - packages
 - import
 - wrangle wikipedia data night 1
 - wrangle wikipedia data night 2
 - wrangle polling criterion data
 - export wiki tables
 - wrangle Google trend data
 - bind
 - gender
 - join Gtrend with wikipedia data
 - poll_perc_cat
 - mapping data (by region)
 - create processed data folder
 - export GtrendDems2020Intere...
 - export GtrendWikiOTAirTime

We have different sources of data in RStudio right now (Wikipedia and Google trend data). They both have information on Candidates though. Often we'll want to join two (or more) data sets on a common column (like candidates). We will perform an example of this in the section outlined below.

¹⁶²<https://raw.githubusercontent.com/rstudio/cheatsheets/master/data-import.pdf>

¹⁶³<http://adv-r.had.co.nz/Data-structures.html>

```

# join Gtrend with wikipedia data -----
# sort alphabetically, join on id
WikiDemAirTime01 <- WikiDemAirTime01 %>% dplyr::arrange(desc(candidate))
# add id
WikiDemAirTime01 <- WikiDemAirTime01 %>%
  mutate(candidate_id = row_number())

# create candidate_id for GTrendDems2020Debate01IOT
GtrendDems2020Debate01IOT <- GTrendDems2020Debate01IOT %>%
  dplyr::mutate(candidate_id = case_when(
    stringr::str_detect(string = keyword, pattern = "Warren") ~ 1,
    stringr::str_detect(string = keyword, pattern = "Ryan") ~ 2,
    stringr::str_detect(string = keyword, pattern = "Beto") ~ 3,
    stringr::str_detect(string = keyword, pattern = "Klobuchar") ~ 4,
    stringr::str_detect(string = keyword, pattern = "Inslee") ~ 5,
    stringr::str_detect(string = keyword, pattern = "Gabbard") ~ 6,
    stringr::str_detect(string = keyword, pattern = "Delaney") ~ 7,
    stringr::str_detect(string = keyword, pattern = "de Blasio") ~ 8,
    stringr::str_detect(string = keyword, pattern = "Castro") ~ 9,
    stringr::str_detect(string = keyword, pattern = "Booker") ~ 10) ) %>%
  dplyr::arrange(desc(candidate_id))

# Join WikiDemAirTime01 to GTrendDems2020Debate01IOT on candidate_id
GtrendWikiIOTAirTime <- GTrendDems2020Debate01IOT %>%
  dplyr::left_join(x = .,
                   y = WikiDemAirTime01,
                   by = "candidate_id")

```

gender
join Gtrend with wikipedia data
poll_perc_cat
mapping data (by region)
create processed data folder
export GtrendDems2020Interest...
export GtrendWikiIOTAirTime

Prep work to get a common id in each data set

Joining two data sets 'candidate_id'

The process usually isn't so involved, but we included extra to give more explicit instructions. Be sure to check out the [relational data chapter](#)¹⁶⁴ the R for Data Science book.

We'll also be creating a map with the Google (or Twitter) data. Doing this requires another common task, which is loading a dataset from a package in R. The code below loads a state-level map into RStudio.Cloud and joins it to the Google trend data.

```

# mapping data (by region) -----
# convert to tibble (another data structure in R)
GtrendDems2020IBRGroup1 <- tibble::as_tibble(GtrendDems2020Night1G1$interest_by_region)
GtrendDems2020IBRGroup2 <- tibble::as_tibble(GtrendDems2020Night1G2$interest_by_region)
# bind Dems2020IBRGroup1 Dems2020IBRGroup2 together
GtrendDems2020IBR <- bind_rows(GtrendDems2020IBRGroup1,
                                 GtrendDems2020IBRGroup2, .id = "data")

# convert the region to lowercase
GtrendDems2020InterestByRegion <- GtrendDems2020IBR %>%
  dplyr::mutate(region = stringr::str_to_lower(location))

# create a data set for the states in the US
statesMap = ggplot2::map_data("state")
# now merge the two data sources together
GtrendDems2020InterestByRegion <- GtrendDems2020InterestByRegion %>%
  dplyr::inner_join(x = .,
                    y = statesMap, |
                    by = "region")

```

wrangle Google trend data
bind
gender
join Gtrend with wikipedia data
poll_perc_cat
mapping data (by region)
create processed data folder
export GtrendDems2020Interest...
export GtrendWikiIOTAirTime

Load the data from the ggplot2 package and join it to the Google trend data

We also export the Google trend data with time-stamp into the processed/ folder. We should continue adding the code into the README.Rmd file until we're confident all the functions will run and we don't see any errors.

¹⁶⁴<https://r4ds.had.co.nz/relational-data.html>

```

131 ## Wrangle Google trend data
132
133 ````{r wrangle-google}
134 # wrangle Google trend data -----
135 # convert Dems2020Group1 to tibble
136 GTrendDems2020Group1IOT <- GTrendDems2020Night1G1$interest_over_time %>%
137   as_tibble()
138 # convert Dems2020Group2 to tibble
139 GTrendDems2020Group2IOT <- GTrendDems2020Night1G2$interest_over_time %>%
140   as_tibble()
141
142 # create numeric hits
143 GTrendDems2020Group1IOT <- GTrendDems2020Group1IOT %>%
144   dplyr::mutate(hits = as.numeric(hits))
145 GTrendDems2020Group2IOT <- GTrendDems2020Group2IOT %>%
146   dplyr::mutate(hits = as.numeric(hits))
147
148 # bind -----
149 GTrendDems2020Debate01IOT <- bind_rows(GTrendDems2020Group1IOT,
150                                         GTrendDems2020Group2IOT,
151                                         .id = "data")
152
153 # gender -----
154 GTrendDems2020Debate01IOT <- GTrendDems2020Debate01IOT %>%
155   dplyr::mutate(gender = case_when(
156     str_detect(keyword, "Elizabeth Warren") ~ "Women",
157     str_detect(keyword, "Amy Klobuchar") ~ "Women",
158     str_detect(keyword, "Tulsi Gabbard") ~ "Women",
159     TRUE ~ "Men"))
160
161 # get distinct
162 GTrendDems2020Debate01IOT <- GTrendDems2020Debate01IOT %>% distinct()
163 # GTrendDems2020Debate01IOT %>% glimpse(78)
164 ``
165

```

Motivation
Data files
Importing data
Wrangle wikipedia data
Wrangle Google trend data
Join Google trend and Wikipedia data
Add polling percentage categories variable
Mapping data
Export wrangled data

These sections correspond to similar sections in the 02-wrangle.R script.

Each section can get a new code block and whatever description you want to include.

Note: The `02-wrangle.R` file is in the `code/` folder, but you won't have to alter the file paths because you're using an RStudio project file. Read more about how these are so helpful to your workflow [here¹⁶⁵](#).

Visualize

OK, we've completed our section for the wrangling the data. We are going to insert a divider (***) and start a new visualize section (## Visualize) in the `README.Rmd` file.

We've created a `03-visualize.Rmd` file for you to download from Github. You can do this by typing the following code into your `README.Rmd` file:

¹⁶⁵<https://support.rstudio.com/hc/en-us/articles/200526207-Using-Projects>

The screenshot shows the RStudio interface with the README.Rmd file open. Line 308 contains the instruction `## Visualize copy this code into your README.Rmd file`. Below it, line 310 provides a download link: `Download this file for the visualization script: http://bit.ly/viz-2019-data`. Lines 312 through 315 show the R code for downloading the file: ````{r download-vis-rmd}`, `download.file(url = "http://bit.ly/viz-2019-data",`, `destfile = "03-visualize.Rmd")`, `````.

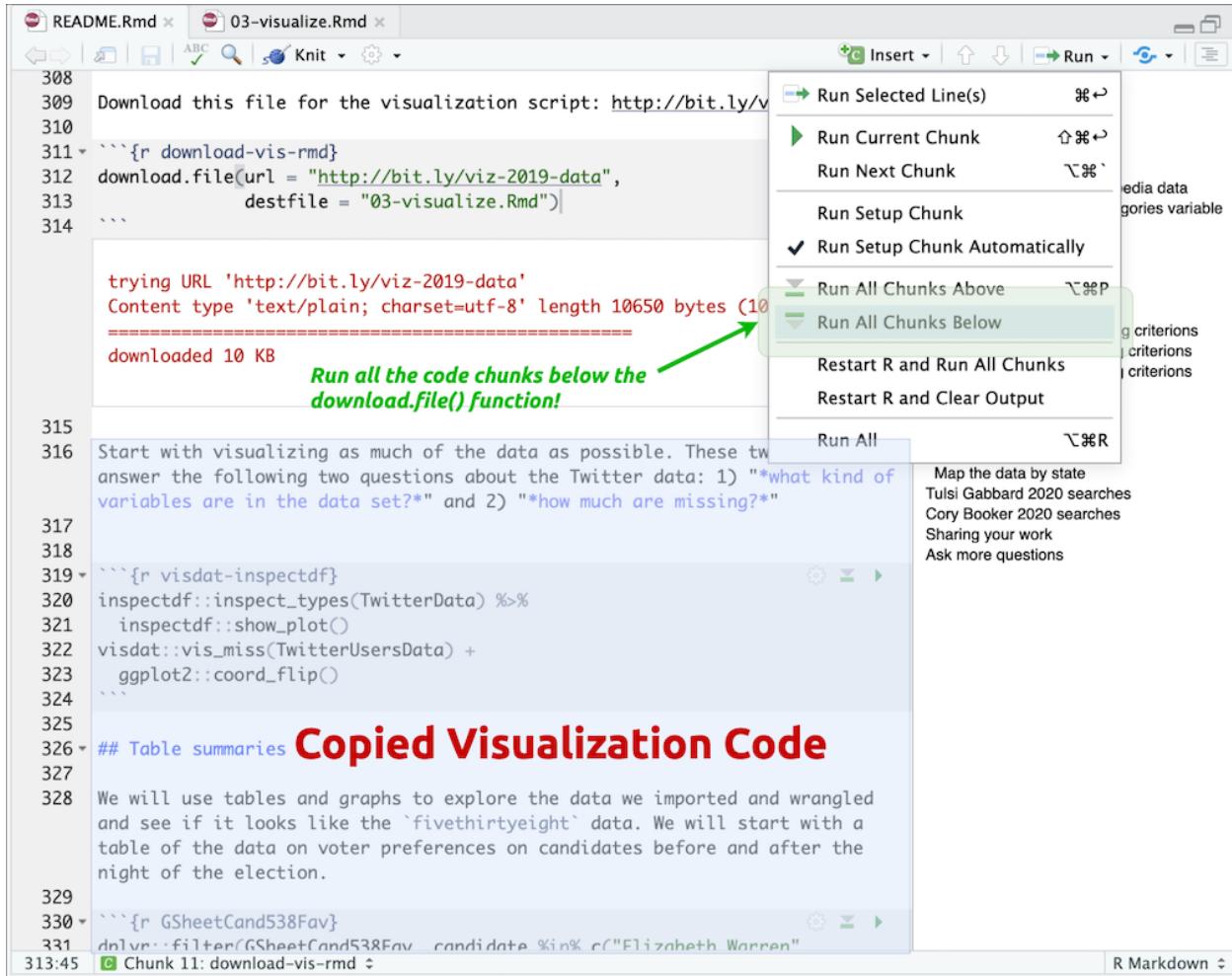
The hyperlink is here: <http://bit.ly/viz-2019-data>

After we've downloaded our `03-visualize.Rmd` file, we can open this file and start copying the code. We should begin at the line just below the `# Visualize data` header (it should be on about line 65) and extending to the end of the file.

The screenshot shows the RStudio interface with the `03-visualize.Rmd` file open. The code has been highlighted in blue. A large red box highlights the text **Visualization Code** and the instruction **Copy all the code highlighted in blue into the README.Rmd file.** The code itself includes several R code chunks and data wrangling commands like `source("code/02-wrangle.R")`, `wrangle`, `inspectdf`, `visdat`, and `GSheetCand538Fav`.

After selecting the code from `03-visualize.Rmd`, we should click on the line directly under the previous code chunk we used to download the `.Rmd` file.

After pasting the code from `03-visualize.Rmd` into the `README.Rmd` file, we can click on the *Run > Run All Chunks Below* (this will run all the code starting at line 319 until the end of the document).



The screenshot shows the RStudio Cloud interface with two tabs open: `README.Rmd` and `03-visualize.Rmd`. The `03-visualize.Rmd` tab contains R code for data visualization. A green arrow points from the text "Run all the code chunks below the download.file() function!" to the `Run All Chunks Below` option in the *Run* menu. The menu also includes options like `Run Selected Line(s)`, `Run Current Chunk`, `Run Next Chunk`, `Run Setup Chunk`, `Run Setup Chunk Automatically`, `Run All Chunks Above`, `Run All Chunks Below` (which is selected), `Restart R and Run All Chunks`, `Restart R and Clear Output`, and `Run All`. The `README.Rmd` tab shows the copied visualization code.

Running the code will create multiple tables and figures in the `README.Rmd` file. We'll go over these in more depth below. For now, we'll follow the directions at the bottom of the pasted code and “*Click knit to get the markdown file to share*.”

Knitting RMarkdown files

Clicking *Knit* (or clicking `shift+cmd+k`) activates the **R Markdown** pane in RStudio.Cloud and we see the code chunks compiling for the entire document.

The screenshot shows the RStudio Cloud interface with the R Markdown tab active. The document content includes R code for downloading data and inspecting it. Progress bars indicate the percentage of the document knitted: 45%, 47%, 49%, 51%, and 53%. Annotations explain these as the percentage of the document knitted and the labels/messages from the knitting process.

```

!.....
label: download-vis-rmd
trying URL 'http://bit.ly/viz-2019-data'
!.....
ordinary text without R code

!.....
label: visdat-inspectdf
Content type 'text/plain; charset=utf-8' length 10650 bytes (10 KB)
=====
downloaded 10 KB

!.....
ordinary text without R code

```

When the knitting process completes, a new browser window will pop up with our README.md document. The README.md will have sections of formatted text (from the Markdown), R code, and the various outputs.

The top of the file should list the title and the packages (just like our README.Rmd file)

The screenshot shows a browser window displaying the generated README.html file. The title is "2019 Democratic Debates data project". Below the title, the packages loaded are listed: tidyverse, Hmisc, janitor, magrittr, rtweet, visdat, inspectdf, ggthemes, and scales. An annotation points to the package list with the text "These are the packages we loaded".

2019 Democratic Debates data project

Jane Doe

```

library(tidyverse)
library(Hmisc)
library(janitor)
library(magrittr)
library(rtweet)
library(visdat)
library(inspectdf)
library(ggthemes)
library(scales)

```

Motivation

If we scroll down to the visualize section and look at the portion titled, *Candidates with high polling criterions*). We can see the different parts of the Rmarkdown file in the image below:

Header & text for this section

Code and comments for this graph

Graph output

More text interpreting the output

Candidates with high polling criterions

When we look at the candidates with the highest percent of polling criterion on the 26th, we see the following:

```
GtrendwikiIOTAirTime %>%  
  # limit to only the candidates with more than 10% of voters  
  dplyr::filter(poll_perc_fct == "> 10.0% of voters") %>%  
  # put data on the x,  
  ggplot2::ggplot(aes(x = date,  
                      # hits on the y  
                      y = hits,  
                      # color it by keyword  
                      color = keyword)) +  
  # use a line  
  ggplot2::geom_line(aes(group = keyword), show.legend = FALSE) +  
  # labels  
  ggplot2::labs(  
    x = "Date",  
    y = "Google search hits",  
    subtitle = "Google trends for candidates \nwith > 10.0% polling") +  
    # add themes from ggthemes  
    ggthemes::theme_wsj(base_size = 9.5,  
                        title_family = "mono") +  
    # use facets for both candidates on the same graph  
    ggplot2::facet_wrap(~ keyword, ncol = 2)
```

This shows Warren doing well after the first night. If we narrow this down to the week of the debates and widen the list of candidates, we see the following:

The file output is a *Preview* of our markdown file (README.md). Our browser renders the markdown as a webpage (README.html).

Rmarkdown to Word

Working in Rmarkdown give us other options for output, too. For example, we can knit this document to word (.docx) using the YAML header below.

```
1 output:  
2   word_document:  
3     fig_caption: yes  
4     highlight: kate  
5     keep_md: yes  
6     toc: yes
```

Rmarkdown to PDF

We can also create a .pdf file using this YAML header. This also requires making a few adjustments to the setup code chunk, and the summary of the mapping data under Summarize Interest Data.

```

1 # changes to setup
2 # create image folder ----
3 if (!file.exists("figs/")) {
4   dir.create("figs/")
5 }
6 # create data folder ----
7 if (!file.exists("data/")) {
8   dir.create("data/")
9 }
10 knitr::opts_chunk$set(
11   warning = FALSE,
12   message = FALSE,
13   echo = TRUE, # show all code
14   tidy = FALSE, # cleaner code printing
15   size = "small", # smaller code
16   fig.path = "figs/")
17 knitr::opts_knit$set(
18   width = 78)
19 base::options(tibble.print_max = 25,
20               tibble.width = 78)

```

Changes to Summarize Interest Data section

```

1 # recheck the structure
2 GtrendDems2020InterestByRegion %>%
3   skimr::skim_to_wide() %>%
4   dplyr::filter(type %in% c("integer", "numeric")) %>%
5   dplyr::select(variable,
6                 n,
7                 mean,
8                 sd,
9                 median = p50) # remove the hist

```

And use this YAML header.

```

1 output:
2   pdf_document:
3     highlight: kate
4     toc: yes

```

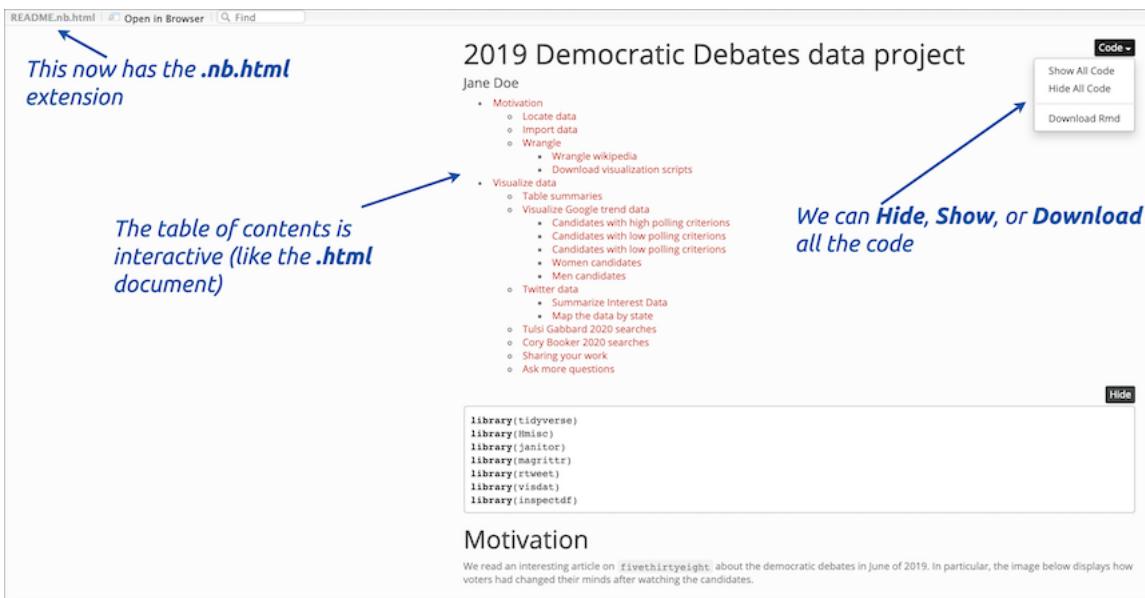
Rmarkdown Notebooks

Rmarkdown has its own dynamic document for data analysis called the **Rmarkdown Notebook**¹⁶⁶. These are unique because they create an output file with the `nb.html` extension, which > “*is a self-contained HTML file which contains both a rendered copy of the notebook with all current chunk outputs (suitable for display on a website) and a copy of the *.Rmd file itself.*”

You can create a Rmarkdown document using the following YAML header.

```
1 output:
2   html_notebook:
3     fig_height: 5
4     fig_width: 7
5     highlight: kate
6     theme: simplex
7     toc: yes
```

Note that this output changes the *Knit* icon to *Preview*. The output should look like the image below:



Being able to produce such a wide variety of documents from our Rmarkdown file makes it handy if you have collaborators who aren’t using our workflow. We’ve only touched the surface of available options for Rmarkdown documents. Please check the appendix in chapter 7 for more.

Extracting the .R from the .Rmd

But now we have all our visualize code in the `03-visualize.Rmd` file—what if we wanted this code in an `.R` script?

¹⁶⁶<https://bookdown.org/yihui/rmarkdown/notebook.html>

We can run the following code in the **Console** pane.

```
1 knitr::purl("03-visualize.Rmd")
```

We'll see the following script file gets generated.

```
1 processing file: 03-visualize.Rmd
2 | ..... | 100%
3 output file: 03-visualize.R
4
5 [1] "03-visualize.R"
```

The `knitr::purl()` function produces an R code file with all the code chunks from an existing `.Rmd` file, so it usually needs a little editing to be a well documented `.R` file.

Remember when we mentioned labeling the code chunks in the Rmarkdown files? This pays off when we run `knitr::purl()` because the labels become our section headers.

Part 6: Putting your project on Github

In the previous chapters, we've set up a Github account, learned how to download the files from a Github repository, upload files into RStudio.Cloud, create and run R code, and commit changes using Git.

In this final chapter, we are going to create some figures and graphs for this project, then put these changes on Github in a way for people to find and share.

Push the changes to Github

We've come to a point where we want to share our project with friends and colleagues. We need to commit the changes we've made, and then push these changes to a Github repository.

To do this, first we need to create a new Github repository. Follow the steps below to complete this:

1) Set up Github repository

Name your new Github repository the same thing as your RStudio.Cloud project.



Create a new repository

A repository contains all project files, including the revision history. Already have a project repository elsewhere? [Import a repository.](#)

Owner  **mjfrigaard2** / **Repository name *** 

Name this the same thing as your RStudio.Cloud project

Great repository names are short and memorable. Need inspiration? How about [fictional-carnival?](#)

Description (optional)

Democratic candidate project

Public
Anyone can see this repository. You choose who can commit.

Private
You choose who can see and commit to this repository.

Skip this step if you're importing an existing repository.

Initialize this repository with a README
This will let you immediately clone the repository to your computer.

Add .gitignore: **None** | Add a license: **None** | ⓘ

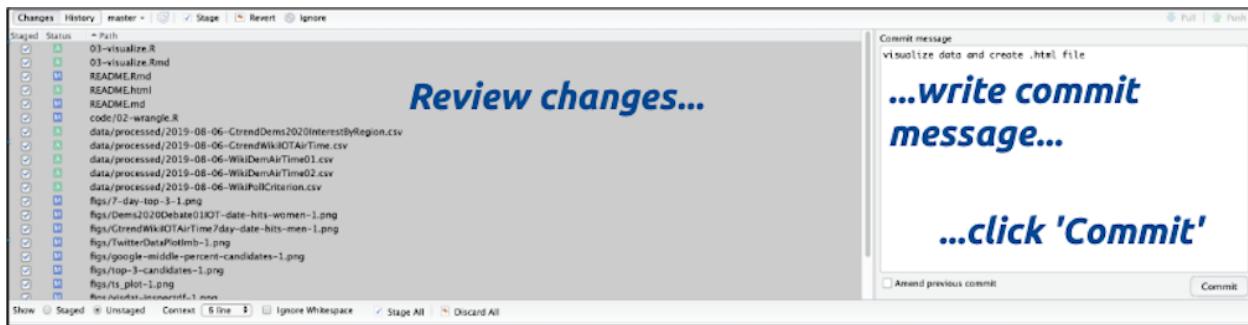
Create repository 

2) Add and commit file changes

After you have a Github repository, we will want to add and commit the changes we've made to our files. We can do this in the **Git** pane by clicking on the yellow question marks (and making them green As or blue Ms). When all the files have been added, click on the *Commit* icon. We've included an icon key so we know what's happening with each file.



Next we want to review the changed files, write a commit message, and commit these changes to Git.



3) Push changes to remote

After these changes have been committed in our local Git repository, we want to sync the local files with the Github repository we created just created. Fortunately, when we setup our Github repository, all the information we needed was on the landing page.

The screenshot shows a GitHub repository page for 'mjfrigaard2 / dem-pres-debate-2019'. At the top, there are tabs for Code, Issues (0), Pull requests (0), Projects (0), Wiki, Security, Insights, and Settings. Below the tabs, there's a section titled 'Quick setup — if you've done this kind of thing before' with options for 'Set up in Desktop' (selected), HTTPS, and SSH, along with the URL <https://github.com/mjfrigaard2/dem-pres-debate-2019.git>. It also suggests creating a new file or uploading an existing one, and recommends including a README, LICENSE, and .gitignore. Another section below says '...or create a new repository on the command line' with the following commands:

```
echo "# dem-pres-debate-2019" >> README.md
git init
git add README.md
git commit -m "first commit"
git remote add origin https://github.com/mjfrigaard2/dem-pres-debate-2019.git
git push -u origin master
```

Finally, there's a section '...or push an existing repository from the command line' with the command:

```
git remote add origin https://github.com/mjfrigaard2/dem-pres-debate-2019.git
git push -u origin master
```

We're going to type this information directly into the RStudio.Cloud *Terminal* pane. We will see a few error messages, but we can disregard them.

```

1 $ git remote add origin https://github.com/mjfrigaard2/dem-pres-debate-2019.git
2 $ git push -u origin master
3 # the error below can be disregarded--enter your username
4 error: cannot run rpostback-askpass: No such file or directory
5 Username for 'https://github.com': mjfrigaard2
6 # you will see another error--disregard again and enter password
7 error: cannot run rpostback-askpass: No such file or directory
8 Password for 'https://mjfrigaard2@github.com':
9 # now we see the objects being pushed
10 Counting objects: 82, done.
11 Delta compression using up to 16 threads.
12 Compressing objects: 100% (79/79), done.
13 Writing objects: 100% (82/82), 6.82 MiB | 2.82 MiB/s, done.
14 Total 82 (delta 16), reused 0 (delta 0)
15 remote: Resolving deltas: 100% (16/16), done.
16 To https://github.com/mjfrigaard2/dem-pres-debate-2019.git
17 * [new branch]      master -> master
18 Branch master set up to track remote branch master from origin.
```

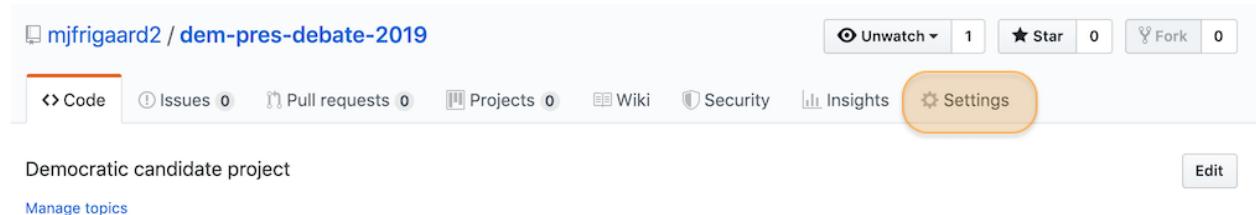
The final message tells us the changes have been pushed from the RStudio.Cloud Git repository to

the remote repository on Github. We can verify this by going to the repository landing page.

Sharing our work (with Github pages)

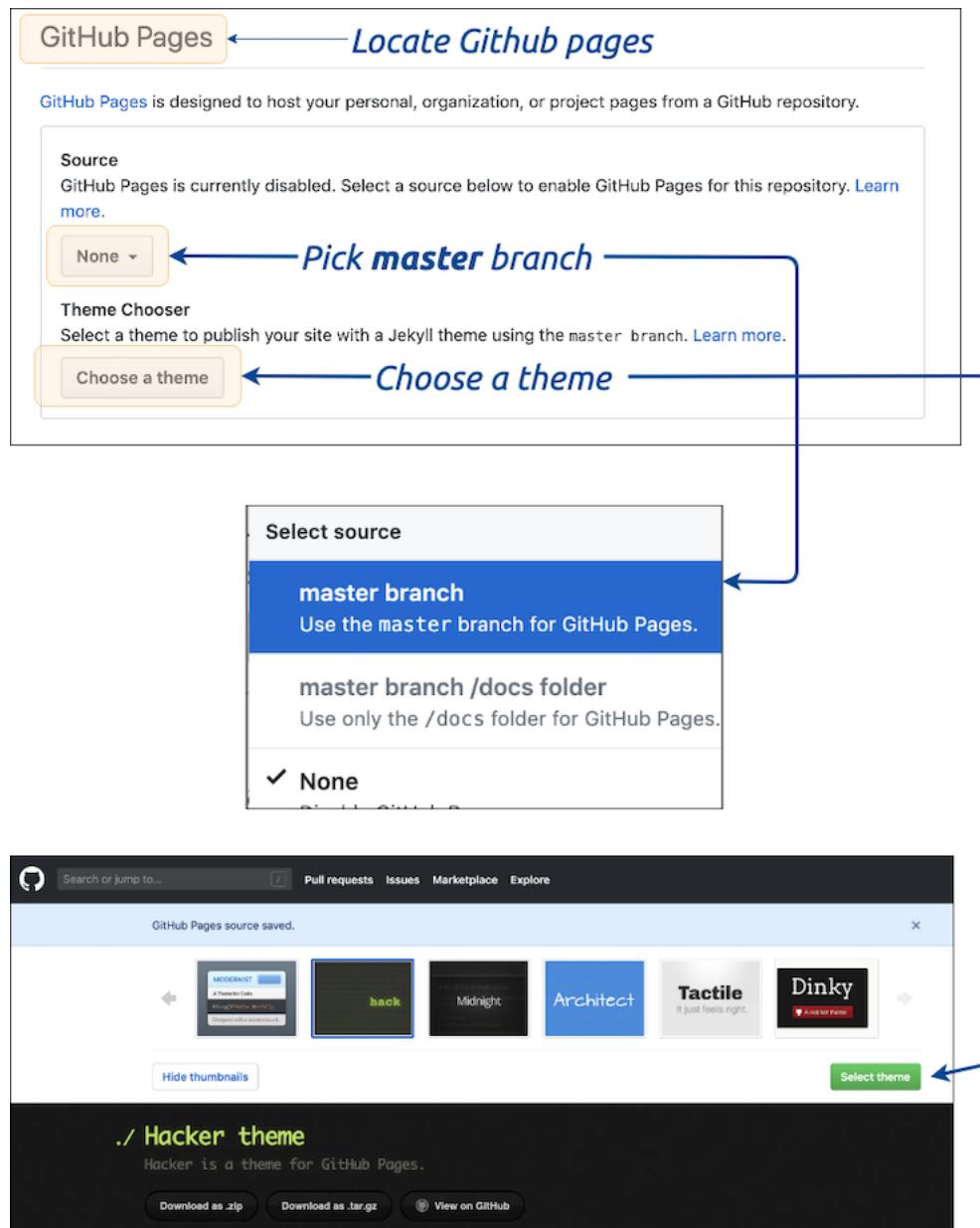
We're finally ready to put our stuff online and share it with everyone on Twitter, LinkedIn, etc. We'll be using [Github pages](#)¹⁶⁷ to make our project more discoverable.

Github pages are literally "Websites for you and your projects", so they're a perfect fit for what we want to create. To turn a repository into a Github page website, we will need to go into our *Settings*.



Scroll down the settings options until you find the Github pages. This is where we will tell Github that we want to use the master branch and pick a theme (we chose the Slate theme).

¹⁶⁷<https://pages.github.com/>



After completing these steps, Github will tell us the website url.

GitHub Pages *Here is our URL!*

GitHub Pages is designed to host your personal, organization, or project pages from a GitHub repository.

✓ Your site is published at <https://mjfrigaard2.github.io/dem-pres-debate-2019/>

This give us a website that uses our README.md file as the landing page.

The screenshot shows a web browser displaying a GitHub Pages site at mjfrigaard2.github.io/dem-pres-debate-2019/. The page title is "dem-pres-debate-2019" and the subtitle is "Democratic candidates". A "View on GitHub" button is visible in the top right. The main content features a large heading "2019 Democratic Debates data project" and a byline "Jane Doe". Below the heading is a code block containing R library imports:

```
library(tidyverse)
library(Hmisc)
library(janitor)
library(magrittr)
library(rtweet)
library(visdat)
library(inspectdf)
library(ggthemes)
library(scales)
```

Motivation

We read an interesting article on fivethirtyeight about the democratic debates in June of

Fortunately, our `README.md` file documents our entire project! But using Rmarkdown means we're not limited to the themes and options provided by Github pages.

1) Create a project website index file

Let's say we don't like the `Slate` theme we selected from Github pages. Fortunately, Rmarkdown and RStudio.Cloud give us the ability to create a new `.html` file we can use as our `index` file.

We'll head back over to RStudio.Cloud, change some of the settings on our Rmarkdown `README.Rmd`, rename the output file to `index.html` and commit and push all these changes to Github.

The `README.Rmd` file can be accessed in the `Files` pane. We are going to alter the `YAML` header by change the output: from `github_document` to `html_document`. Copy the `yaml` header below into the `README.Rmd` file

```
1 output:  
2   html_document:  
3     df_print: kable  
4     highlight: zenburn  
5     keep_md: yes  
6     theme: united  
7     toc: yes
```

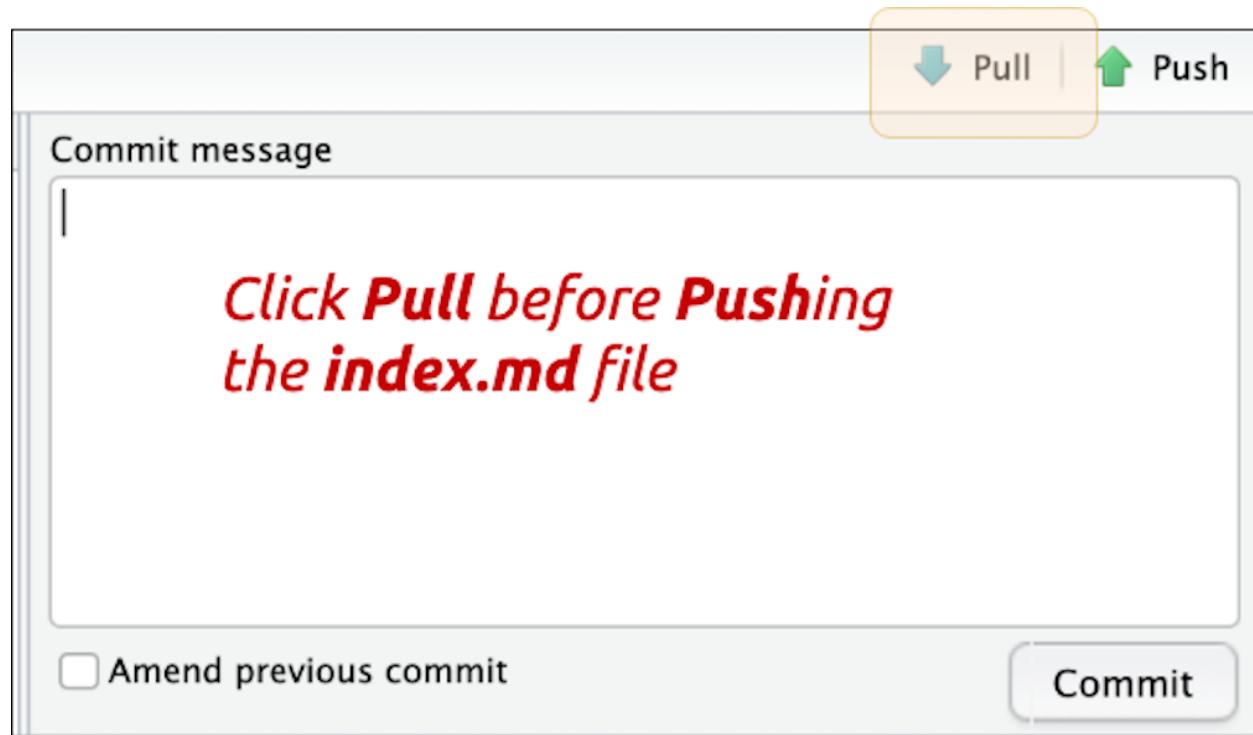
Save this file under the name `index.Rmd` after making the changes.

The settings displayed above are some of the ways we can customize our `.Rmd` files. Read more about the html documents [here in the Rmarkdown guide¹⁶⁸](#).

Notice the changes to the files in the **Git** pane. Follow the previous steps for adding and committing the file (don't push the changes yet!).

2) Pull the changes from Github (before pushing)

We made a few changes to the repository on Github when we set up the Github pages. Specifically, we added a `_config.yml` file. We need to make sure the local repository has the same files as the remote repository on Github, and we can do this by clicking on the *Pull* button (it's right next to the *Push* button) in the upper right corner of the review changes window.

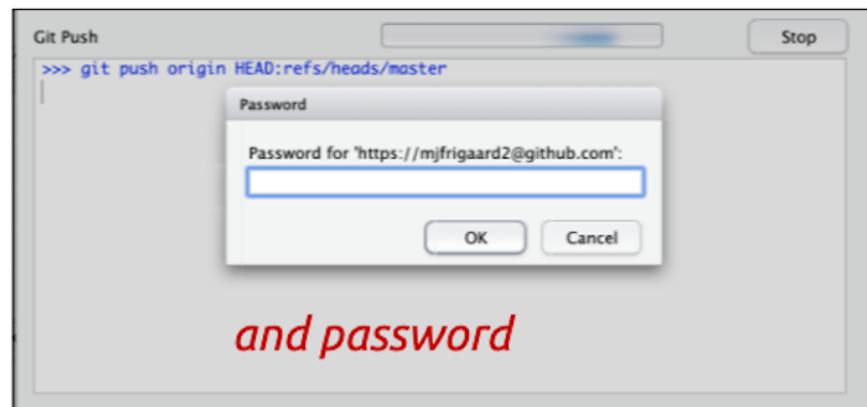


¹⁶⁸<https://bookdown.org/yihui/rmarkdown/html-document.html>

After clicking on the *Pull* button, the results should display something like the image below.

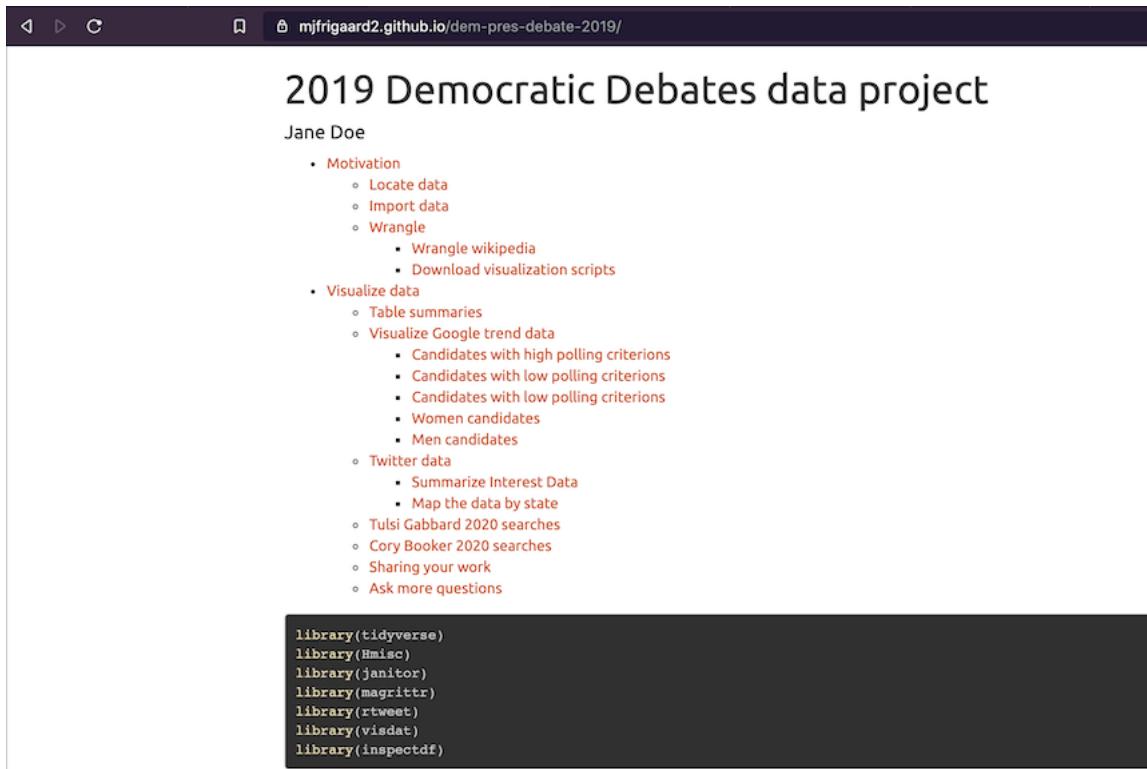


Now we can push the changes to Github.



After pushing these changes to Github, we can go back to the project website link. This will now take us to a nice website for our project.

Open the `index.md` file in the **Source** pane, click on the small gear next to the *Preview* button, and follow the directions in the figure below to setup the `index.html` output file.



When we click *Ok* and the window closes, we should see a new YAML header at the top of the `index.md` file.

```
1  -----
2  output:
3    html_document:
4      df_print: kable
5      fig_height: 5.5
6      fig_width: 7.5
7      highlight: kate
8      keep_md: yes
9      theme: simplex
10     toc: yes
11     toc_depth: 6
12  -----
```

The settings displayed above are some of the ways we can customize our `.Rmd` files. Read more about

the html documents [here in the Rmarkdown guide¹⁶⁹](#).

To see what these settings do, we will click the *Preview* button on the `index.md` file.

New table of contents (click on each section to jump to it)

- 2019 Democratic Debates data project
- Motivation
 - Data files
 - Importing data
 - Wrangle Wikipedia data
 - Wrangle Google trend data
 - Join Google trend and Wikipedia data
 - Add polling percentage categories variable
 - Mapping data
 - Export wrangled data
 - Visualize
 - Table summaries
 - Visualize Google trend data
 - Candidates with high polling criterions
 - Candidates with medium and high polling criterions (last seven days)
 - Candidates with low polling criterions (last seven days)
 - Women candidates
 - Men candidates
- Twitter data
 - Summarize Interest Data
 - Map the data by state
- Tulsi Gabbard 2020 searches
- Cory Booker 2020 searches
- Sharing your work
- Ask more questions

2019 Democratic Debates data project

Jane Doe

```
library(tidyverse)
library(htmllc)
library(janitor)
library(magrittr)
library(rtweet)
library(visdat)
library(inspectdf)
library(ggthemes)
library(scales)
```

Now we can add, commit, and push these changes to the Github remote following the steps outlined above to get a new project webpage.

The new `index.html` file has a few neat qualities: first, we can navigate this file like a webpage (the back button works when we click through the links in the document). Second, just about every computer will have a browser, so we don't have to worry about what version to save. Third, these files are easily transported to the web, which we will do in the next section.

Conclusion

There you have it! We've covered how to add and extract the code in an `.Rmd` file, how to push these changes to a Github repository, how to knit these files into `.html` documents,

Being able to create files in `.html` format is incredibly handy for interweaving various formatted text, code, and media (tables, images, and graphs).

¹⁶⁹<https://bookdown.org/yihui/rmarkdown/html-document.html>

Appendix

These are additional resources and notes from the preceding chapters.

Intro

We wrote this with no intention of anyone actually memorizing all this information. Our goal was for you to have a reference you could refer back to and check when you need to brush up on how these tools fit together.

Chapter 1

These aren't new concepts—we're really just repeating things we've learned from people smarter than us on these topics.

- Check out the [Data Carpentry lessons¹⁷⁰](#) for another great introduction to these tools
- The author of the paper above, Greg Wilson, has an [excellent blog¹⁷¹](#) on teaching technology and all things awesome.

Chapter 2

- The Ford Foundation report, "[Roads and Bridges](#)"¹⁷², outlines some other reason you should be using open-source software.
- Read these articles on attentional residue and multitasking (then try to stop doing it):
 - 1) [Why is it so hard to do my work? The challenge of attention residue when switching between work tasks¹⁷³](#)
 - 2) [Information, Attention, and Decision Making¹⁷⁴](#)
 - 3) [Causes, effects, and practicalities of everyday multitasking¹⁷⁵](#)
- See [Baumer et al.¹⁷⁶](#) for an in-depth summary of why you should abandon a copy + paste workflow

¹⁷⁰<https://datacarpentry.org/lessons/>

¹⁷¹<http://third-bit.com/>

¹⁷²<https://www.fordfoundation.org/about/library/reports-and-studies/roads-and-bridges-the-unseen-labor-behind-our-digital-infrastructure/>

¹⁷³<https://www.sciencedirect.com/science/article/pii/S0749597809000399>

¹⁷⁴https://aom.org/uploadedFiles/Publications/AMJ/June_2015_FTE.pdf

¹⁷⁵<https://www.sciencedirect.com/science/article/pii/S0273229714000513>

¹⁷⁶<https://arxiv.org/abs/1402.1894>

- This text is also an *opinionated technical manual*¹⁷⁷, and covers topics left out of typical statistics texts,

“Statisticians have long shied away from teaching process, with the complaint that it might limit the creativity necessary to tackle different analytical problems. However, by not teaching opinionated analysis development, we subject fledgling data to each individually spin their wheels in coming up with process for avoiding common and generalized problems.”

Chapter 3

Here are some excellent resources for safely adding the command line to your wheelhouse.

- The Unix Workbench¹⁷⁸
- Data Science at the Command Line¹⁷⁹
- Software Carpentry Unix Workshop¹⁸⁰

Chapter 4

There are a ton of resources on how to get started on a data analysis project. Some of our favorites are listed below:

- Data organization in spreadsheets¹⁸¹
- Best practices organizing data science projects¹⁸²
- Organizing Data Science Projects¹⁸³

Chapter 5

Chapter 6

¹⁷⁷<https://peerj.com/preprints/3210/>

¹⁷⁸<https://seankross.com/the-unix-workbench/>

¹⁷⁹<https://www.datascienceatthecommandline.com/>

¹⁸⁰<https://swcarpentry.github.io/shell-novice/>

¹⁸¹<https://www.tandfonline.com/doi/full/10.1080/00031305.2017.1375989>

¹⁸²<https://www.thinkingondata.com/how-to-organize-data-science-projects/>

¹⁸³<https://leanpub.com/universities/courses/jhu/cbds-organizing>