# Setup

1. Go to `www.github.com` and make a free account

2. Make sure you have a recent version (v1.1 or later) of RStudio `https://www.rstudio.com/products/rstudio/download/#download`

3. Keep `www.happygitwithr.com` open

4. Download these slides via: `https://github.com/kuriwaki/github-demo/raw/master/presentation-slides/kuriwaki_github.pdf`
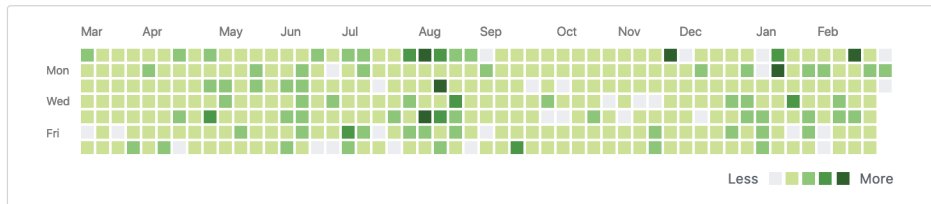
**Happy Git and GitHub for the useR**

*Jenny Bryan, the STAT 545 TAs, Jim Hester*

**Let's Git started**

2,412 contributions in the last year

# Introduction to git for social science students

(not software developers)

Shiro Kuriwaki

March 5, 2019

# Thanks for having me

**About me**

- ▶ G-4 in Government

- ▶ American Politics, elections and representation

- ▶ Before: Political data analytics (where I learned git from Annie Wang)

# Thanks for having me

**About me**

- ▶ G-4 in Government

- ▶ American Politics, elections and representation

- ▶ Before: Political data analytics (where I learned git from Annie Wang)

- ▶ I do some software development,

- ▶ but most of my work is applied ("substantive")

# Thanks for having me

**About me**

- ▶ G-4 in Government

- ▶ American Politics, elections and representation

- ▶ Before: Political data analytics (where I learned git from Annie Wang)

- ▶ I do some software development,

- ▶ but most of my work is applied ("substantive")

**My perspective**

- ▶ Version control is mandatory for programmers (and professional data scientists)

# Thanks for having me

**About me**

► G-4 in Government

► American Politics, elections and representation

► Before: Political data analytics (where I learned git from Annie Wang)

► I do some software development,

► but most of my work is applied ("substantive")

**My perspective**

► Version control is mandatory for programmers (and professional data scientists)

► but does it make sense for *applied* researchers who ...

► work with datasets that are

# Thanks for having me

## About me

- ▶ G-4 in Government

- ▶ American Politics, elections and representation

- ▶ Before: Political data analytics (where I learned git from Annie Wang)

- ▶ I do some software development,

- ▶ but most of my work is applied ("substantive")

## My perspective

- ▶ Version control is mandatory for programmers (and professional data scientists)

- ▶ but does it make sense for *applied* researchers who ...

- ▶ work with datasets that are **large**,

# Thanks for having me

**About me**

- ▶ G-4 in Government

- ▶ American Politics, elections and representation

- ▶ Before: Political data analytics (where I learned git from Annie Wang)

- ▶ I do some software development,

- ▶ but most of my work is applied ("substantive")

**My perspective**

- ▶ Version control is mandatory for programmers (and professional data scientists)

- ▶ but does it make sense for *applied* researchers who ...

- ▶ work with datasets that are **large**, **unstructured**,

# Thanks for having me

**About me**

- ► G-4 in Government

- ► American Politics, elections and representation

- ► Before: Political data analytics (where I learned git from Annie Wang)

- ► I do some software development,

- ► but most of my work is applied ("substantive")

**My perspective**

- ► Version control is mandatory for programmers (and professional data scientists)

- ► but does it make sense for *applied* researchers who ...

- ► work with datasets that are **large**, **unstructured**, **prone to change**,

# Thanks for having me

**About me**

- ▶ G-4 in Government

- ▶ American Politics, elections and representation

- ▶ Before: Political data analytics (where I learned git from Annie Wang)

- ▶ I do some software development,

- ▶ but most of my work is applied ("substantive")

**My perspective**

- ▶ Version control is mandatory for programmers (and professional data scientists)

- ▶ but does it make sense for *applied* researchers who ...

- ▶ work with datasets that are **large**, **unstructured**, **prone to change**, **with collaborators**

# Setting Expectations: Is it worth it?

**What do Gentzkow and Shapiro say?**

Definitely:

> *"It will probably take you a couple days to set up a repository and learn how you want to interact with [Version Control]. You will break even on that time investment within a month or two."*[1]

but also see[2]

R Studio Community

**Version control with Google Drive**

Brett-Johnson                                    2018-01-08

I've experimented using Google Drive and GitHub with my team (a small ecological research team) for version control and collaboration. I've found that both have there uses and I'm keen to share how I've been doing it so that I can hear from others how they are doing things, and whether I'm on the right track.

I initially started off committing everything I worked on to Github in different sub folders in the same repo. All of my internal analyses that aren't meant for a public report or peer reviewed paper went into different folders in the same general 'internal' private repo. This worked all right when it was just me using the repo. But when I brought a co-worker into the mix, we quickly realized what a pain it actually is to try to collaborate on GitHUB on a day to day basis. We were spending a load of time messing around with merge conflicts and all sorts of other un-intuitive issues. We felt GitHUB was cumbersome for day to day analysis collaboration internally.

So now I would like to move back to simply using Google Drive for internal analyses. Google drive is great for version controlling (especially now that you can 'name versions' in Google Drive similar to a GitHUB commit). I sometimes rely on the revision history of Google Drive to actually roll back a script, because it's way more intuitive than doing that in Git not to mention that every time you save your script in, it gets an un-named version in Google Drive, so the chances of not losing your work is actually greater using Google Drive. Google Drive allows you share all the files you and data you need, and using the here() package we shouldn't have to worry about working directories.

---

[1] Code and Data for Social Sciences: A Practioners Guide. 2014. https://perma.cc/5J9D-BTD6. Although I'm not sure about learning version control in "a couple of days" (I certainly couldn't!), I can guarantee reading their guide in its entirety *is* a time investment you'll break even on immediately.

[2] https://community.rstudio.com/t/version-control-with-google-drive/4032

**My (recommended) setup**

# Terminology

- ► **Git** is a particular type of software for version control (Subversion is another)

- ► **GitHub** is an app (recently bought by Microsoft) to host git on the web (Bitbucket is another)

- ► A **desktop client** is an app that connects a webhost like Github to your computer and facilitates simple tasks (here I use **RStudio**, there are many others)

- ► A **repository** is the fundamental unit of a version control, like a project folder.

# Terminology

- ► **Git** is a particular type of software for version control (Subversion is another)

- ► **GitHub** is an app (recently bought by Microsoft) to host git on the web (Bitbucket is another)

- ► A **desktop client** is an app that connects a webhost like Github to your computer and facilitates simple tasks (here I use **RStudio**, there are many others)

- ► A **repository** is the fundamental unit of a version control, like a project folder. Do not make a repository within a repository!

# Keep Track of how your results changed

*Problem: You tweak a regression specification and re-run your script, re-writing dozens of tables. You need to know how much your results changed*

# Keep Track of how your results changed

*You collect more data and re-run the regressions. Now how did the results change?*

# Tracking your text changes

*Problem: You start writing up your paper,* `draft.tex`

- ▶ The next day, you make a new draft. Do you overwrite?

# Tracking your text changes

*Problem: You start writing up your paper,* `draft.tex`

- ▶ The next day, you make a new draft. Do you overwrite?
- ▶ Or do you call it `draft_0305.tex` ? `draft_03052019.tex`?

# Tracking your text changes

*Problem: You start writing up your paper*, `draft.tex`

- ▶ The next day, you make a new draft. Do you overwrite?
- ▶ Or do you call it `draft_0305.tex` ? `draft_03052019.tex`?
- ▶ The next week, you find a single typo. Do you "Save As" with a new date?

# Tracking your text changes

*Problem: You start writing up your paper*, draft.tex

- ▶ The next day, you make a new draft. Do you overwrite?
- ▶ Or do you call it draft_0305.tex ? draft_03052019.tex?
- ▶ The next week, you find a single typo. Do you "Save As" with a new date?
- ▶ Three weeks later, you return to your paper. Your computer indicates that the file named draft_0305.tex was "Last modified March 12, 2019".

# Tracking your text changes

*Problem: You start writing up your paper,* `draft.tex`

- ▶ The next day, you make a new draft. Do you overwrite?
- ▶ Or do you call it `draft_0305.tex` ? `draft_03052019.tex`?
- ▶ The next week, you find a single typo. Do you "Save As" with a new date?
- ▶ Three weeks later, you return to your paper. Your computer indicates that the file named `draft_0305.tex` was "Last modified March 12, 2019".

# And more cool stuff like

**Getting a free, customizable, add-free website**

(instead of a click-and-drag Wordpress/Squarespace website)

# And more cool stuff like

**Work on a collaborative workbook**

(instead of needing to add people to your Dropbox)

# And more cool stuff like

**Work on a collaborative workbook**

(instead of needing to add people to your Dropbox)

# And more cool stuff like

**Contributing to / getting the latest on actual software packages**

(Github issues is the de facto communication of open-source developers)

# Terminology

► Files increment by **commit**s. The line-by-line changes from commits are called **diffs**.

# Terminology

- ► Files increment by **commit**s. The line-by-line changes from commits are called **diffs**.

- ► Commits have a human-readable **message**, and a serial code called a **SHA** (like `992bb07`).

# Terminology

- ▶ Files increment by **commit**s. The line-by-line changes from commits are called **diffs**.

- ▶ Commits have a human-readable **message**, and a serial code called a **SHA** (like `992bb07`).

- ▶ At least two copies of your repository exist: the **local** on your computer, and a **remote** (hosted on Github, with URL `https://github.com/user/repo.git`), which has the name **origin**



b13c7f7: Render as report
2736ead: Formula method
992bb07: Coauthor prefers str()

"commit"   A   B

Δ = "diff"

draft-01

Them                     You

origin

push ↑  ↓ pull

not your problem

daily work, your stuff

# Terminology

► Files increment by **commit**s. The line-by-line changes from commits are called **diffs**.

► Commits have a human-readable **message**, and a serial code called a **SHA** (like `992bb07`).

► At least two copies of your repository exist: the **local** on your computer, and a **remote** (hosted on Github, with URL `https://github.com/user/repo.git`), which has the name **origin**

► Once you make commits on your local, you **push** them to your remote. (The opposite of this is a **pull**)

b13c7f7: Render as report

2736ead: Formula method

992bb07: Coauthor prefers str()

"commit"   A   B

Δ = "diff"

draft-01

Them                          You

origin

push

pull

not your problem

daily work, your stuff

# Now, some caveats

**Only plain-text files get tracked**

This rules out:

# Now, some caveats

**Only plain-text files get tracked**

This rules out: PDFs,

# Now, some caveats

**Only plain-text files get tracked**

This rules out: PDFs, jpg's,

# Now, some caveats

**Only plain-text files get tracked**

This rules out: PDFs, jpg's, ,
Microsoft Word,

# Now, some caveats

**Only plain-text files get tracked**

This rules out: PDFs, jpg's, ,
Microsoft Word, any Microsoft
Office product,

# Now, some caveats

**Only plain-text files get tracked**

This rules out: PDFs, jpg's, ,
Microsoft Word, any Microsoft
Office product,Google Docs

# Now, some caveats

### Only plain-text files get tracked

This rules out: PDFs, jpg's, ,
Microsoft Word, any Microsoft
Office product,Google Docs, `.sav` ,
`.por` , `.dta` , `.Rds` , `RData` ...

### Which requires a switch to plain-text

# Now, some caveats

### Only plain-text files get tracked

This rules out: PDFs, jpg's, ,
Microsoft Word, any Microsoft
Office product,Google Docs, `.sav` ,
`.por` , `.dta` , `.Rds` , `RData` ...

### Which requires a switch to plain-text

Markdown ( `.md` )and TeX ( `.tex` )
for writing,

# Now, some caveats

### Only plain-text files get tracked

This rules out: PDFs, jpg's, , Microsoft Word, any Microsoft Office product,Google Docs, `.sav` , `.por` , `.dta` , `.Rds` , `RData` ...

### Which requires a switch to plain-text

Markdown ( `.md` )and TeX ( `.tex` ) for writing, Code ( `.R` , `.py` ) centered dataset generation, small datasets in `.csv` or `.txt` ,

# Now, some caveats

### Only plain-text files get tracked

This rules out: PDFs, jpg's, , Microsoft Word, any Microsoft Office product,Google Docs, `.sav` , `.por` , `.dta` , `.Rds` , `RData` ...

### Which requires a switch to plain-text

Markdown ( `.md` )and TeX ( `.tex` ) for writing, Code ( `.R` , `.py` ) centered dataset generation, small datasets in `.csv` or `.txt` , interweavers like `.Rmd` ,
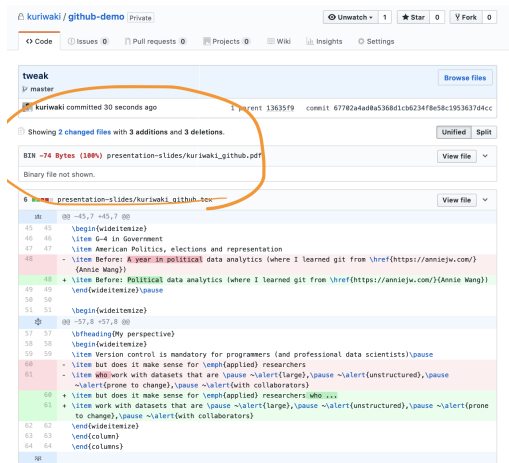
# Now, some caveats

### Only plain-text files get tracked

This rules out: PDFs, jpg's, , Microsoft Word, any Microsoft Office product,Google Docs, `.sav` , `.por` , `.dta` , `.Rds` , RData ...

### Which requires a switch to plain-text

Markdown ( `.md` )and TeX ( `.tex` ) for writing, Code ( `.R` , `.py` ) centered dataset generation, small datasets in `.csv` or `.txt` , interweavers like `.Rmd` , *Kieran Healy, "The Plain Person's Guide to Plain Text Social Science."*



Git is **not** built for storing data!

# Now, some caveats

### Only plain-text files get tracked

This rules out: PDFs, jpg's, , Microsoft Word, any Microsoft Office product,Google Docs, `.sav` , `.por` , `.dta` , `.Rds` , `RData` ...

### Which requires a switch to plain-text

Markdown ( `.md` )and TeX ( `.tex` ) for writing, Code ( `.R` , `.py` ) centered dataset generation, small datasets in `.csv` or `.txt` , interweavers like `.Rmd` , *Kieran Healy, "The Plain Person's Guide to Plain Text Social Science."*



Git is **not** built for storing data!

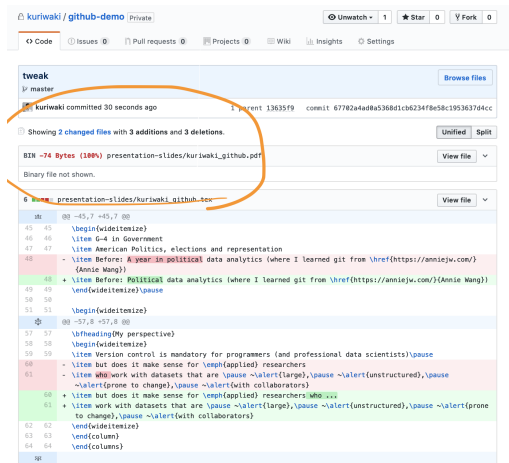(Rely on the usual Dropbox / Google Drive / Dataverse / Cloud Servers for that)

# Now, some caveats

### Only plain-text files get tracked

This rules out: PDFs, jpg's, , Microsoft Word, any Microsoft Office product, Google Docs, `.sav`, `.por`, `.dta`, `.Rds`, `RData` ...

### Which requires a switch to plain-text

Markdown (`.md`) and TeX (`.tex`) for writing, Code (`.R`, `.py`) centered dataset generation, small datasets in `.csv` or `.txt`, interweavers like `.Rmd`, *Kieran Healy, "The Plain Person's Guide to Plain Text Social Science."*



Git is **not** built for storing data!

(Rely on the usual Dropbox / Google Drive / Dataverse / Cloud Servers for that)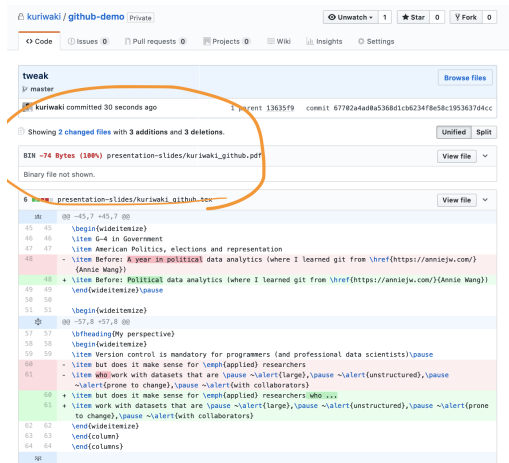