

Quantitative Research Workflow

Francisco Villamil

UC3M – February 8, 2024

Reason to think about a workflow

1. **Automate stuff:** you spend a lot of time on the computer so make it work for you
2. **Avoid errors:** we should not trust ourselves

Problems

- I revise a Word document too many times and end up with `final.docx`, `final2.docx`, `finalnov23.docx`, `finalFINAL.docx`
- After four months, I go back to a data project and have one 5000-thousand R file that I completely do not understand anything of it
- I run an old R file and suddenly it doesn't run because a file is missing - and I don't know where it is
- The project is ready and instead of using data since 1991, I want to use data since 1989: do I have to run everything again?
- I have a book or a dissertation (or a MA thesis) ready, with 60+ tables, and after it's ready, I found a little mistake that just changes the second decimal in the analyses tables - good luck updating all those tables
- Mistake in the code because it doesn't tell me where it's wrong (Real: journals.sagepub.com/doi/10.1177/20531680221126454)

Problems

- I revise a Word document too many times and end up with `final.docx`, `final2.docx`, `finalnov23.docx`, `finalFINAL.docx`
- After four months, I go back to a data project and have one 5000-thousand R file that I completely do not understand anything of it
- I run an old R file and suddenly it doesn't run because a file is missing - and I don't know where it is
- The project is ready and instead of using data since 1991, I want to use data since 1989: do I have to run everything again?
- I have a book or a dissertation (or a MA thesis) ready, with 60+ tables, and after it's ready, I found a little mistake that just changes the second decimal in the analyses tables - good luck updating all those tables
- Mistake in the code because it doesn't tell me where it's wrong (Real: journals.sagepub.com/doi/10.1177/20531680221126454)

Problems

- I revise a Word document too many times and end up with `final.docx`, `final2.docx`, `finalnov23.docx`, `finalFINAL.docx`
- After four months, I go back to a data project and have one 5000-thousand R file that I completely do not understand anything of it
- I run an old R file and suddenly it doesn't run because a file is missing - and I don't know where it is
- The project is ready and instead of using data since 1991, I want to use data since 1989: do I have to run everything again?
- I have a book or a dissertation (or a MA thesis) ready, with 60+ tables, and after it's ready, I found a little mistake that just changes the second decimal in the analyses tables - good luck updating all those tables
- Mistake in the code because it doesn't tell me where it's wrong (Real: journals.sagepub.com/doi/10.1177/20531680221126454)

Problems

- I revise a Word document too many times and end up with `final.docx`, `final2.docx`, `finalnov23.docx`, `finalFINAL.docx`
- After four months, I go back to a data project and have one 5000-thousand R file that I completely do not understand anything of it
- I run an old R file and suddenly it doesn't run because a file is missing - and I don't know where it is
- The project is ready and instead of using data since 1991, I want to use data since 1989: do I have to run everything again?
- I have a book or a dissertation (or a MA thesis) ready, with 60+ tables, and after it's ready, I found a little mistake that just changes the second decimal in the analyses tables - good luck updating all those tables
- Mistake in the code because it doesn't tell me where it's wrong (Real: journals.sagepub.com/doi/10.1177/20531680221126454)

Problems

- I revise a Word document too many times and end up with `final.docx`, `final2.docx`, `finalnov23.docx`, `finalFINAL.docx`
- After four months, I go back to a data project and have one 5000-thousand R file that I completely do not understand anything of it
- I run an old R file and suddenly it doesn't run because a file is missing - and I don't know where it is
- The project is ready and instead of using data since 1991, I want to use data since 1989: do I have to run everything again?
- I have a book or a dissertation (or a MA thesis) ready, with 60+ tables, and after it's ready, I found a little mistake that just changes the second decimal in the analyses tables - good luck updating all those tables
- Mistake in the code because it doesn't tell me where it's wrong (Real: journals.sagepub.com/doi/10.1177/20531680221126454)

Problems

- I revise a Word document too many times and end up with `final.docx`, `final2.docx`, `finalnov23.docx`, `finalFINAL.docx`
- After four months, I go back to a data project and have one 5000-thousand R file that I completely do not understand anything of it
- I run an old R file and suddenly it doesn't run because a file is missing - and I don't know where it is
- The project is ready and instead of using data since 1991, I want to use data since 1989: do I have to run everything again?
- I have a book or a dissertation (or a MA thesis) ready, with 60+ tables, and after it's ready, I found a little mistake that just changes the second decimal in the analyses tables - good luck updating all those tables
- Mistake in the code because it doesn't tell me where it's wrong (Real: journals.sagepub.com/doi/10.1177/20531680221126454)

Principles

1. Using computers

- Use *plain text* files as much as possible
- Customize your work tool
- Use a code/text editor and make it *yours*
- Learn how to use the Terminal (unix commands) and automate
- Use version control (git)

2. Coding and empirical projects

- Separate code in specific tasks, be tidy
- Integrate different parts of same project (R, tex...)
- **Automate output** (tables, plots ...)
- Use **functions** (automate code), i.e. do not do the same thing twice
- Checks and warnings in code
- (Optional: consider using *Makefile*)

Some resources

- Hadley Wickham's **R Style guide** (and the whole **Advanced R book** later on)
- Software Carpentry's lessons:
<https://software-carpentry.org/lessons/>
 - Especially **Unix Shell** and **Version Control with Git**
- Kieran Healy's *The Plain Person's Guide to Plain Text Social Science*: <https://plain-text.co/>
 - Although emacs is perhaps a bit too hardcore
- The best Git course I know is this: <https://gitexercises.fracz.com/>
- MIT's *The Missing Semester of Your CS Education*:
<https://missing.csail.mit.edu/>

Roadmap

Using computers

Coding better and organizing data projects

Plain text

- What's plain text?
- Quicker and easier to work with
- Cross-platform and does not depend on proprietary software
- Much better for the things you want to do
 - You can use version control on it
 - Closer to how machines work it - so easier for whatever related to machines (e.g. syncing two computers) → **example1**, **example2**
 - It's a base ingredient you can convert into whatever (e.g. with R, LaTeX, etc)

Customizing your computer

- [https://franvillamil.github.io/posts/setup_{macos}.html](https://franvillamil.github.io/posts/setup_macos.html)
- <https://github.com/franvillamil/templates>
- <https://github.com/franvillamil/configfiles>
- Examples: mdtopdf/docxtopdf, baserepos, Spectacle, ...

Code editor

- Choose and get used to some code editor
 - You're probably using the editor in RStudio, that's fine, but there are reasons to use better and more general tools
- You can customize these so suit your needs, e.g.:
 - Edit & run languages you use (R, Latex, whatever)
 - Small stuff that saves time, like snippets
 - Navigate a project
 - And much more complicated stuff we're not going to talk about and that I do not know so much about
- I use Sublime Text: <https://www.sublimetext.com/>
- Anyway, **don't use MS Word** (as much as possible)

Using the command line

- Think of it as the language to communicate with the OS
- No need that you become a computer wizard, but I personally think it pays off to learn a little bit
- Why?
 - Automate stuff in the computer (e.g. from updating local files to converting .docx into pdf)
 - Navigate and work with files faster
 - Version control, installing stuff, solving issues
 - Virtual machines
- **Note:** Unix/Mac vs Windows

Version control (Git)

- `final1.docx`, `finalfinal.docx`... but in a proper way

Version control (Git)

- `final1.docx`, `finalfinal.docx`... but in a proper way
- You want to keep control of all versions of a file, something like MS Word's 'Tracked changes' but just much better

Version control (Git)

- `final1.docx`, `finalfinal.docx`... but in a proper way
- You want to keep control of all versions of a file, something like MS Word's 'Tracked changes' but just much better
 - Keep a time machine of all versions of a file

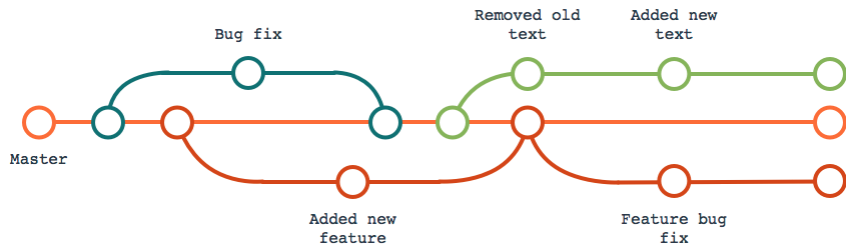
Version control (Git)

- `final1.docx`, `finalfinal.docx`... but in a proper way
- You want to keep control of all versions of a file, something like MS Word's 'Tracked changes' but just much better
 - Keep a time machine of all versions of a file
 - Allow collaboration between different people (or between two computers)

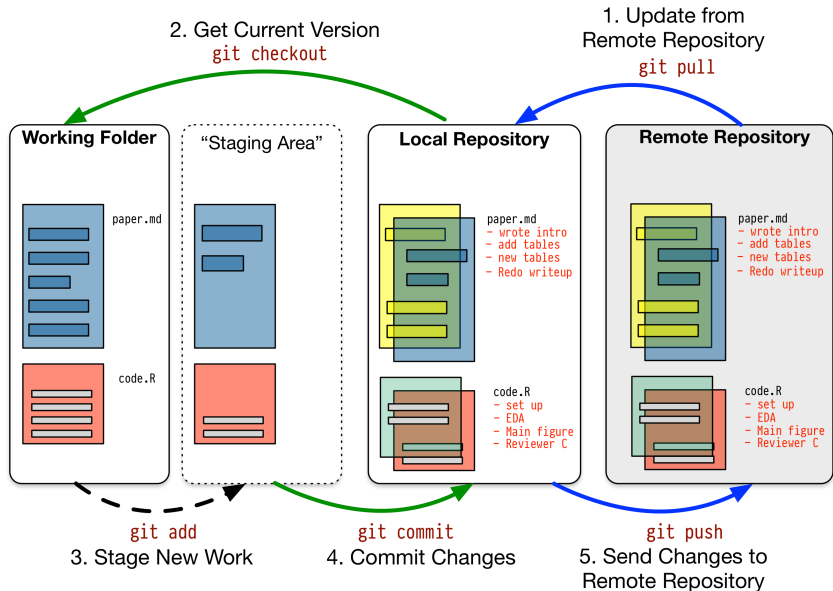
Version control (Git)

- `final1.docx`, `finalfinal.docx`... but in a proper way
- You want to keep control of all versions of a file, something like MS Word's 'Tracked changes' but just much better
 - Keep a time machine of all versions of a file
 - Allow collaboration between different people (or between two computers)
- There is more than one system, but most people use Git (and Github)

Version control



Version control



Version control - a note

- Version control works **much** better if you work with other people who also use version control, which is often not the case (at least not mine)
- Yet, there are two advantages to use it in my view:
 - Obvious one: keep older versions of a file
 - If you work with two computers, perhaps Google Drive/Dropbox do not work that well
 - Virtual machines (e.g. Google Cloud Computing, Amazon Web Services)

Roadmap

Using computers

Coding better and organizing data projects

Coding projects I: tasks as folders

- This applies especially to the R part of projects
- Do not create one huge R code file, use different files for different tasks
- Ideally, you probably want to do the same with the folder structure
 - And optionally, use *Makefile*

Coding projects I: tasks as folders

Coding projects II: filenaming

Coding projects III: integrate different parts

Coding projects IV: Makefile

Writing code I: automate via functions

Writing code II: write checks and warnings

- As you write code, always write checks using `stop()` or `warning()`, e.g.:
 - if a new data frame is built from merging, what should be the number of rows in the final df? or columns?
 - should two objects be identical?
 - do we have duplicated values by some ID?
 - do you expect `‘‘145‘‘` or `145` (character vs integer)?
 - ...

```
## Load

# Pre-invasion data
pre_data = read_dta("data/test_revisado.dta", encoding = "latin1") %>%
  left_join(read_dta("data/soft_300.dta", encoding = "latin1")[, c("response_id", "TS")]) %>%
  mutate(date = as.Date(str_sub(TS, 1, 10), "%m/%d/%Y"), post = 0) %>%
  rename(trust_army = Q49) %>%
  rename(Q41 = Q41_h)# labeling error

# Post-invasion data
post_data = read_dta("data/test5_revisado.dta", encoding = "latin1") %>%
  mutate(date = as.Date(str_sub(ts, 1, 10), "%m/%d/%Y"), post = 1) %>%
  # filter(date <= as.Date("2022-03-10")) %>%
  rename(trust_army = Q46)

## Checks
if(!identical(attr(pre_data$Q11, 'label'), attr(post_data$Q11, 'label')) &
  identical(attr(pre_data$Q42, 'label'), attr(post_data$Q42, 'label')) &
  identical(attr(pre_data$Q48, 'label'), attr(post_data$Q48, 'label')) &
  identical(attr(pre_data$Q43, 'label'), attr(post_data$Q43, 'label')) &
  identical(attr(pre_data$trust_army, 'label'), attr(post_data$trust_army, 'label')) &
  identical(attr(pre_data$Q47, 'label'), attr(post_data$Q47, 'label'))){stop("!")}
```



```
mentions_by_url = function(filename, keywords){  
  
  # Set up  
  month = gsub("output/webs_(\\d+-\\d+)\\.rds", "\\1", filename)  
  df = url_df[url_df$month == month,]  
  
  # Read  
  raw = readRDS(filename)  
  
  # Check  
  if(length(raw) != nrow(df)){stop("diff length df/raw! (1)")}  
}
```

Writing code II: write checks and warnings

- Also try to minimize errors, e.g. that you have visuals of real output, e.g.:
1. I use `print()` all the time to show length of stuff, number of missing data, etc
 2. `modelsummary` vs `stargazer` example
 - journals.sagepub.com/doi/10.1177/20531680221126454
 - github.com/franvillamil/streets_vox/blob/master/robust/rob.R