

Project management

John M. Drake & Andrew W. Park

Project management for reproducible research

This exercise is about practical stuff – managing scientific projects and workflows efficiently to avoid mistakes and keep the research moving forward.

R Projects

TBD

Version control, Git, and Github

Create a Github account

1. Go to website <http://github.com> and click on “sign up”
2. Create a profile (username, email, password)
3. Click “start a project”
4. Verify email address
5. Create a project using the provided form

Link Github account to R/Rstudio

(There are lots of ways to link your installation of R with Github. These instructions for use with R/Rstudio installed on a Windows computer such as in the teaching lab. Note that to use Github you also have to have Git installed on your computer)

1. Open RStudio
2. Select File -> New Project -> Version Control -> Git
3. Go to the website for the project you created and click the green “Clone or download” button; copy the URL
4. Paste the copied URL in the RStudio dialog box navigate to the directory where you would like your project to be cloned. This will create a local copy of all the files associated with your project

You should now have a “project” open. In the file explorer you should see the file `README.md` (if you created one), an `.Rproj` file, and another file `.gitignore` (which, conveniently, you can ignore). To include files (e.g. data, scripts) in the project, simply create or copy them into the directory. For instance:

Exercise. Copy the MERS data file `cases.csv` and paste it into your working directory.

Exercise. Create a new script following the prototype introduced in class yesterday. Your script should load the mers data and make a plot.

Now you have two files in your local directory that you don’t have on the Github repository. To synchronize your local copy and the repository, you need to *Commit* the changes and then *Push* the changes to Github.

Exercise. To *commit* a change, navigate to the “Git” tab in the top right explorer window. You will see a list of files in your work directory. Select the files that need to be pushed to Github and click on “Commit”. A dialog box will open. In the top right there is an editing window where you can register a comment describing the nature of the commit.

Exercise. Once you have committed one or more changes (and documented with comments), you need to *push* the changes to the archived version. To do this, click on “Push”. You will

need to enter your Github credentials in the dialog box that pops up. Now refresh the website for your project. The latest versions of your files should appear.

Github is very useful for collaboration. Multiple users can contribute code to the same project. Projects can be “branched” and “merged” and Github provides tools to identify and resolve conflicts when multiple programmers are working on the same project. Although it’s fine to use the Github repository as the definitive version and create a new local copy (a new project) every time, you don’t have to. Instead, you can leave your local working directory as it is and simply *pull* your collaborators contributions every time you start a new session.

Statistically literate programming with R Markdown

Literate programming is a programming paradigm due to Linux founder Donald Knuth in which natural language explanations of a program’s logic are interspersed with the code *snippets* that actually perform the computation. *Statistically literate programming* applies this paradigm to data analysis. Statistically literate programming is the idea that the thought process of the data analyst can be captured in a report that contains explanation and interpretation, the code used to perform an analysis, and the products of that analysis such as tables of data, quantities (e.g. p-values) and graphs.

There are a handful of ways that one can do statistically literate programming with R/RStudio. In this exercise, we will use R/Markdown and `knitr`. Markdown is a lightweight markup language that allows documents to be rendered in html and other formats (e.g. pdf) with a minimum of special formatting. Knitr is a system for dynamic report generation in R. Both should already be installed on your computer.

In general, a project developed with R/Markdown will consist of a markdown document (with extension `.Rmd`) and the compiled report. You should learn a little more about the functions of R/Markdown and knitr, but first we will engage in a little learning-by-doing.

Exercise. Create a basic R Markdown document. Go to File -> New File -> R Markdown. Enter a title and author into the dialog box. Select the desired default output format. Save the resulting, automatically generated file. To compile the document, click on “Knit” in the R Studio editor. Study the code and the resulting report.

You now know enough R programming that you should be able to discern what is commentary and what is code in the `.Rmd` document. First, notice that the `.Rmd` document begins with a header and how it is set off from the rest of the document with three hyphens at the start and end of the header. These are editable properties. Indeed, there are many more options that we could encode in the header, but do not need to bother with now.

Next, observe that *R chunks* are set off from commentary by three back ticks and curly braces (with the argument “r”). Everything within this section must be properly formatted R code. R code chunks are indicated in the editor window with a different background color. Notice that the third R code chunk includes code to generate a plot. The compiled document shows this code, but also the figure it produces! This is the beauty of statistically literate programming: it keeps the explanation, code, and results all together in a way that can be inspected and used both by the data analyst and by others who come later.

Finally, notice that natural language explanation is interspersed with the code chunks. These chunks don’t require any special designation (like the hyphens for the header or the back ticks for the R code), but are basically just “everything else”. This natural language commentary can nonetheless be “marked up” with special symbols that cause the formatted document to display sections titles, bold face or italics, clickable URLs, etc. In fact, the Markdown markup language enables quite a lot of reasonably sophisticated typesetting.

In closing, we note that we have only scratched the surface of what can be done with R/Markdown and knitr.

Exercise. Visit the R Markdown website and look around. Especially, read through the Authoring Basics.

This document – and all the others used for this workshop – were written in R/Markdown.