

# Research Tools

---

Michael Noonan

Biol 520C: Statistical modelling for biological data

1. Housekeeping
2. R Markdown
3. GitHub
4. Cluster Computing

# Housekeeping

---

- Great job on the talks.
- Term papers are due on Sunday (Dec. 12th), I will be submitting your grades on (Dec. 19th). Any assignments still missing by then will receive a grade of 0.
- Course/Instructor evaluations are due on Dec. 10th. Please take some time to fill them out.

# R Markdown

---

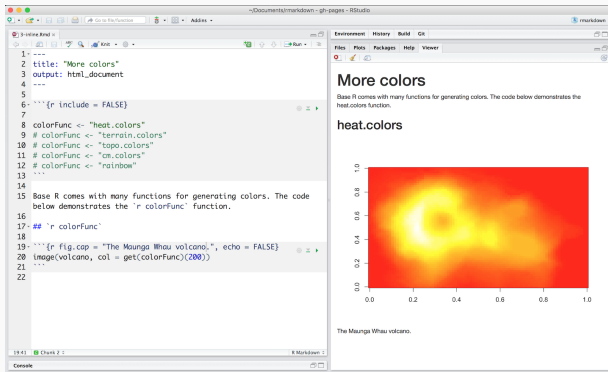
R Markdown allows you to do two main things:

- save and execute code (how we have been using it in this course);  
and
- generate high quality reports that can be shared with an audience.

Importantly, you can do both of these things in a single R Markdown file.

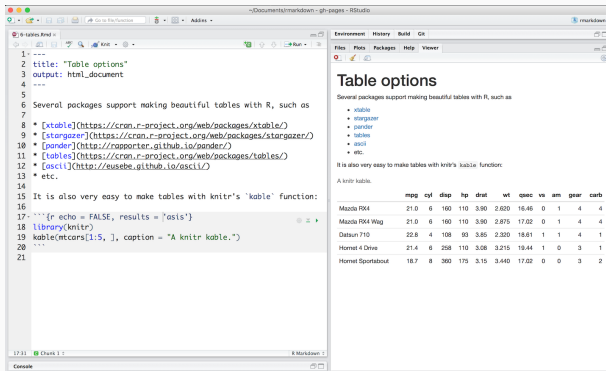
This allows for easier reproducibility, since both the computing code and narratives are in the same document, and results are automatically generated from the source code/data.

A useful feature of R Markdown is that code results can be inserted directly into the text of a .Rmd file.



This makes updating your documents easier, and prevents transcription errors.

Figures and tables can also be generated directly from your data and places into your document.



The screenshot shows the RStudio interface. The left pane contains R code for generating a table using the `kable` function. The right pane shows the rendered HTML document, which includes a title, a list of packages, and a table of car data.

```

1 ---
2 title: "Table options"
3 output: html_document
4 ---
5
6 Several packages support making beautiful tables with R, such as
7
8 * [xtable](https://cran.r-project.org/web/packages/xtable/)
9 * [stargazer](https://cran.r-project.org/web/packages/stargazer/)
10 * [pander](http://rapporter.github.io/pander/)
11 * [tables](https://cran.r-project.org/web/packages/tables/)
12 * [asci](http://eusebe.github.io/asci/)
13 * etc.
14
15 It is also very easy to make tables with knitr's 'kable' function:
16
17 ```{r echo = FALSE, results = 'asis'}
18 library(knitr)
19 kable(mtcars[1:5, ], caption = "A knitr kable.")
20 ```
21

```

Table options

Several packages support making beautiful tables with R, such as

- xtable
- stargazer
- pander
- tables
- asci
- etc.

It is also very easy to make tables with knitr's `kable` function:

A knitr kable.

	mpg	cyl	displ	hp	drat	wt	qsec	vs	am	gear	carb
Mazda RX4	21.0	6	160	110	3.90	2.620	16.46	0	1	4	4
Mazda RX4 Wag	21.0	6	160	110	3.90	2.875	17.02	0	1	4	4
Datsun 710	22.8	4	108	93	3.85	2.320	18.61	1	1	4	1
Hornet 4 Drive	21.4	6	258	110	3.08	3.215	19.44	1	0	3	1
Hornet Sportabout	18.7	8	360	175	3.15	3.440	17.02	0	0	3	2



Just like using `#` in R scripts to add comments, you can wrap text in the following: `<!-- -->` to add comments to your narratives that won't show up in the knit document.

R Markdown is free, open source, and allows you to house all of your work (data cleaning and analysis, figures, tables, writing) in a single document.

It has built in spellchecking, commenting, and reference formatting and renders high quality reports on par with many other word processing software (e.g. Microsoft word, Pages,  $\text{\LaTeX}$ ).

Allows you to generate publication quality documents that are robust to issues like transcription errors and file disorganisation, ensuring reproducibility.

Also allows you to build web pages (e.g., the 520C course website was built entirely using R Markdown).

If you ever have any questions about how to do something in R Markdown, the R Markdown Definitive Guide is a good place to look:

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

You can also access the cheatsheet from the help tab in R Studio.

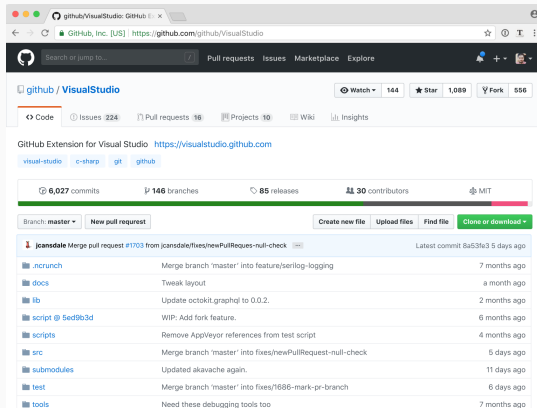
**GitHub**

---

GitHub is a web-based interface that provides access to open source version control software that lets multiple people make separate changes to projects at the same time.



GitHub is centred around repositories (or repos). A repo is a folder in which all files/folders associated with your project and their version histories are stored.



github / VisualStudio

6,027 commits 146 branches 85 releases 30 contributors MIT

Branch: master New pull request

Create new file Upload files Find file Clone or download

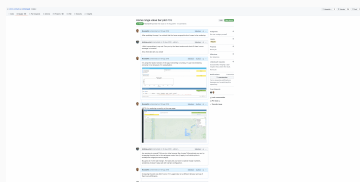
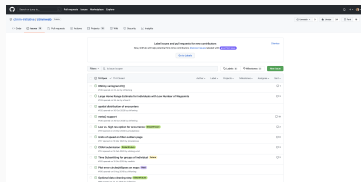
jcansdale	Merge pull request #1703 from jcansdale/fixes/newPullRequest-null-check	Latest commit 8a53fa3 5 days ago
.ncrunch	Merge branch 'master' into feature/serilog-logging	7 months ago
docs	Tweak layout	a month ago
lib	Update octokit.graphql to 0.0.2.	2 months ago
script @ 5ed9b3d	WIP: Add fork feature.	6 months ago
scripts	Remove Appveyor references from test script	4 months ago
src	Merge branch 'master' into fixes/newPullRequest-null-check	5 days ago
submodules	Updated akavache again.	11 days ago
test	Merge branch 'master' into fixes/1686-mark-pr-branch	6 days ago
tools	Need these debugging tools too	7 months ago

Working in GitHub involves pushing, pulling, merging, and committing changes to branches.

- **Branch** — a workspace in which you can make changes (can be the main branch, or a personal branch).
- **Commit** — a saved record of a change made to a file within the repo.
- **Pull Request (PR)** — the way to ask for changes made to a branch to be merged into another branch that also allows for multiple users to see, discuss and review work being done.
- **Merge** — after a pull request is approved, the commit will be pulled in (or merged) from one branch to another.

GitHub also provides a place for team members to discuss 'Issues' that need addressing.

In essence the issues feature allow people to identify new tasks that need to be tackled, and to track progress on the task from beginning to end.





GitHub repositories can always be changed or deleted (meaning they do not function as permanently stable archives).

GitHub repositories have limited data storage capacity (2Gb), not suitable for some types of research.

GitHub repositories can be public, which can make them inappropriate for sensitive data.

GitHub is free, open source, and allows you to house all of your work (data cleaning and analysis, figures, tables, writing) in a single place.

It has built in version control, meaning you can track your changes and recover files (also protects you from computer failures).

A repository can be public or private, giving you full control over how your work is being shared and who can collaborate on your project.

Also allows you to host web pages (e.g., the 520C course website is hosted in GitHub).

Can be integrated into Gitter, permitting a chat interface.

# Cluster Computing

---

Most standard analyses are fairly quick and can be run on the order of seconds to minutes.

Large datasets, complex models, bootstrapping, bayesian models, and/or simulation based experiments are becoming increasingly common in biological research.

This can easily push the capacity of standard desktop computers, requiring days to months of computation time.

Cluster computing lets you spread your calculations out across multiple, linked computers (nodes).

Leaning on a cluster can dramatically cut down on run times (e.g., from months to days, or days to hours).

As UBC students you have access to Sockeye and Compute Canada (you also have access to Chinook for secure data storage).

If you think you will need to use a cluster as part of your research, UBC's Advanced Research Computing team have a series of webinars on how to use Sockeye: <https://osf.io/wpcg6/>, and there is a good wiki for Compute Canada.