

R-projects

- Everything in its right place



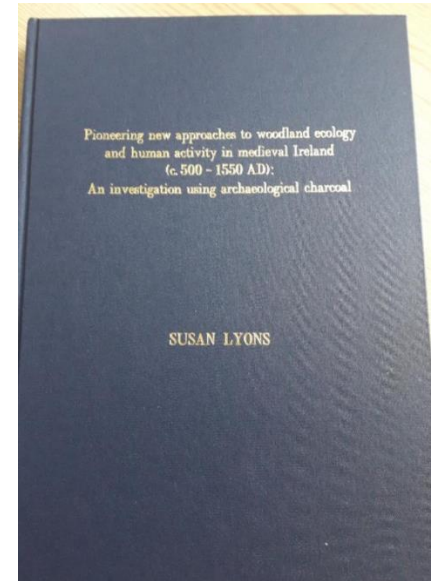
Brendan Palmer,

Statistics & Data Analysis Unit,

Clinical Research Facility - Cork

How is research presented?

Theses



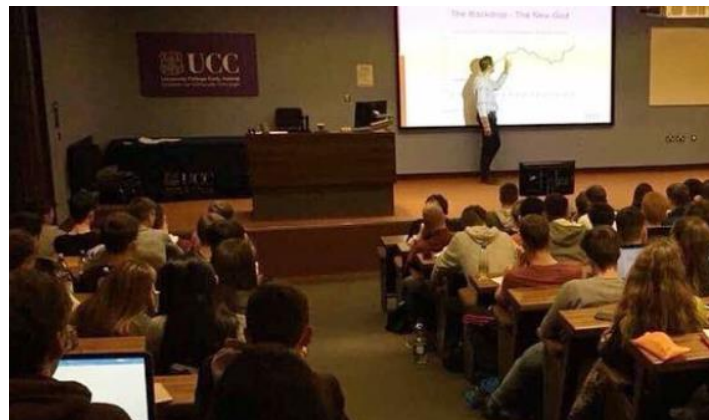
Books



Posters



Talks



Papers

Journal of Virology

Network Analysis of the Chronic Hepatitis C Virome Defines Hypervariable Region 1 Evolutionary Phenotypes in the Context of Humoral Immune Responses

Brendan A. Palmer,* Daniel Schmidt-Martin,* Zoya Dimitrova,* Pavel Skums,* Orla Crosbie,* Elizabeth Kenny-Walsh,* Liam J. Fanning*

Molecular Biology Diagnostic & Research Laboratory, Department of Medicine, University College Cork, Ireland† Division of Viral Hepatitis, Centers for Disease Control and Prevention, Atlanta, Georgia, USA* Department of Hepatology, Cork University Hospital, Cork, Ireland†

ABSTRACT
Hypervariable region 1 (HVR1) of hepatitis C virus (HCV) comprises the first 27 N-terminal amino acid residues of E2. It is classically seen as the most heterogeneous region of the HCV genome. In this study, we assessed HVR1 evolution by using ultradeep pyrosequencing for a cohort of treatment-naïve, chronically infected patients over a short, 16-week period. Organization of the sequence set into connected components that represented single nucleotide substitution events revealed a network dominated by highly connected, centrally positioned master sequences. HVR1 phenotypes were observed to be under strong purifying (stationary) and strong positive (antigenic drift) selection pressures, which were coincident with advancing patient age and cirrhosis of the liver. It followed that stationary viromes were dominated by a single HVR1 variant surrounded by minor variants comprised from conservative single amino acid substitution events. We present evidence to suggest that neutralization antibody efficacy was diminished for stationary-virome HVR1 variants. Our results identify the HVR1 network structure during chronic infection as the preferential dominance of a single variant within a narrow sequence space.

IMPORTANCE
HCV infection is often asymptomatic, and chronic infection is generally well established in advance of initial diagnosis and subsequent treatment. HVR1 can undergo rapid sequence evolution during acute infection, and the variant pool is typically seen to diverge away from ancestral sequences as infection progresses from the acute to the chronic phase. In this report, we describe HVR1 viromes in chronically infected patients that are defined by a dominant epitope located centrally within a narrow variant pool. Our findings suggest that weakened humoral immune activity, as a consequence of persistent chronic infection, allows for the acquisition and maintenance of host-specific adaptive mutations at HVR1 that reflect virus fitness.

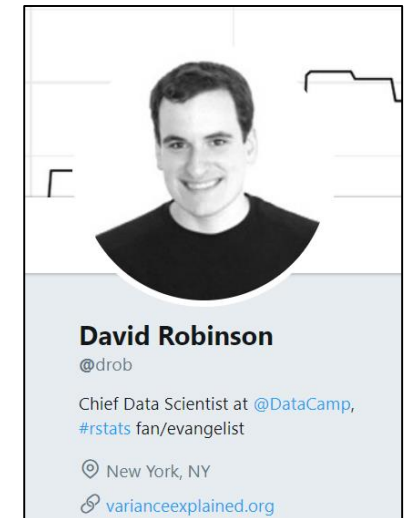
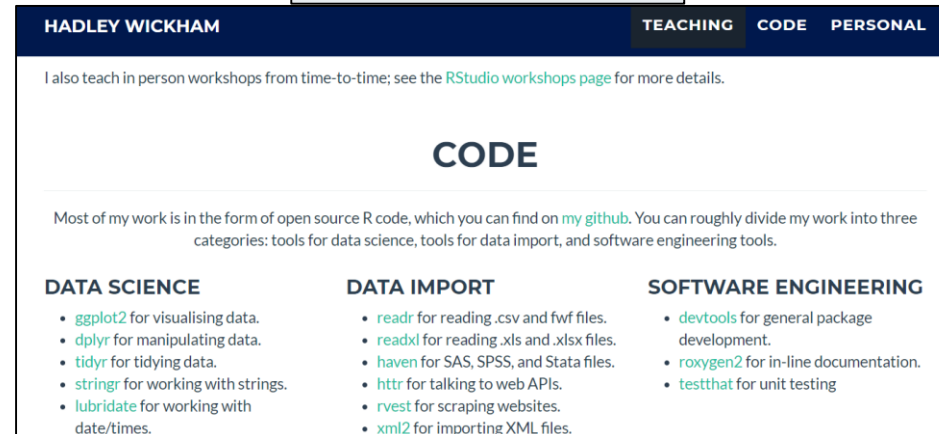
Hepatitis C virus (HCV) infection is a global health issue and is recognized as a major etiological agent of liver-related diseases (1). It has been estimated that the current prevalence of HCV represents approximately 2% of the global adult (15 years of age and older) population (2). Following transmission, HCV infection may remain asymptomatic for decades, resulting in the majority of infections initially passing undetected (3). It is estimated that up to 1 million Americans are living with the virus, the majority of whom became infected prior to the isolation and identification of the virus (4, 5). Consequently, the U.S. Centers for Disease Control and Prevention now recommend that Americans born from 1945 to 1965 be screened for the presence of the virus notwithstanding the presence of clinical symptoms (3, 5). HCV is a single-stranded positive-sense RNA virus of considerable genomic heterogeneity. A recent reclassification defined the HCV global distribution into 7 genotypes and 67 subtypes, with genotypes 1 and 3 accounting for the majority of infections worldwide (6, 7). An error-prone RNA-dependent RNA polymerase, together with an inherent tolerance of defined hypervariable regions (HVR), accounts for much of this variability. Three HVRs are located within the envelope glycoprotein E2. The greatest heterogeneity has been identified at the 27-amino-acid HVR1 (residues 384 to 410 of the H77 reference strain), located at the amino-terminal end of the E2 glycoprotein (8). Recent studies indicated that the central region of E2 (residues 456 to 656) is globular and surprisingly compact, whereas the first 80 amino acids (including

Received 21 November 2015; Accepted 22 December 2015
Accepted manuscript posted online 10 December 2015
Citation: Palmer BA, Schmidt-Martin D, Dimitrova Z, Skums P, Crosbie O, Kenny-Walsh E, Fanning LJ (2016) Network analysis of the chronic hepatitis C virome defines hypervariable region 1 evolutionary phenotypes in the context of humoral immune responses. J. Virol. 90:3218–3228. doi:10.1128/JVI.02090-15
Editor: M. S. Diamond
Address correspondence to Liam J. Fanning, lfanning@ucc.ie.
BA.P. and D.S.M. contributed equally to this article.
Copyright © 2016, American Society for Microbiology. All Rights Reserved.

But what does it really look like?



Disclaimer



STAT
545

Home FAQ Syllabus Topics People

Data wrangling, exploration, and analysis with R

UBC STAT 545A and 547M

Learn how to

- explore, groom, visualize, and analyze data
- make all of that reproducible, reusable, and shareable
- using R

VARIANCE EXPLAINED ABOUT ME POSTS LEARN R TEXT MINING IN R INTRODUCTION TO EMPIRICAL BAYES



David Robinson

Chief Data Scientist at DataCamp, works in R and Python.

- ✉ Email
- ✉ Twitter
- ✉ Github
- ✉ Stack Overflow

This is the homepage and blog of David Robinson, Chief Data Scientist at DataCamp. For more about me, [see here](#).

Recent Posts

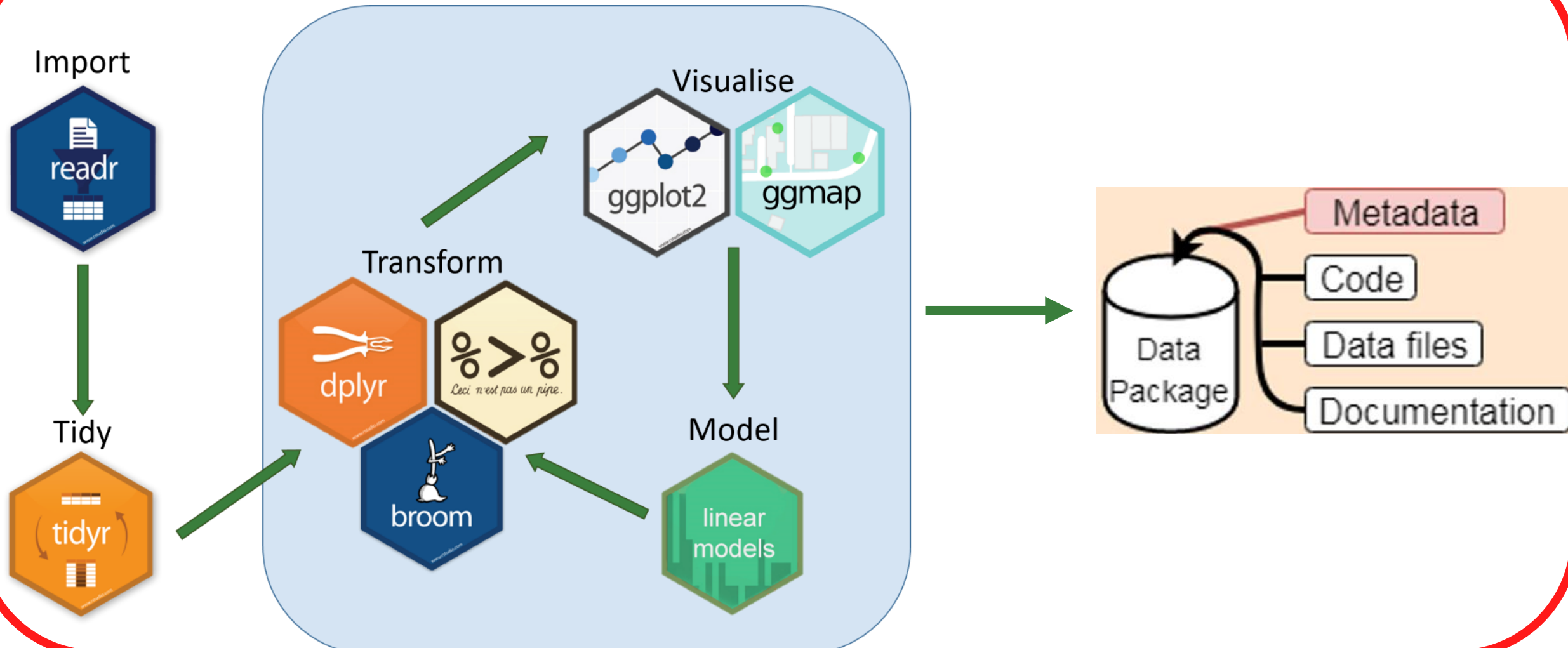
Exploring college major and income: a live data analysis in R October 16, 2018
A live screencast of an exploratory data analysis from the Tidy Tuesday series. This one explores college major and income data from 538.

Who wrote the anti-Trump New York Times op-ed? Using tidytext to find document similarity September 06, 2018
An analysis of an anonymous op-ed in the New York Times, using document similarity metrics to match it to Twitter accounts.

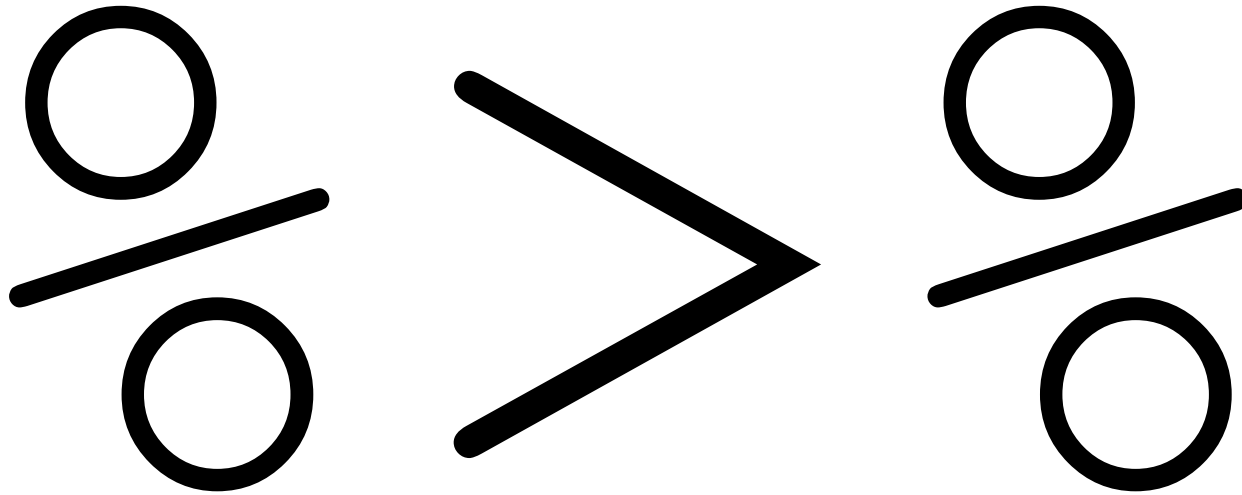
Scientific debt May 10, 2018
Introducing an analogy to 'technical debt' for data scientists.

Putting the pieces together

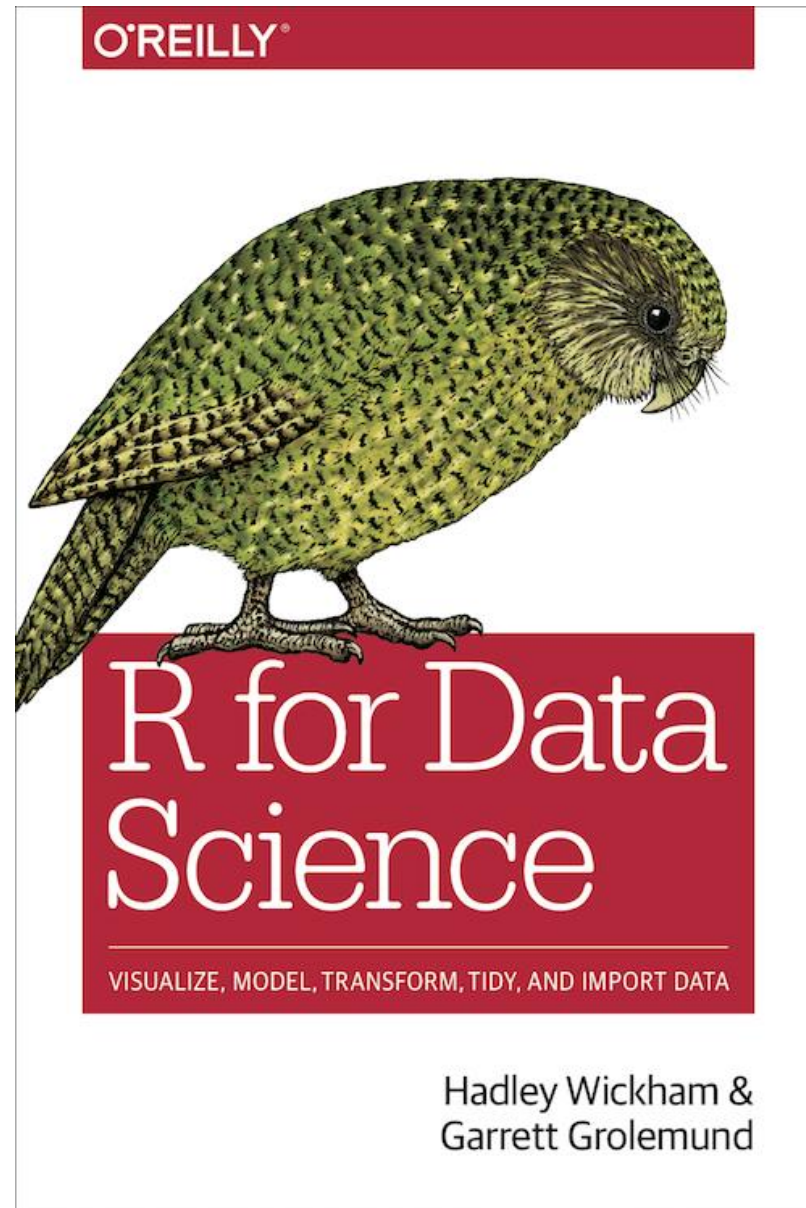
- Data analysis in a tidyverse nutshell



Putting the pieces together



You could write a book on that!!



[R for Data Science webpage](#)

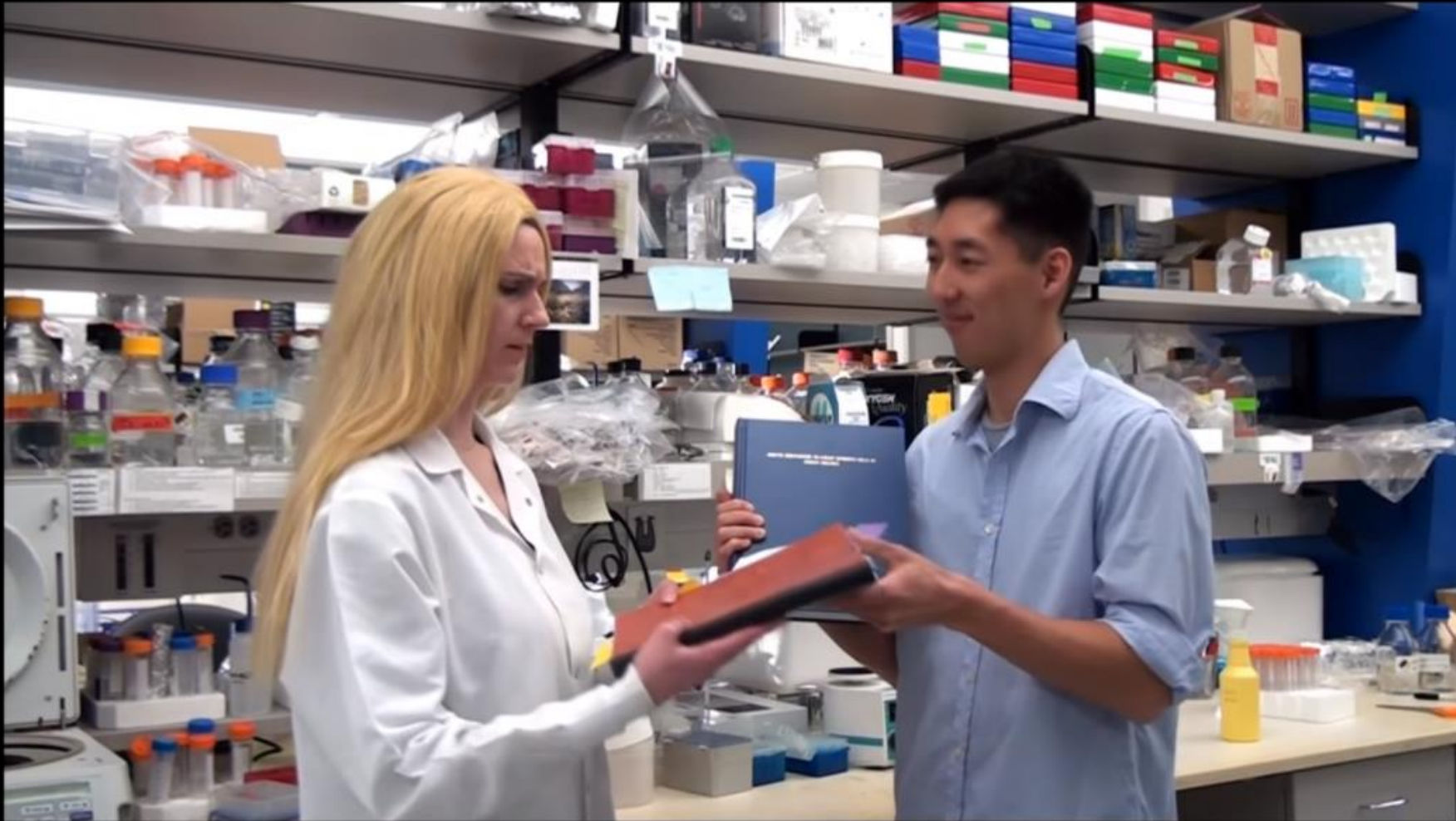
Putting the pieces together

A: Define a project structure

B: Set a naming convention

C: Use scripted workflows

D: Reproducible research



You were defending, one foot out the door

[Link to the video on YouTube](#)



I got your project and its problems galore

[Link to the video on YouTube](#)



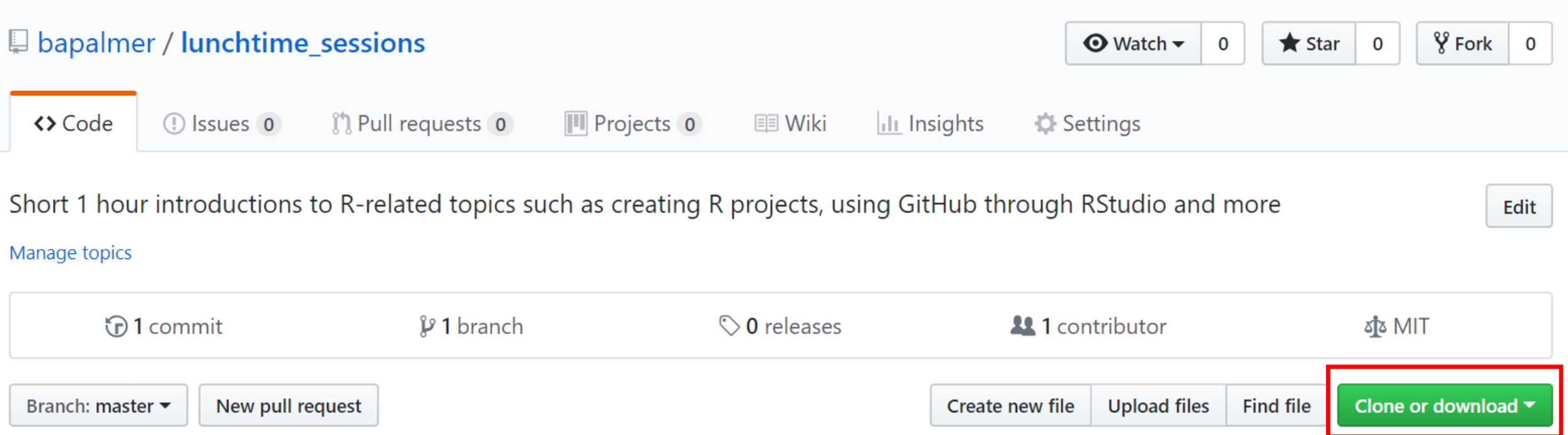
I hate my life,

[Link to the video on YouTube](#)

R projects

- Here's one I made earlier.....

https://github.com/bapalmer/lunchtime_sessions



The screenshot shows the GitHub repository page for 'bapalmer / lunchtime_sessions'. The repository name is at the top left. To the right are buttons for 'Watch', 'Star', and 'Fork', each with a count of 0. Below these are tabs for 'Code', 'Issues', 'Pull requests', 'Projects', 'Wiki', 'Insights', and 'Settings'. The 'Code' tab is selected. Below the tabs is a description: 'Short 1 hour introductions to R-related topics such as creating R projects, using GitHub through RStudio and more'. Below the description is a 'Manage topics' link. At the bottom, there is a bar with repository statistics: '1 commit', '1 branch', '0 releases', '1 contributor', and 'MIT' license. Below this bar are buttons for 'Branch: master', 'New pull request', 'Create new file', 'Upload files', 'Find file', and 'Clone or download'. The 'Clone or download' button is highlighted with a red box and a red arrow pointing to it.

bapalmer / lunchtime_sessions

Watch 0 Star 0 Fork 0

Code Issues 0 Pull requests 0 Projects 0 Wiki Insights Settings

Short 1 hour introductions to R-related topics such as creating R projects, using GitHub through RStudio and more

Manage topics







1 commit 1 branch 0 releases 1 contributor MIT

Branch: master New pull request Create new file Upload files Find file Clone or download

- Open the R project folder and explore it

R projects

Documents > R_Users_Workshop > lunchtime_sessions-master > R-projects

Name ^		Date modified
 data		14/11/2018 23:03
 docs		15/11/2018 00:19
 figs		14/11/2018 22:53
 scripts		14/11/2018 23:48
 tables		15/11/2018 00:19
 R-projects		15/11/2018 00:27

A: Still haven't found what I'm looking for

- Help your future-self

Final Final version

File Home Share View

← → ↕ ↑ > This PC > B_Palmer_Medicine_Files > 4a Project > Pyrosequencing_analysis > Pyrosequencing_Paper > Draft_Paper_incl_Figs > Submission > JVI_Resubmission > JVI_resubmission_files > Final Final version

	Name	Date modified	Type	Size
Quick access	Cover_Letter_B_A_Palmer_Sept_2014	10/09/2014 17:05	Microsoft Word 97 - 200...	559 KB
	Fig_1_Sept_14	11/09/2014 10:31	Adobe Acrobat Docum...	25 KB
	Fig_1_Sept_14	10/09/2014 23:07	Microsoft PowerPoint 9...	158 KB
	Fig_2_Sept_14	11/09/2014 10:31	Adobe Acrobat Docum...	12 KB
	Fig_2_Sept_14	10/09/2014 23:07	Microsoft PowerPoint 9...	212 KB
	Fig_3_Sept_14	11/09/2014 10:31	Adobe Acrobat Docum...	173 KB
	Fig_3_Sept_14	10/09/2014 23:07	Microsoft PowerPoint 9...	527 KB
	Fig_4_Sept_14	11/09/2014 10:31	Adobe Acrobat Docum...	40 KB
	Fig_4_Sept_14	10/09/2014 23:07	Microsoft PowerPoint 9...	342 KB
	Fig_5_Sept_14	11/09/2014 10:33	Adobe Acrobat Docum...	12 KB
	Fig_5_Sept_14	10/09/2014 23:07	Microsoft PowerPoint 9...	178 KB
	HCV_UDPS_B_A_Palmer_Sept_14	17/09/2014 12:21	Microsoft Word 97 - 200...	442 KB
	Response_to_Reviewer_Sept_14	10/09/2014 22:42	Microsoft Word Docum...	559 KB
	Supplementary_Figure_B_A_Palmer_Sept_14	29/08/2014 13:21	Microsoft Word Docum...	378 KB
	Supplementary_Figure_B_A_Palmer_Sept_14	10/09/2014 22:31	Adobe Acrobat Docum...	224 KB
	Tables_B_A_Palmer_Sept_2014	10/09/2014 22:09	Microsoft Word 97 - 200...	185 KB

Desktop Downloads Documents Pictures Projects Google Drive House Google Drive File Stream (G:) FAIR_workshop Icon Files R_Users_Workshop subgroup_4-2_drafts OneDrive This PC

Define a generic project structure



- STEP 1: Give your research projects a shared structure

File Home Share View			
← → ↕ ↑ This PC > Documents > Projects > generic.project > analysis			
	Name	Date modified	Type
Quick access	.Rproj.user	09/04/2018 20:28	File folder
	data	26/04/2017 06:43	File folder
	docs	26/04/2017 06:43	File folder
	plots	26/04/2017 06:43	File folder
	scripts	09/04/2018 20:28	File folder
	tables	26/04/2017 06:43	File folder
	.Rhistory	22/03/2018 14:09	RHISTORY File
	generic	25/08/2017 15:46	RMD File
	genericProject	22/03/2018 14:05	R Project
	style.1	06/07/2017 13:33	Microsoft Word D...
Desktop			
Downloads			
Documents			
Pictures			
Projects			
Google Drive			
House			
Google Drive File Stream (G:)			
FAIR_workshop			

B: Give your files informative names

- STEP 2: Include metadata in the file names











› This PC › Documents › R_Users_Workshop › lunchtime_sessions-master › R-projects › data

	Name ^	Date modified
	raw_data	14/11/2018 23:02
	2018-11-04_clean_who_tb_data	04/11/2018 15:15













Come back to what you know

- STEP 3: Make you file names machine readable, human readable and work with default ordering

NO

 Epistatic_change
 Epistatic_change_match_discovery
 Epistatic_change_match_discovery_fig_2_point_1
 Epistatic_change_v2
 epistatic_codon_change_tracking
 Epistatic_connection_network
 Heatmap_for_epistatic_syn
 Heatmap1_for_epi_site_co-change
 Heatmap2_for_epi_fdr_adjusted_p-value
 Heatmap2_for_epi_p-value

Yes

Documents > Projects > 18.04.27-WP3_Feeding_Trial > analysis > scripts		
	Name	Date modified
	 01_data_import_and_tidying_master_file	02/10/2018 18:51
	 02_data_import_and_tidying_nutritics_grouped	19/10/2018 19:47
	 03_figures	17/10/2018 16:40
	 04_tables	22/05/2018 12:26
	 <u>05_study_overview</u>	19/10/2018 23:06
	 functions	13/05/2018 23:13

Outline a file naming convention

Machine readable:

- Inherent order
- Avoid spaces
- Avoid punctuation
- Remove case-sensitivity

Human readable:

- Contains info on content
- Avoid spaces
- Avoid punctuation
- Remove case sensitivity

Metadata:

Separate with underscores ("_")

- Avoid punctuation
- Remove case-sensitivity

01_marshal-data.r

02_pre-dea-filtering.r

03_dea-with-limma-voom.r

04_explore-dea-results.r

90_limma-model-term-name-fiasco.r

helper01_load-counts.r

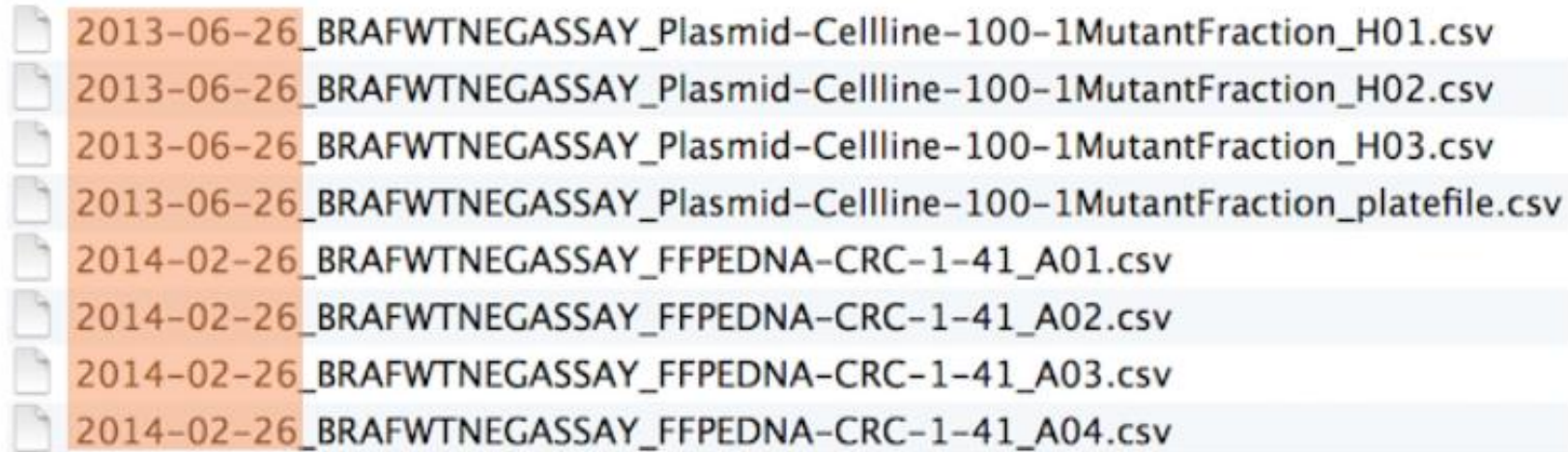
helper02_load-exp-des.r

helper03_load-focus-statinf.r

helper04_extract-and-tidy.r

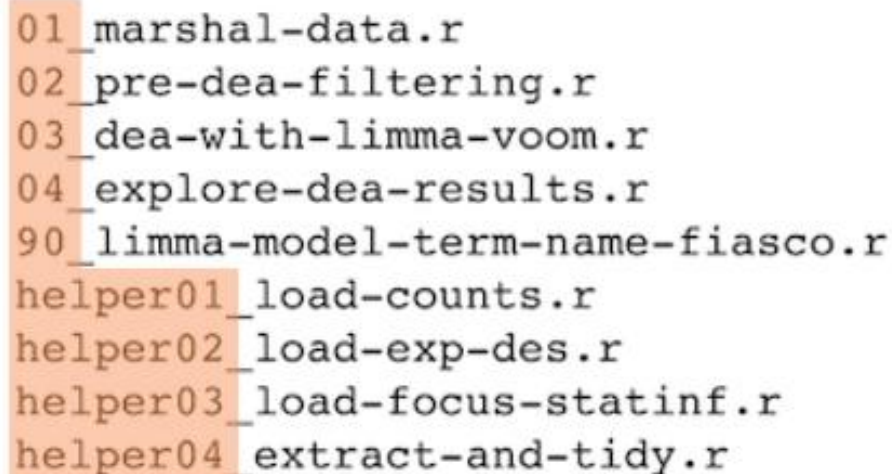
Outline a file naming convention

Chronological order:



```
2013-06-26_BRAFWTNEGASSAY_Plasmid-Cellline-100-1MutantFraction_H01.csv
2013-06-26_BRAFWTNEGASSAY_Plasmid-Cellline-100-1MutantFraction_H02.csv
2013-06-26_BRAFWTNEGASSAY_Plasmid-Cellline-100-1MutantFraction_H03.csv
2013-06-26_BRAFWTNEGASSAY_Plasmid-Cellline-100-1MutantFraction_platefile.csv
2014-02-26_BRAFWTNEGASSAY_FFPEDNA-CRC-1-41_A01.csv
2014-02-26_BRAFWTNEGASSAY_FFPEDNA-CRC-1-41_A02.csv
2014-02-26_BRAFWTNEGASSAY_FFPEDNA-CRC-1-41_A03.csv
2014-02-26_BRAFWTNEGASSAY_FFPEDNA-CRC-1-41_A04.csv
```

Logical order:



```
01_marshall-data.r
02_pre-dea-filtering.r
03_dea-with-limma-voom.r
04_explore-dea-results.r
90_limma-model-term-name-fiasco.r
helper01_load-counts.r
helper02_load-exp-des.r
helper03_load-focus-statinf.r
helper04_extract-and-tidy.r
```

Yesterday

- Open the script 01_bad_habits.R

The screenshot shows the RStudio interface. The script editor on the left contains the following code:

```
20
21 sites1<-as.list(unique(RL6.7$Var1))
22 sites2<-as.list(unique(RL6.7$Var2))
23
24 sites<-as.data.frame(t(merge(sites1,sites2)))
25 colnames(sites)[1]<-"Position"
26
27 for(i in 1:nrow(sites)){
28   ans<-(sites$Position[i]<=65)
29   sites$E1[i]<-ans
30 }
31
32 # Start building network
33 RL6.7_topology<-subset(RL6.7[2:3])
34 g2<-graph.data.frame(RL6.7_topology,vertices=sites,directed=FALSE)
35 g<-simplify(g2)
36 V(g)$color<-ifelse(V(g)$E1==TRUE,"white","grey")
37 V(g)$color<-ifelse(V(g)$E1==TRUE,"white","grey")
```

The Environment pane on the right shows a cluttered list of objects:


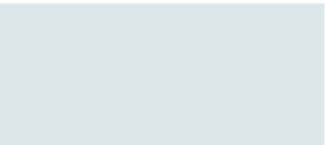
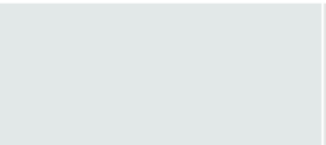





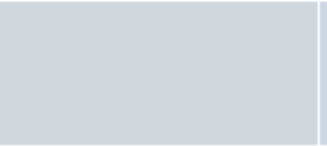
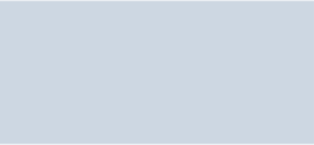




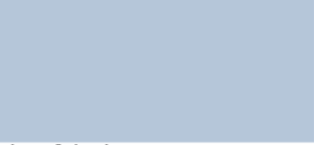




















Name	Type	Length	Size	Value
ans	logical	1	56 B	TRUE
df	tbl_df	4	5.4 KB	39 obs. of 4 variables
g	igraph	10	2.5 KB	List of 10
g2	igraph	10	2.1 KB	List of 10
i	integer	1	56 B	4L
RL1.2	tbl_df	4	1.3 KB	3 obs. of 4 variables
RL2.3	tbl_df	4	1.3 KB	3 obs. of 4 variables
RL3.4	tbl_df	4	1.3 KB	1 obs. of 4 variables
RL4.5	tbl_df	4	1.6 KB	9 obs. of 4 variables
RL5.6	tbl_df	4	1.8 KB	20 obs. of 4 variables
RL6.7	tbl_df	4	1.3 KB	2 obs. of 4 variables
RL6.7_topolo...	tbl_df	2	1000 B	2 obs. of 2 variables
sites	data.frame	2	1.1 KB	4 obs. of 2 variables
sites1	list	2	176 B	List of 2
sites2	list	2	176 B	List of 2

Two red arrows point from the text below to the script and environment panes.

Lack of annotation
Poor naming conventions
Poor readability
Spacing absent

Cluttered environment
Intermediate objects

Is too much choice good or bad?

				
Blue Horizon SW 6497	Sky High SW 6504	Snowdrop SW 6511	Ski Slope SW 6518	Rarified Air SW 6525
				
Byte Blue SW 6498	Atmospheric SW 6505	Balmy SW 6512	Hinting Blue SW 6519	Icelandic SW 6526
				
Stream SW 6499	Vast Sky SW 6506	Take Five SW 6513	Honest Blue SW 6520	Blissful Blue SW 6527
				
Open Seas SW 6500	Resolute Blue SW 6507	Respite SW 6514	Notable Hue SW 6521	Cosmos SW 6528
				
Manitou Blue SW 6501	Secure Blue SW 6508	Leisure Blue SW 6515	Sporty Blue SW 6522	Scanda SW 6529
				
Loch Blue SW 6502	Georgian Bay SW 6509	Down Pour SW 6516	Denim SW 6523	Revel Blue SW 6530
				
Bosporus SW 6503	Loyal Blue SW 6510	Regatta SW 6517	Cammodore SW 6524	Indigo SW 6531

Inconsistent function names, inconsistent syntax

- R is a very versatile language
 - Sometimes it can be too versatile
 - Do you want to use.....

`Names` or `colnames`

`row.names` or `rownames`

`rowSums` or `rowsum`

`Sys.time`, `system.time`

- Is it written as.....

`newobject` or `new.Object`

`x = 5` or `x <- 5`

`mapping=aes(x,y)` or `mapping = aes(x, y)`

Writing clearer code

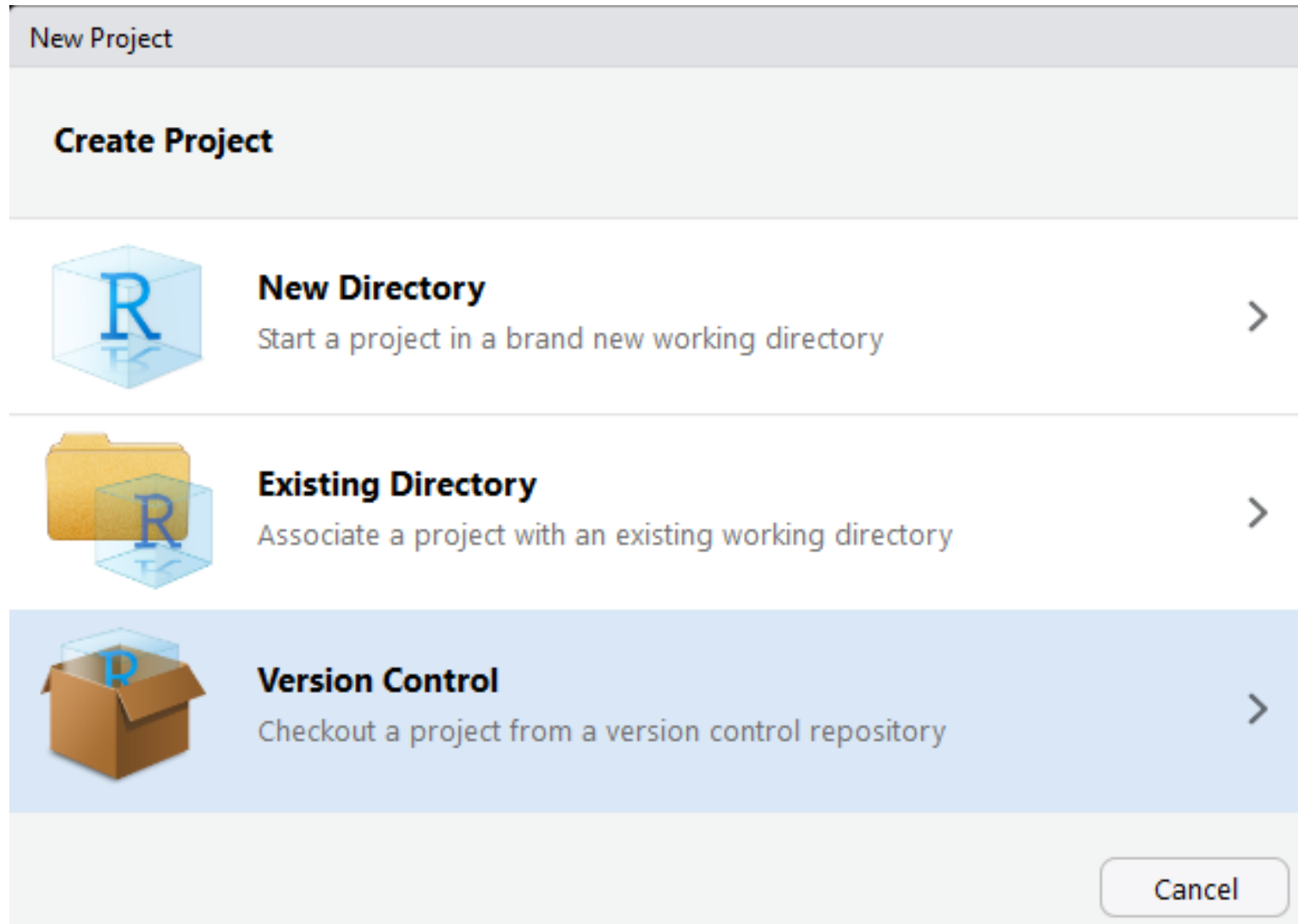
- Annotation
- Object names
 - should use only lowercase letters, numbers, and “_”
- Spacing
 - Put a space before and after =
 - Put a space after a ,
 - Operators should be surrounded by spaces e.g. ==, <-, +
- For a more complete list visit
 - <http://style.tidyverse.org/syntax.html>
- Open the script 02_good_habits.R

C: Joined up thinking

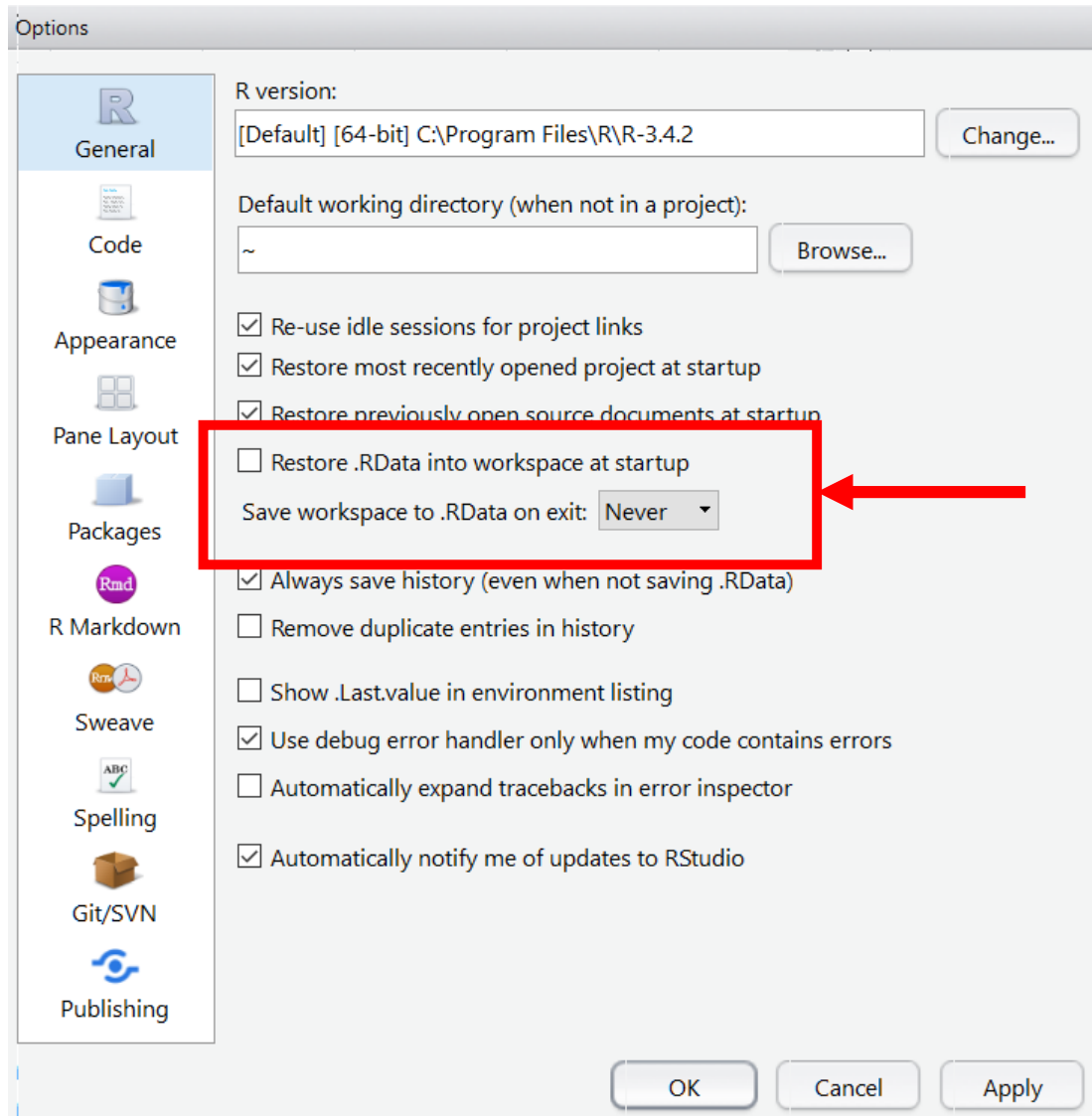
- The R scripts you generate should be human readable
 - Annotate the code
 - Break up the scripts into dedicated tasks
 - Interlink with other within project scripts

```
1 # Data ----
2 # Eight tibbles returned from the 01_data_import_and_tidying_master_file.R
3 # 1. fgf23_data => FGF23 readings from study centres 01-03
4 # 2. food_level_data => Food diary entries
5 # 3. grouped_data => Dialysis and nondialysis diary entries by component
6 # 4. k_data => Serum potassium
7 # 5. master_data_clean => all the clean master file data if required
8 # 6. p_data => Serum phosphate
9 # 7. pth_data => Parathyroid hormone readings
10 # 8. pulses_nuts_data
11
12 source("scripts/01_data_import_and_tidying_master_file.R")
```


D: Setting up an R project



Other points to note



- You might consider your environment as "real"
- If you continue to use R, it is better for you to consider your R scripts as "real", as these should recreate the environment
- You may suffer short term pain
- This will prevent long term agony

Don't Do What Donny Don't Does!!



Donny Don't:

- Start your script with...
`setwd()`

Donny Don't:

- Start your script with...
`rm(list = ls())`

Everything in its right place

- benefits of using R projects for data analysis tasks

▸ This PC ▸ Documents ▸ R_Users_Workshop ▸ lunchtime_sessions-master ▸ R-projects

	Name	Date modified
★	data	14/11/2018 23:03
★	docs	15/11/2018 00:19
★	figs	14/11/2018 22:53
★	scripts	14/11/2018 23:48
★	tables	15/11/2018 00:19
★	R-projects	15/11/2018 00:27

- Open the script 03_data_cleaning.R
- Open the script 04_plots.R
- Open the script 05_tables.R
- Open the script 06_analysis.R

Meetup

Cork (Ireland) R-Users Group

