# R-eproducible workflows

**1-day workshop**
**Afternoon overview**
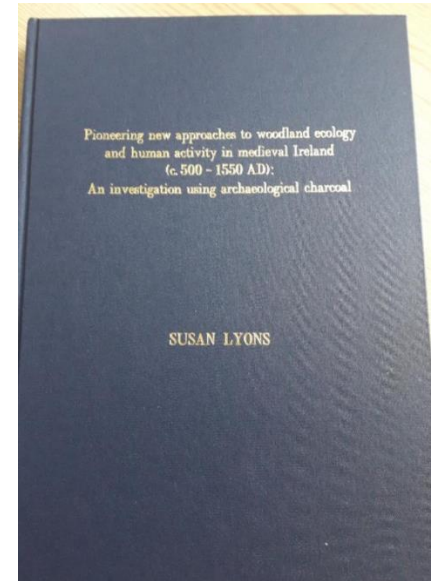
**Brendan Palmer,**

**Statistics & Data Analysis Unit,**

**Clinical Research Facility - Cork**

University College Cork, Ireland
Coláiste na hOllscoile Corcaigh

CRF-C
HRB Clinical Research Facility Cork

HRB
Health Research Board

# How is research presented?

**Papers**



**Books**



**Theses**



**Posters**



**Talks**

# But what does it look like under the bonnet?

You were defending, one foot out the door

I got your project and its problems galore

I hate my life,

# Disclaimer

- Jenny Bryan is the go-to resource for most of the content in this session

# The explosion of data in the life sciences



**Stevens et. al., 2015, Plos Biology**

# This session

- **Project structure**

  - **Naming conventions**

    - **Scripted workflows**

      - **RMarkdown reports**

        - **Reproducible research**

THIS PERSON IS likely to be YOU BTW!!

# Still haven't found what I'm looking for

- Help your future-self

# Define a generic project structure

- STEP 1: Give your research projects a shared structure

# Give your files informative names

- STEP 1: Give your research projects a shared structure

# Everything in its right place

- STEP 2: Make you file names machine readable, human readable and work
with default ordering

## NO

Name

- All unique 4a amino acid Sequences (B-N).fas
- All unique 4a amino acid Sequences (B-N).meg
- All_AA_haplotypes.meg
- All_AA_haplotypes_with_clonal_sequences.meg
- BS100_AA_with_clones
- BS100_AA_with_clones.nwk
- BS1000_AA_pyro&clones
- BS1000_AA_pyro&clones.nwk
- BS1000_AA_pyro_only
- BS1000_AA_pyro_only.nwk
- BS1000_Unique_Clonal_AA
- BS1000_Unique_Clonal_AA.nwk
- BS1000_Unique_Pyro_AA
- BS1000_Unique_Pyro_AA.nwk
- pic

## 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 |
| 05_study_overview | 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
```

datacarpentry.org/rr-organization1/01-file-naming/

# 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_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
```
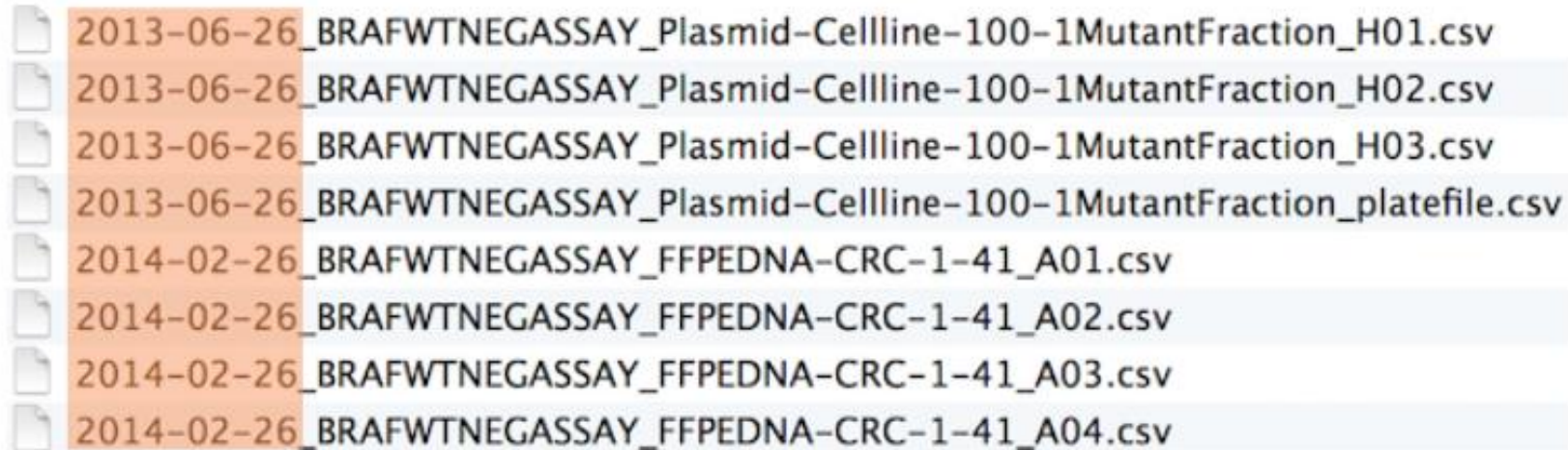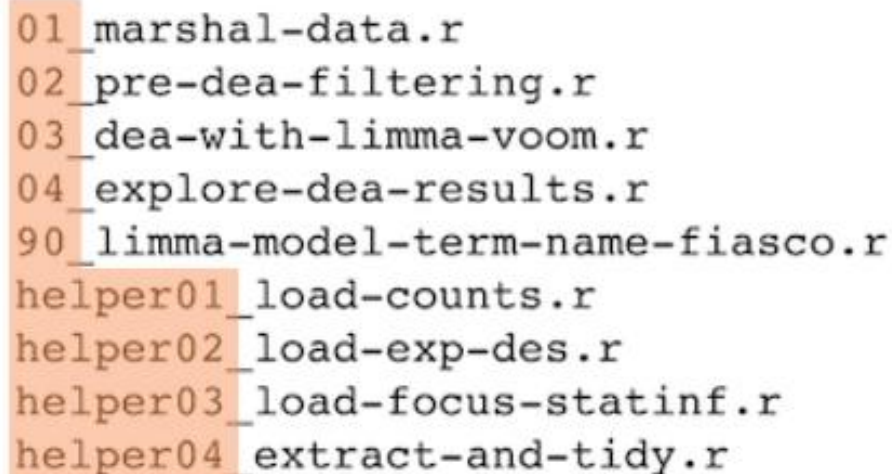
- Reopen ~/morning_session/data
  - Anything we can improve on here?

datacarpentry.org/rr-organization1/01-file-naming/

# 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

```r
 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")
```

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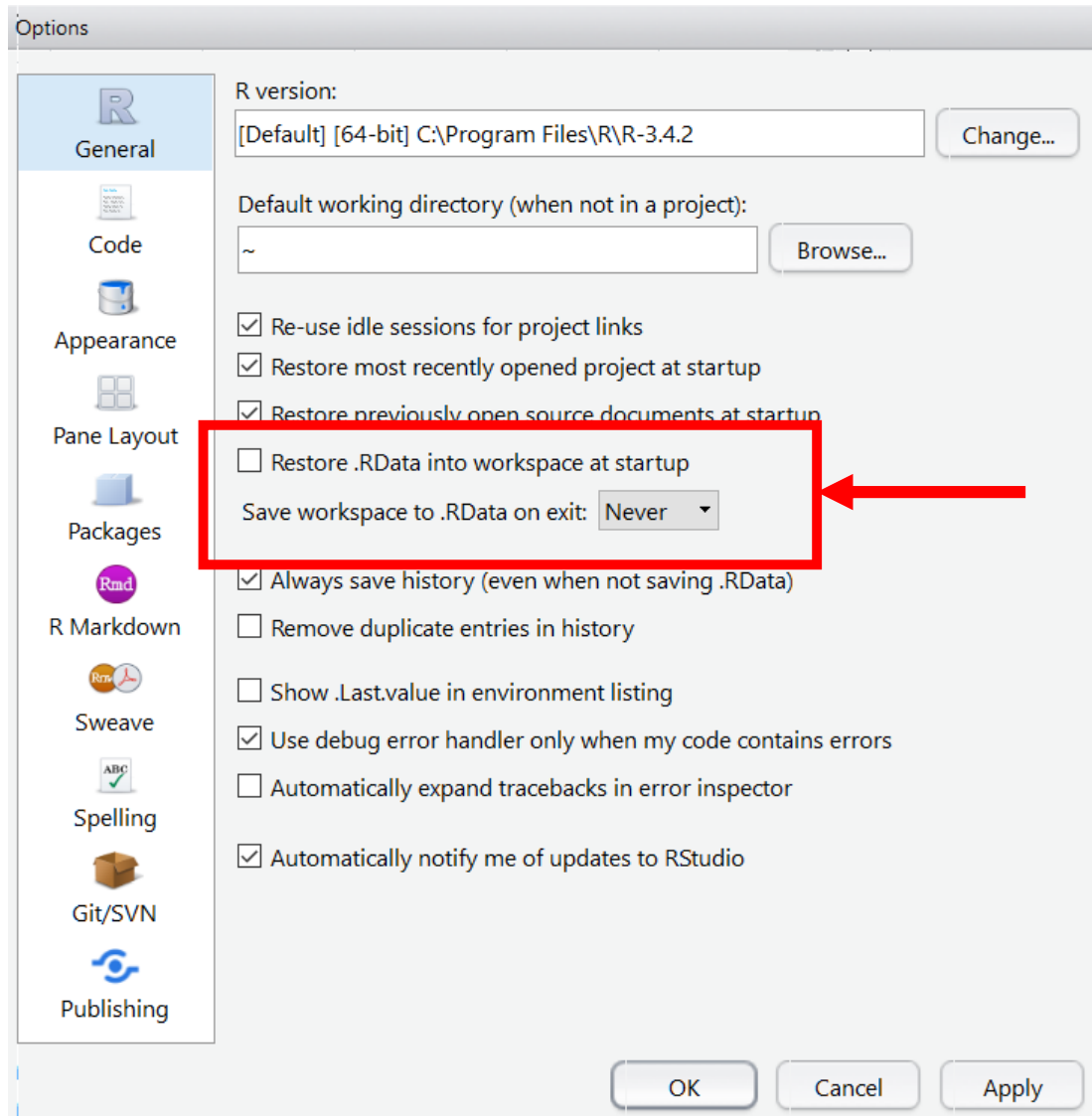


**Donny Don't:**
- Start your script with… setwd()


**Donny Don't:**
- Start your script with… rm(list = ls())
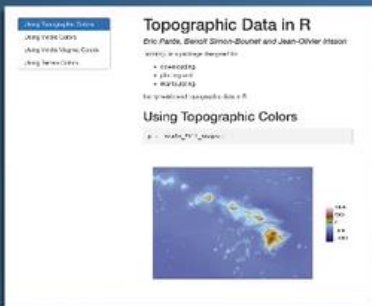
# Other points to note



- 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
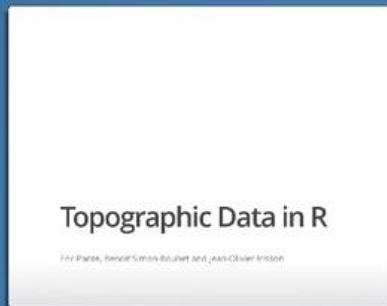
- This will prevent long term agony

# R Markdown

- R Markdown combines the code you wrote, the output produced
  and you own comments

- You can view it as a digital lab notebook, where you are
  both recording what you're doing, and what you were thinking
  while you were thinking it!

- R Markdown outputs can take many forms
          - Word documents, PDFs, slideshows etc.

- Once created the .Rmd file get sent to knitr, which executes
  the chunks of code and creates a new markdown document (.md)
          - this is then processed by pandoc which creates the
            finished file
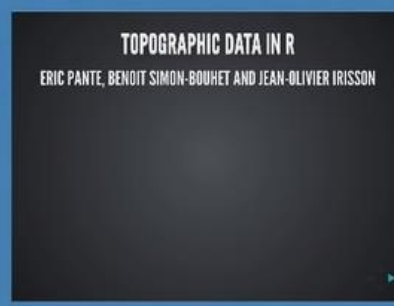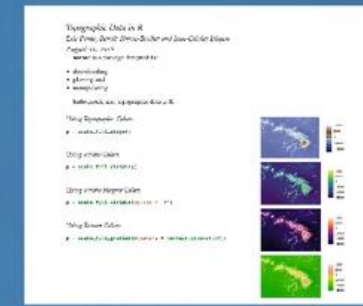              - knitr and pandoc are external websites
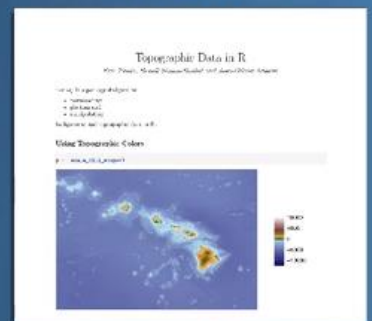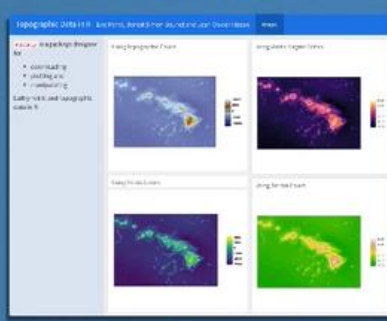
# What has R Markdown ever done for us?
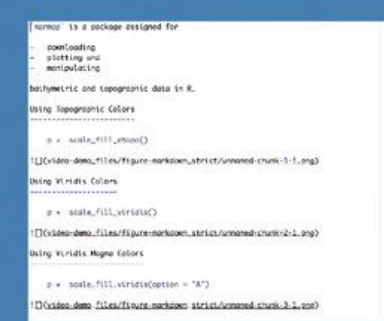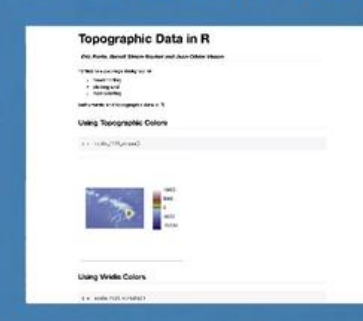


html

ioslides
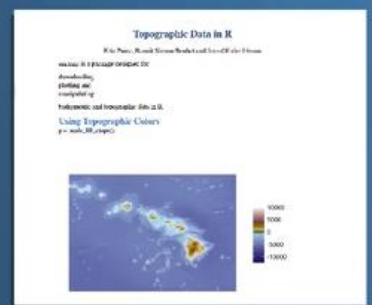
reveal.js

rtf

tufte handout
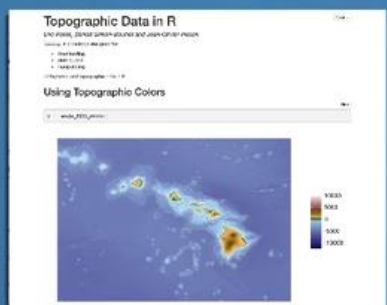
book

pdf

dashboard

slidy
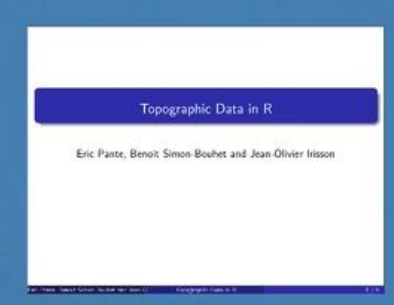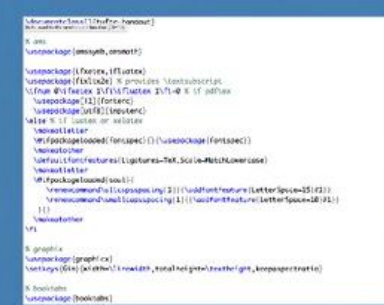
markdown

package vignette

website

Word

notebook

beamer

latex

custom template

shiny app

# R Markdown

YAML header

```
---
title: "Diamond sizes"
date: 2016-08-25
output: html_document
---
```

Chunks of code

````
```{r setup, include = FALSE}
library(ggplot2)
library(dplyr)
smaller <- diamonds %>%
filter(carat <= 2.5)
```
````

Plain text with integrated outputs from R

We have data about `r nrow(diamonds)` diamonds. Only
`r nrow(diamonds) - nrow(smaller)` are larger than
2.5 carats. The distribution of the remainder is shown below:

Chunks of code

````
```{r, echo = FALSE}
smaller %>%
ggplot(aes(carat)) +
geom_freqpoly(binwidth = 0.01)
```
````