

# A Reproducible Workflow



using R and GitHub

Henry Partridge | BRITE event | 3 July 2019



# Who am I?



I'm the manager of the [Trafford Data Lab](#). I have an academic background in German philosophy and crime science. I've previously worked at TfL and MMU.  
I've been a cheerleader for [#rstats](#) since 2013.

# What is reproducibility?

**reproducibility** /ri:pri:dju:sə'biliti/ *noun* to obtain the same results using the method and data of the original study

*which is different from ...*

**replication** /rəpli'keɪʃ(ə)n/ *noun* to obtain the same results using the method of the original study and independently collected data

# Why is reproducibility important?

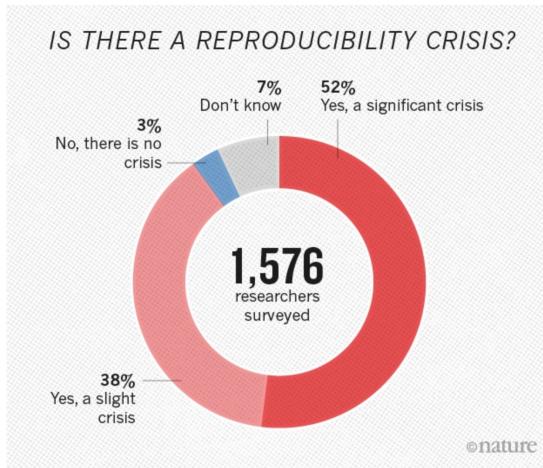
“non-reproducible single occurrences are of no significance to science”

Karl Popper, *The Logic of Scientific Discovery*

# Why is reproducibility important?

- allows checking and double checking by yourself and others
- enables rigorous peer review
- gives confidence in results

# "Reproducibility crisis"



Source: [nature.com](http://nature.com)

100 experimental and correlational studies in psychology were repeated with larger sample sizes. 97% of the original studies had statistically significant results but only 36% of the replications did. The replication effects were on average half the magnitude of the mean effect size of the original effects.

Open Science Collaboration. (2015). Estimating the reproducibility of psychological science. *Science*, 349(6251), aac4716, DOI:10.1126/science.aac4716

# Open science initiatives

- Registered Reports
- study pre-registrations
- many-lab replication projects e.g [ManyLabs](#)
- sharing data
- open access publishing

“ Reproducibility has the potential to serve as a minimum standard for judging scientific claims when full independent replication of a study is not possible. ”

Peng, R. D. (2011). Reproducible research in computational science. *Science*, 334(6060), 1226-1227,  
[DOI:10.1126/science.1213847](https://doi.org/10.1126/science.1213847)

# Practical steps for a reproducible workflow

# Organised projects

- make your project folder self-contained
- quarantine your raw data

```
└── project
    ├── data
    │   ├── raw          # read-only pre-processed datasets
    │   └── processed    # intermediate datasets
    ├── R              # R scripts
    ├── outputs         # tables, charts
    ├── README.md       # project description
    ├── LICENCE.txt
    └── .gitignore
```

# Readable code

- avoid absolute paths
- adopt a consistent style
- comment your code
- write functions

```
# this is an absolute path
df <- read.csv("/Users/henrypartridge/Documents/project/data/foo.csv", string

# this is a relative path
df <- read.csv("data/foo.csv", stringsAsFactors = FALSE)
```

# Literate programming

- avoid word processing software like MS Word
- combine code with human-readable plain text in **R Markdown**

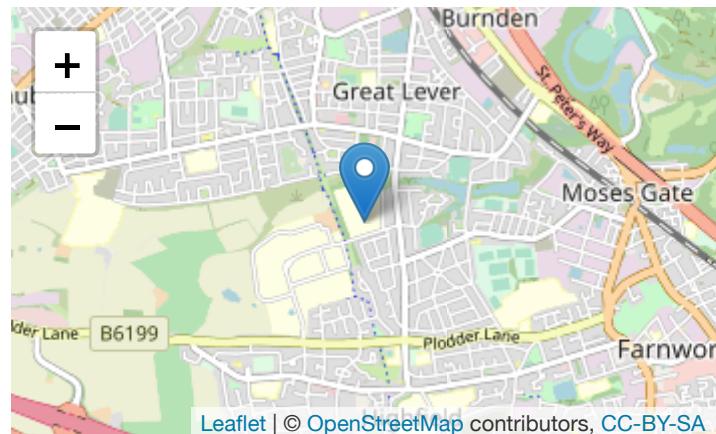
## *R Markdown*

Bolton Science and Technology Centre  
is located on Minerva Road.

```
```{r out.width = '100%',  
fig.height = 3, echo = FALSE}  
leaflet() %>%  
  addTiles() %>%  
  addMarkers(-2.424208, 53.554980,  
             popup = "Bolton Science  
and Technology Centre")  
```
```

## *HTML output*

Bolton Science and Technology  
Centre is located on Minerva Road.



# Version control

- tracks changes to code and plain text files without need for version v0.1 etc.
- timestamps your work
- encourages collaboration
- integrates with RStudio
- remote copies of local projects can be stored on GitHub which also provides issue tracking, wikis and website hosting

Showing 1 changed file with 1 addition and 1 deletion.

Unified Split

View file

2019-03-22\_GMPHIN/R/script.R

```
@@ -2,7 +2,7 @@ library(tidyverse)

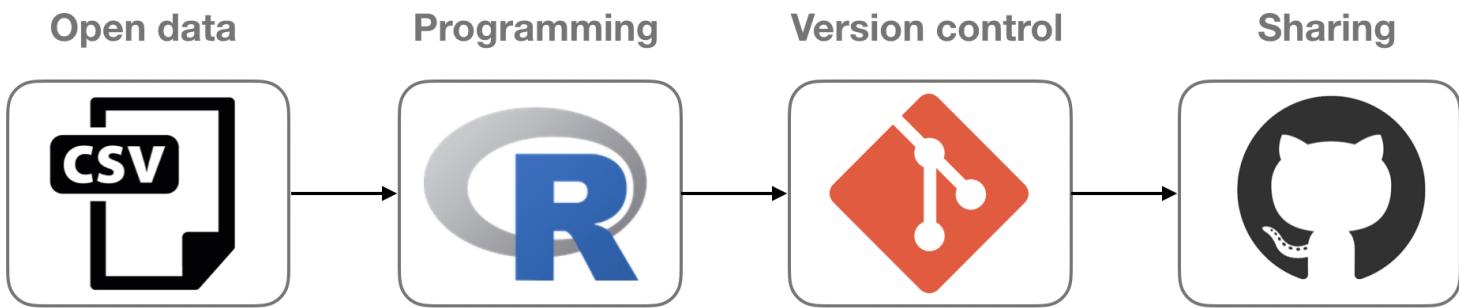
starwars %>%
  select(name, height, mass) %>%
- filter(!is.na(height), !is.na(mass)) %>%
+ drop_na() %>%
  mutate(bmi = mass / (height/100)^2) %>%
  ggplot(aes(x = bmi)) +
  geom_histogram(fill = "#3182bd", colour = "ffffff",
```

# Licensing

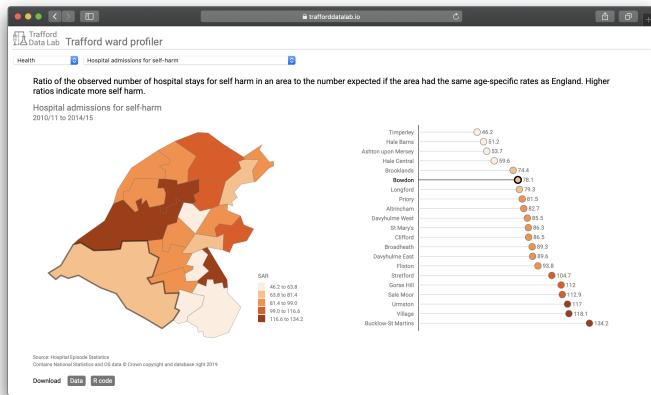
Give people permission to use your data and code:

- Open Government Licence 3.0 for government published data
- CC-BY (Creative Commons Attribution) for media and text
- MIT licence for code

# The Lab's workflow



# Example #1



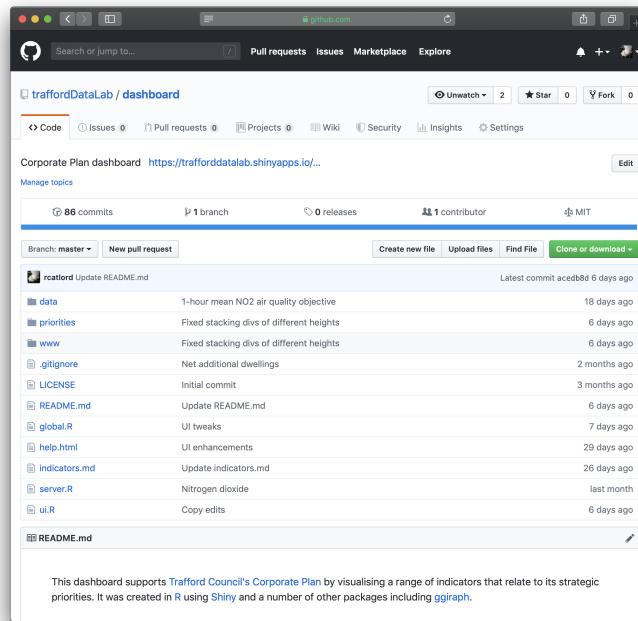
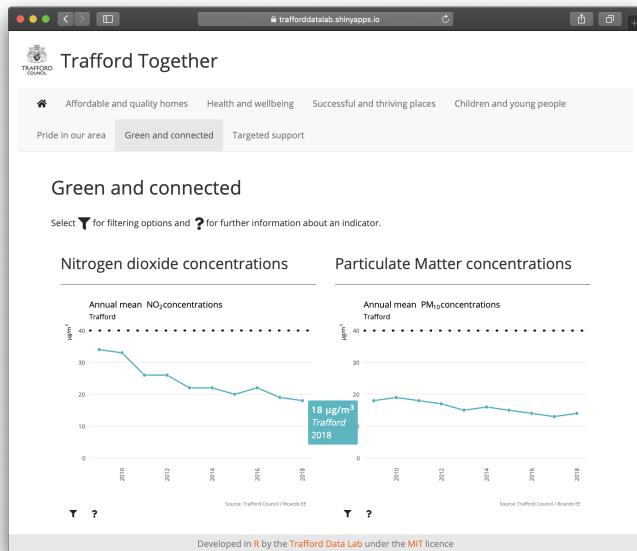
The figure shows a GitHub repository page for 'ward\_data / health / code / hospital\_admissions\_self\_harm.R'. The code is an R script that filters hospital admissions data for Trafford wards based on a specific indicator ID and period, then writes the results to a CSV file.

```
# Health: Hospital admissions for self-harm, 2010/11 - 2014/15 #
# Source: Hospital Episode Statistics
# URL: https://fingertips.phe.org.uk/
# Licence: Open Government Licence

library(tidyverse) ; library(fingertipsR)

df <- fingertips_data(IndicatorID = 92584, AreaTypeID = 8) %>%
  filter(ParentName == "Trafford") %>%
  select(area_code = AreaCode,
         area_name = AreaName,
         value = Value) %>%
  mutate(period = "2010/11 to 2014/15",
        measure = "Hospital admissions for self-harm",
        measure = "SAH",
        unit = "Admissions") %>%
  select(area_code, area_name, indicator, period, measure, unit, value)
write_csv(df, ".../data/hospital_admissions_self_harm.csv")
```

# Example #2



# Example #3

[Edit this page](#)

[back to recipes](#)

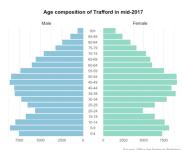
## Population pyramids

A population pyramid visualises the age-gender distribution. This recipe uses the latest mid-year population estimates for a local authority and visualises them by gender and 5-year age bands.

### Ingredients

Data sources Nomis, Office for National Statistics

R packages tidyverse ggplot

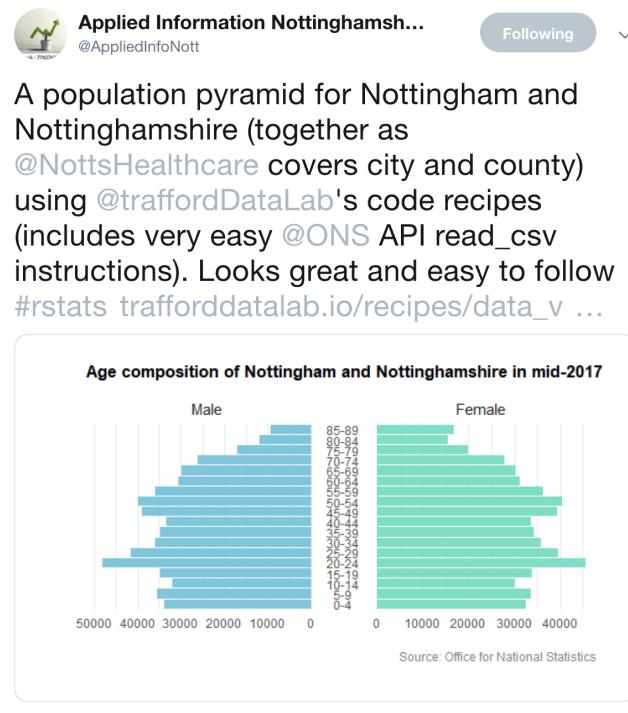


### Instructions

1. Load the necessary R packages.

```
library(tidyverse); library(ggplot)
```
2. Visit [Nomis](#) and select your chosen Geography, Date, Age and Sex. Here we've chosen Trafford (Geography), 2017 (Date), individual ages (Age) and both genders (Sex).
3. Navigate to 'Format / Layout' and choose 'Nomis API' as your format. Punch the 'Download Data' button. Select the 'Tabulation links' tab and right-click the 'Comma Separated Values (csv)' file to obtain the URL path.
4. Paste the URL in the path argument of `read_csv`.

```
df <- read_csv("http://www.nomisweb.co.uk/api/v01/dataset/NM_2002_1.data.csv?geography=80800009&date=latest&genders=M,F")
```



9:42 AM - 15 Jun 2019

1 Like



# take-home message

# GCSE Mathematics Specification (8300/3F)

Paper 3 Foundation tier

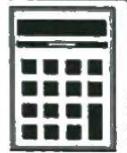
F

Date Morning 1 hour 30 minutes

## Materials

For this paper you must have:

- a calculator
- mathematical instruments.



## Instructions

- Use black ink or black ball-point pen. Draw diagrams in pencil.
- Fill in the boxes at the bottom of this page.
- Answer all questions.
- You must answer the questions in the spaces provided. Do not write outside the box around each page or on blank pages.
- Do all rough work in this book.
- In all calculations, show clearly how you work out your answer.

Source: AQA

# Useful resources

- Broman, K. W., & Woo, K. H. (2018). Data organization in spreadsheets. *The American Statistician*, 72(1), 2-10, DOI:10.1080/00031305.2017.1375989
- Bryan, J. (2018). Happy Git and GitHub for the useR
- Bryan, J. (2017) Project-oriented workflow
- rOpenSci, Reproducibility in Science

# thank you

 @trafforddatalab  
 @trafforddatalab  
 [trafforddatalab.io](http://trafforddatalab.io)

Slides created with **remark.js** and the R package **xaringan**