

R

Greater Manchester Public Health Intelligence Network Session

22 March 2019

Dr Henry Partridge
Trafford Data Lab Manager



Hello



- Studied history, philosophy and crime science
- Data analyst in public sector
- Used R since 2013

Why use R?



ben goldacre

@bengoldacre

Following

I think choice of tool is an important debate in science. I'm always concerned by anyone using point and click for research. Scripts best for reproducibility, transparency, external sense-checking, good audit, and code re-use. And (tenuously...) arguably makes u think more clearly

Andy Field @ProfAndyField

Replies to @LaurentWada @VictorKovalets and 10 others

People get way too snobby about software, it's just a tool to do a job. If you're happy with SPSS use it, if you prefer R use it, love Excel? No problem. Use what works for you (and your collaborators).

1:23 PM - 21 Nov 2017

39 Retweets 108 Likes



19

39

108



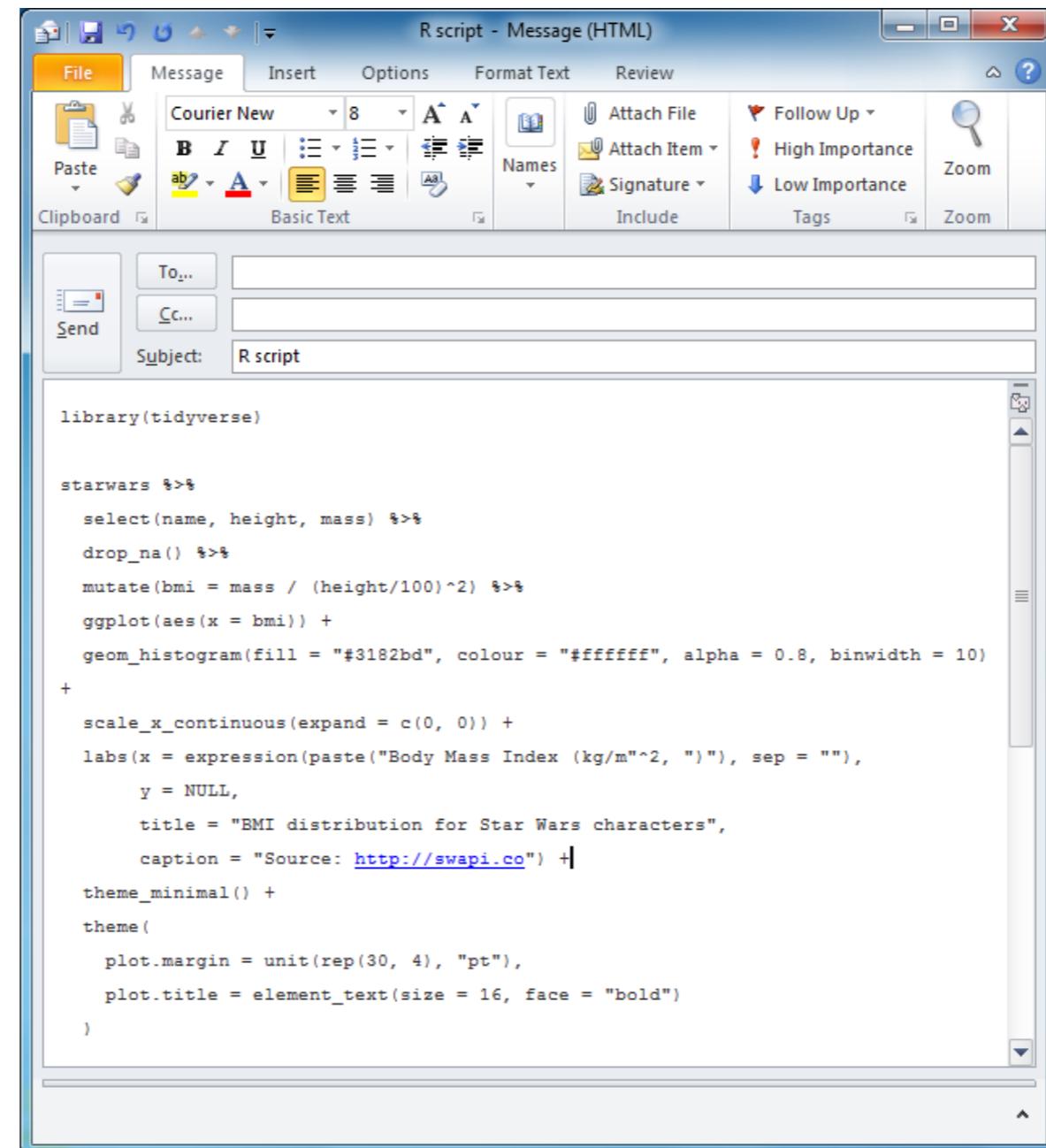
Open



```
script.R
<functions> Help search

1 library(tidyverse)
2
3 starwars %>%
4   select(name, height, mass) %>%
5   drop_na() %>%
6   mutate(bmi = mass / (height / 100)^2) %>%
7   ggplot(aes(x = bmi)) +
8   geom_histogram(fill = "#d9eaf7", colour = "#ffffbf",
9                 alpha = 0.2, binwidth = 10) +
10  scale_x_continuous(expand = c(0, 0))
11  labs(x = expression("Body Mass Index" ~ "kg/m" ^ 2), sep = ""),
12        y = NULL,
13        title = "BMI distribution for Star Wars characters",
14        caption = "Source: Star Wars API")
15  theme_minimal() +
16  theme(
17    plot.margin = unit(c(10, 10, 10, 10), "mm"),
18    plot.title = element_text(size = 10, bold = TRUE))
19
20
21 ggsave("starwars_BMI.png", dpi = 300, scale = 1)
22
```

Shareable



Human readable

```
human_readable.txt  x          script.R          o
1  load the tidyverse R package
2
3  load the starwars dataframe
4  select name, height and mass variables
5  remove rows with missing values
6  calculate the body mass index of each character
7  call a ggplot object with bmi on the x-axis
8  encode the data in a histogram
9  remove padding around zero
10 create an x-axis title
    leave the y-axis title blank
12 add a plot title
13 add a plot caption
14 add a minimalistic theme
15
16 add padding to the plot
17 increase the size of the title
18
19 save the plot as a PNG image in high resolution
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
```

Diffable

Showing 1 changed file with 1 addition and 1 deletion.

Unified Split

2 2019-03-22_GMPHIN/R/script.R

[View file](#)

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

2		2
3	starwars %>%	3 starwars %>%
4	select(name, height, mass) %>%	4 select(name, height, mass) %>%
5	- filter(!is.na(height), !is.na(mass)) %>%	5 + drop_na() %>%
6	mutate(bmi = mass / (height/100)^2) %>%	6 mutate(bmi = mass / (height/100)^2) %>%
7	ggplot(aes(x = bmi)) +	7 ggplot(aes(x = bmi)) +
8	geom_histogram(fill = "#3182bd", colour = "ffffff",	8 geom_histogram(fill = "#3182bd", colour = "ffffff",

@@

Community driven



Manchester R
SUPPORTING THE R COMMUNITY

Supportive

The screenshot shows the Stack Overflow homepage with the following elements:

- Header:** Stack Overflow logo, search bar, navigation icons (User, Flag, Help, Chat), Log In, and Sign Up buttons.
- Left Sidebar (PUBLIC):** Home, Stack Overflow, Tags, Users, Jobs, Teams (Q&A for work) with a Learn More button.
- Section Title:** Top Questions
- Filter Buttons:** Interesting, 386 Featured, Hot, Week, Month.
- Question List:** Four questions are listed:
 - 0 votes, 0 answers, 1 view: "Clean/Build a maven project take too long time" (eclipse, maven) - asked 4 secs ago by Shakti Pravesh 11.
 - 0 votes, 0 answers, 4 views: "Cannot resolve keyword '' into field, but all looks ok" (python, sql, django, orm) - asked 21 secs ago by Tyomik_mnemonic 68.
 - 0 votes, 0 answers, 2 views: "OSX Terminal Pkg Install" (macos, terminal, pkg-file) - asked 22 secs ago by h3tr1ck 301.
 - 0 votes, 0 answers, 2 views: "Python TypeError: derivatives_circ() takes 2 positional arguments but 6 were given" (python, python-3.6, python-3.7) - asked 30 secs ago by Leo 3.

Cutting edge analytics

```
# How many @traffordDataLab followers are bots? #

library(tidyverse) ; library(rtweet) ; library(tweetbotornot)

# retrieve followers of @OpenGovInt
followers <- get_followers("traffordDataLab", n = "all")
followers_info <- lookup_users(followers$user_id) %>%
  select(screen_name, name, followers = followers_count, following = friends_count)

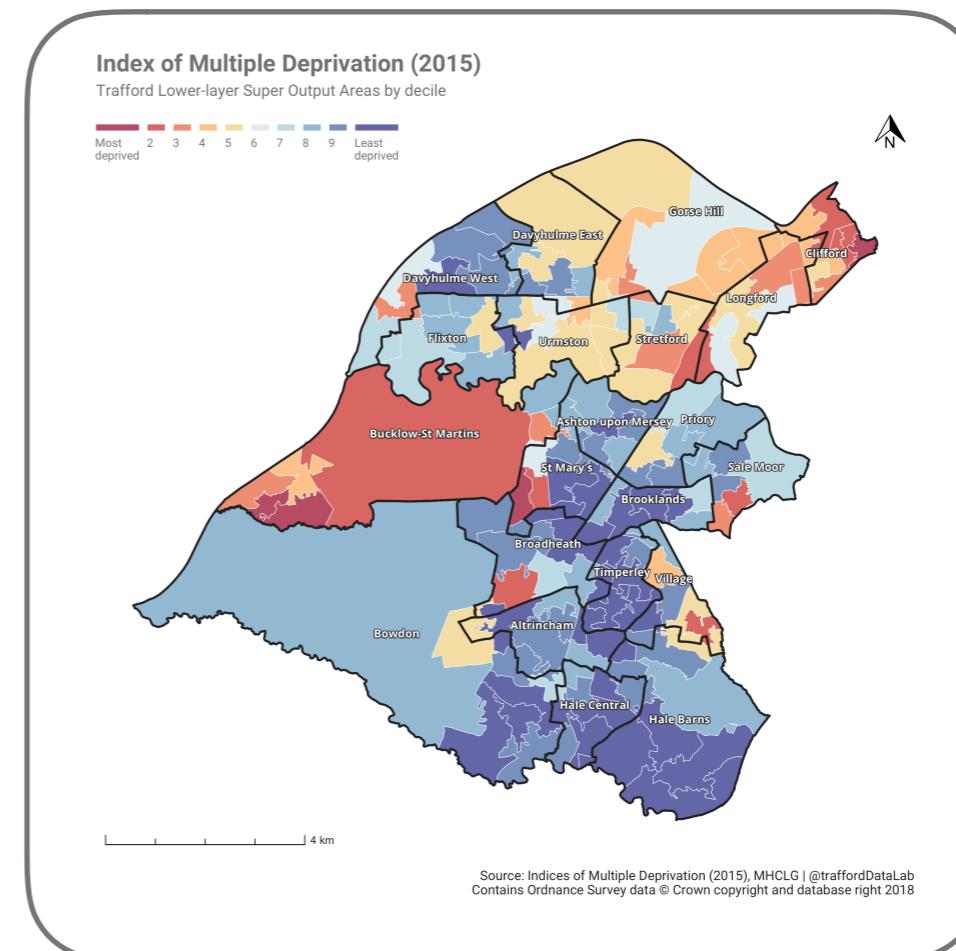
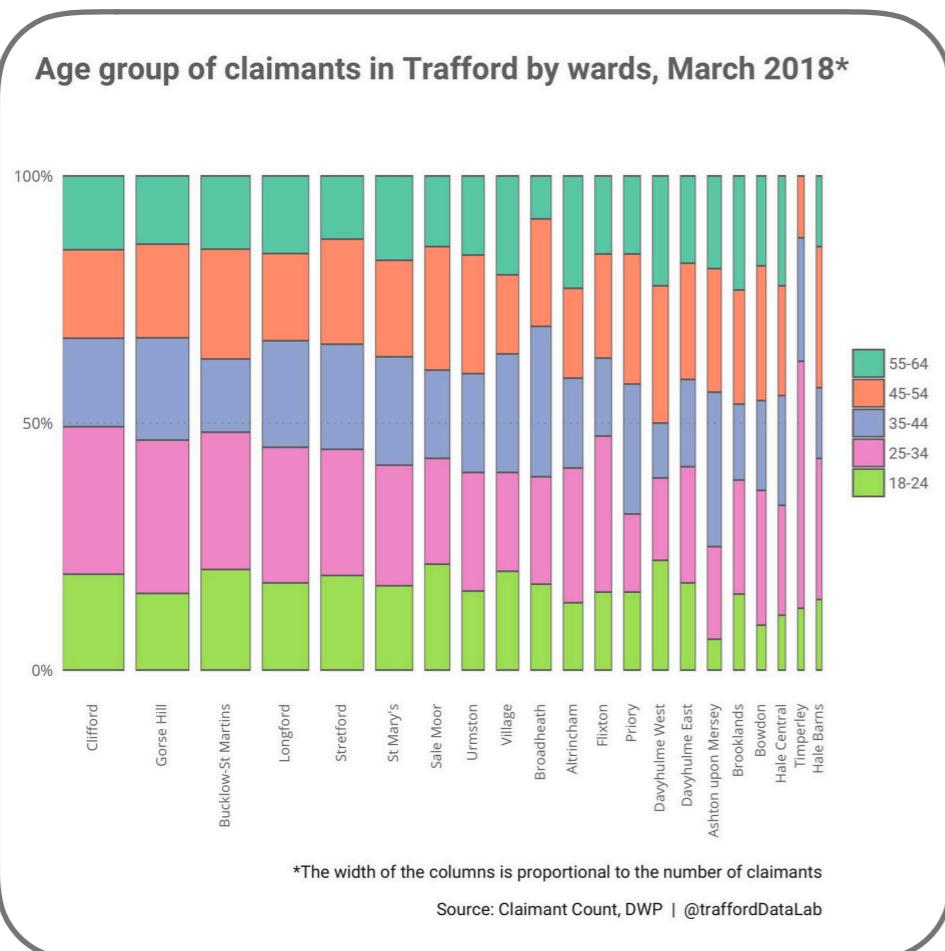
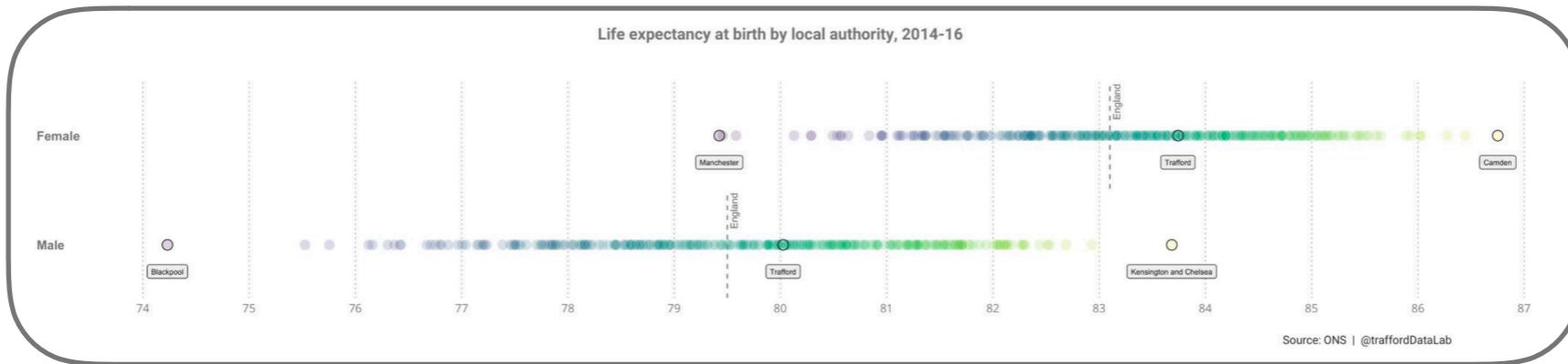
# how many followers are bots?
bot <- tweetbotornot(followers_info$screen_name[1:50], fast = FALSE) %>%
  arrange(prob_bot)

# arrange by probability estimates
bot[order(bot$prob_bot), ]

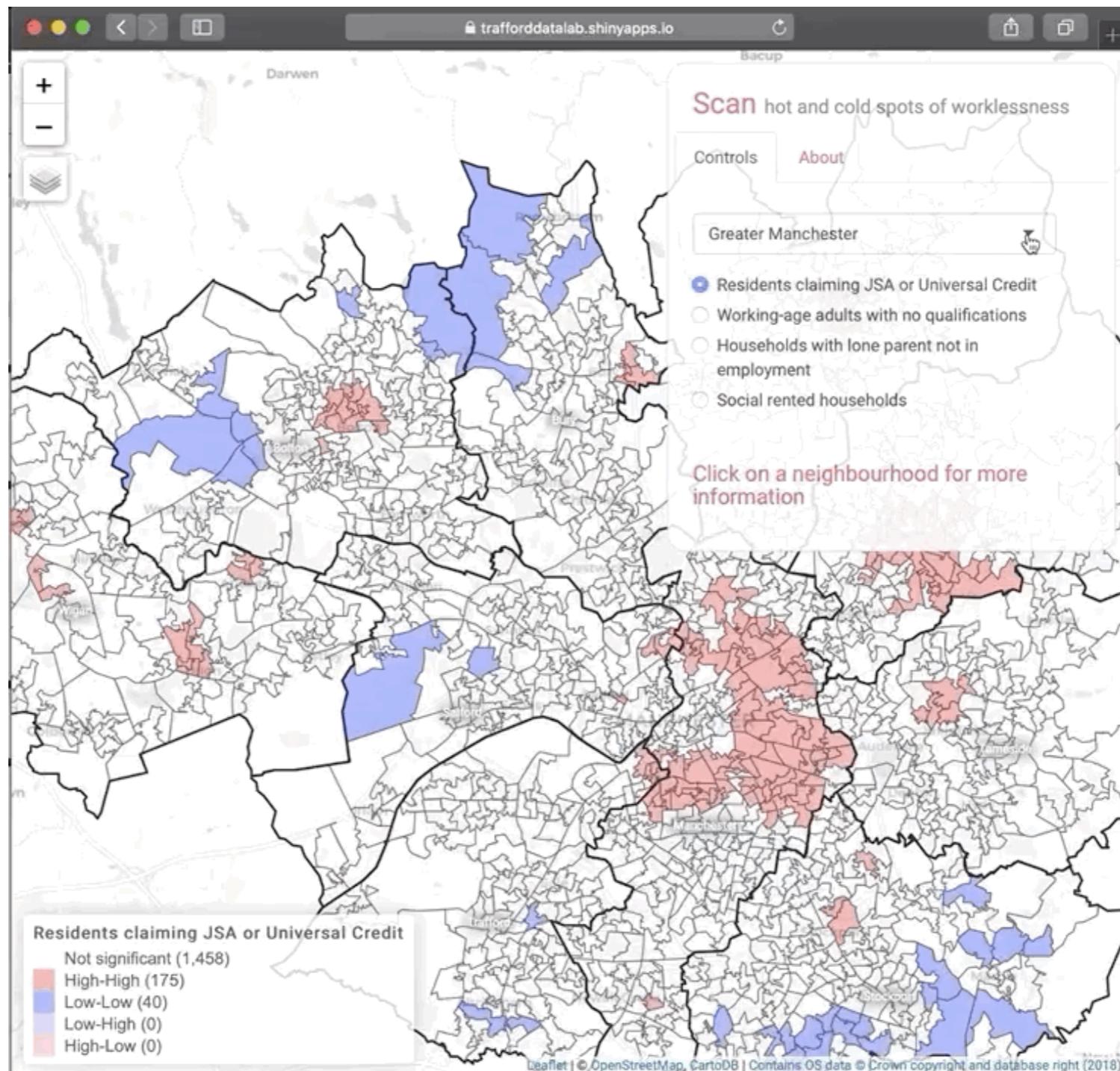
# plot probability estimates
bot %>%
  select(screen_name, prob_bot) %>%
  arrange(prob_bot) %>%
  ggplot() +
  geom_col(aes(x = reorder(screen_name, -prob_bot), y = prob_bot), fill = "#E44690") +
  scale_y_continuous(limits = c(0, 1), expand = c(0,0)) +
  coord_flip() +
  labs(title = "Probability of @traffordDataLab followers being bots",
       x = NULL, y = NULL) +
  theme_minimal() +
  theme(plot.margin=unit(c(1,1,1,1),"cm"),
        panel.grid.major.y = element_blank(),
        panel.grid.minor.y = element_blank(),
        axis.text.y = element_text(hjust = 0))

ggsave("bot_or_not.png", dpi = 300, scale = 1)
```

Graphics

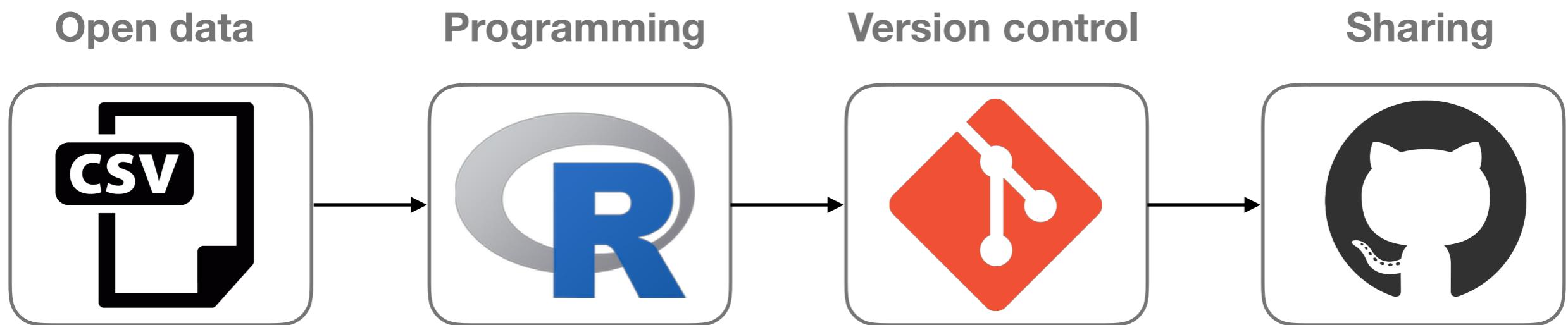


Interactive web apps

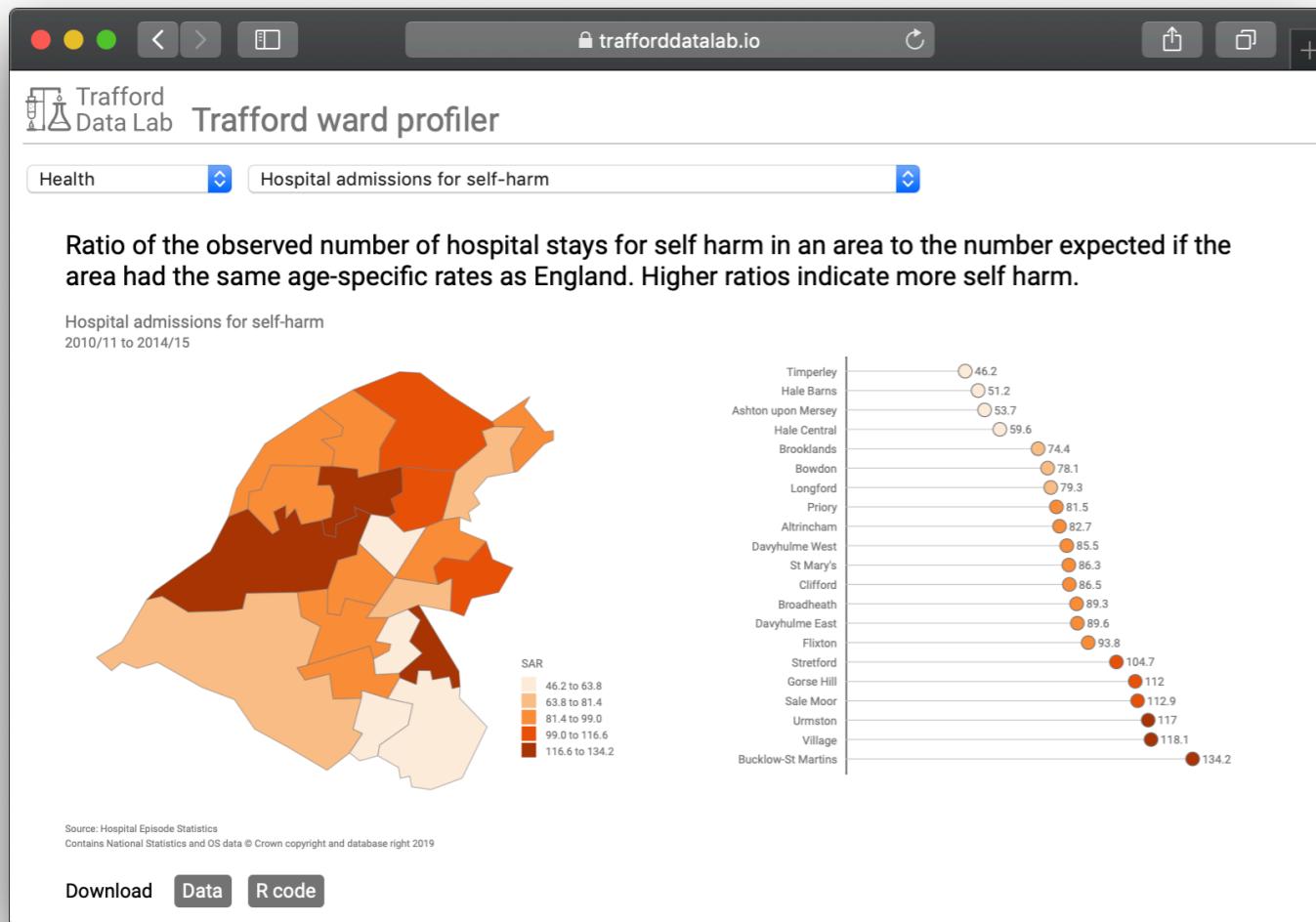


How the Lab uses R

Reproducible workflow



An example



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, extracts specific columns, and writes the results to a CSV file.

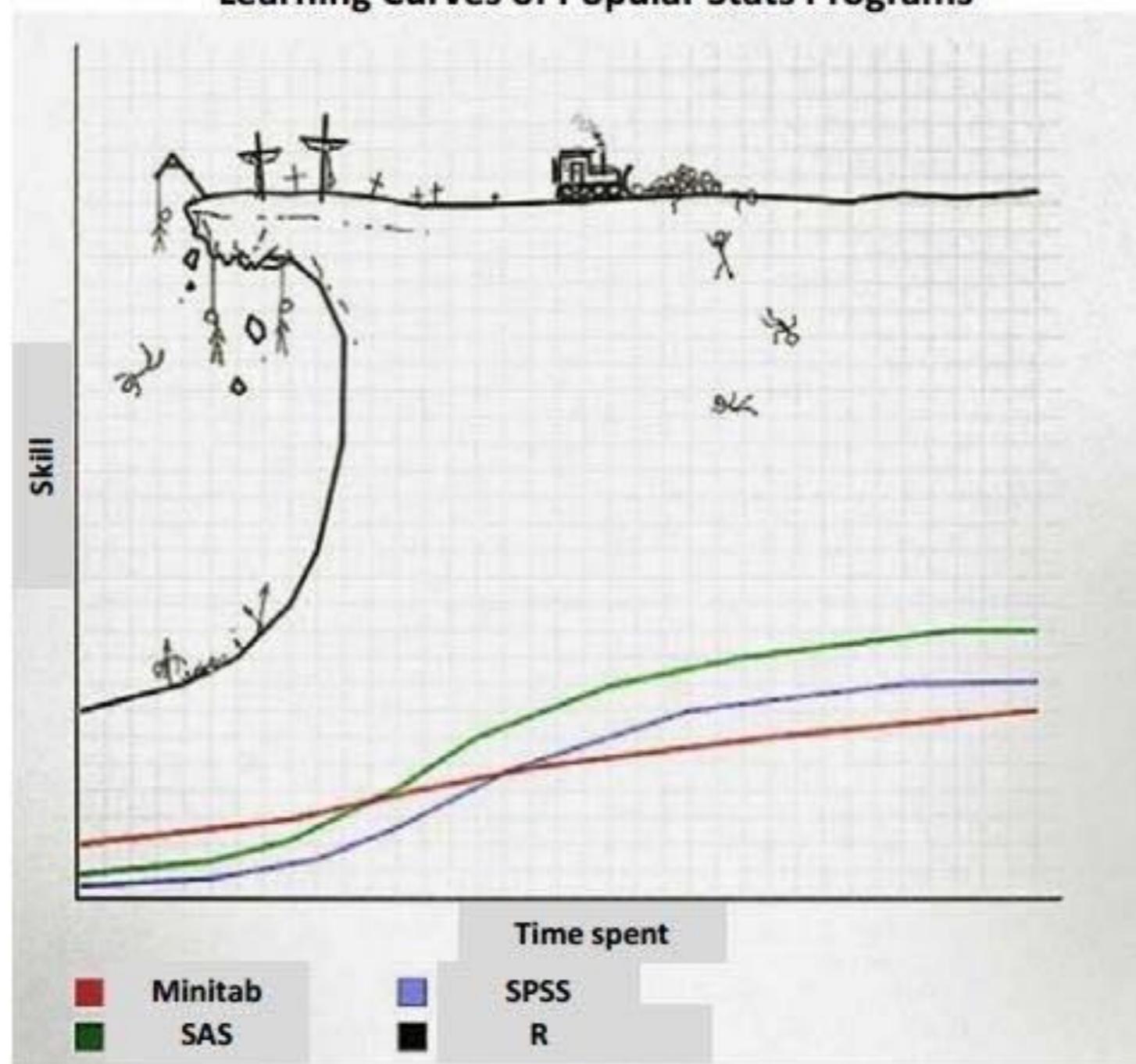
```
1 # Health: Hospital admissions for self-harm, 2010/11 - 2014/15 #
2
3 # Source: Hospital Episode Statistics
4 # URL: https://fingertips.phe.org.uk/
5 # Licence: Open Government Licence
6
7 library(tidyverse) ; library(fingertipsR)
8
9 df <- fingertips_data(IndicatorID = 92584, AreaTypeID = 8) %>%
10   filter(ParentName == "Trafford") %>%
11   select(area_code = AreaCode,
12         area_name = AreaName,
13         value = Value) %>%
14   mutate(period = "2010/11 to 2014/15",
15         indicator = "Hospital admissions for self-harm",
16         measure = "SAR",
17         unit = "Admissions") %>%
18   select(area_code, area_name, indicator, period, measure, unit, value)
19
20 write_csv(df, ".../data/hospital_admissions_self_harm.csv")
```

Scalability

```
1 # Health: Hospital admissions for self-harm, 2010/11 – 2014/15 #
2
3 # Source: Hospital Episode Statistics
4 # URL: https://fingertips.phe.org.uk/
5 # Licence: Open Government Licence
6
7 library(tidyverse) ; library(fingertipsR)
8
9 df <- fingertips_data(IndicatorID = 92584, AreaTypeID = 8) %>%
10   filter(ParentName == "Trafford") %>%
11   select(area_code = AreaCode,
12         area_name = AreaName,
13         value = Value) %>%
14   mutate(period = "2010/11 to 2014/15",
15         indicator = "Hospital admissions for self-harm",
16         measure = "SAR",
17         unit = "Admissions") %>%
18   select(area_code, area_name, indicator, period, measure, unit, value)
19
20 write_csv(df, ".../data/hospital_admissions_self_harm.csv")
```

Learning R

Learning Curves of Popular Stats Programs



Getting started

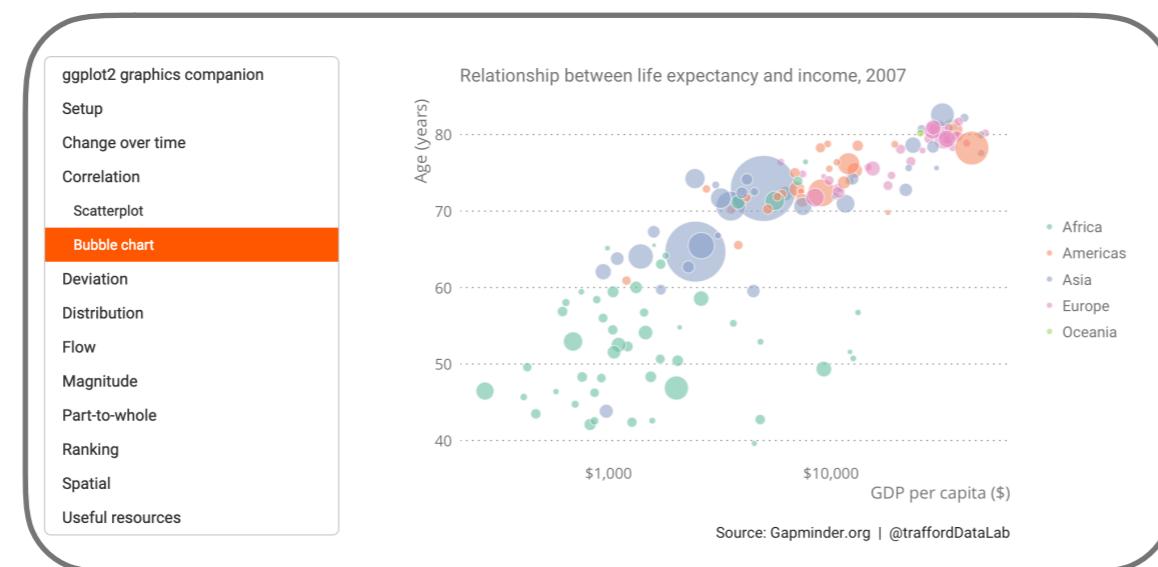
- Reach for R rather than Excel
- Follow [@hadleywickham](#), [@JennyBryan](#), [@dataandme](#), [@NHSrCommunity](#) and monitor #rstats tweets on Twitter
- Find answers or seek help on [stackoverflow](#) and [RStudio Community](#)
- Get a [GitHub](#) profile and commit your code
- Create a [blog using R](#) and post some tutorials
- Set up an R User Group with your colleagues

Lab resources

GitHub



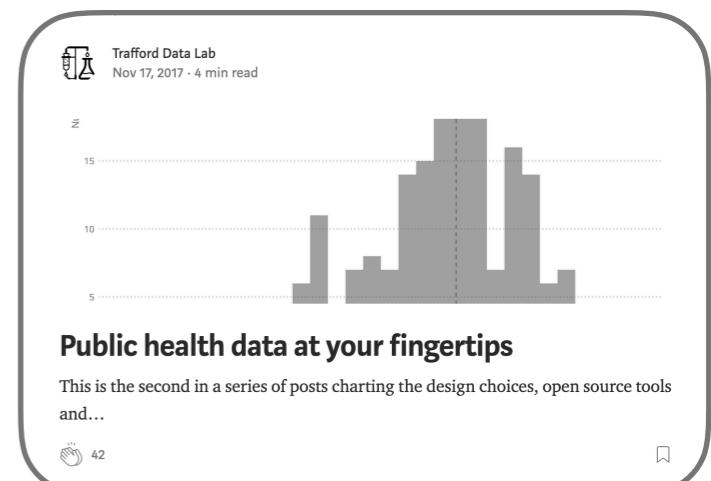
ggplot2 graphics companion



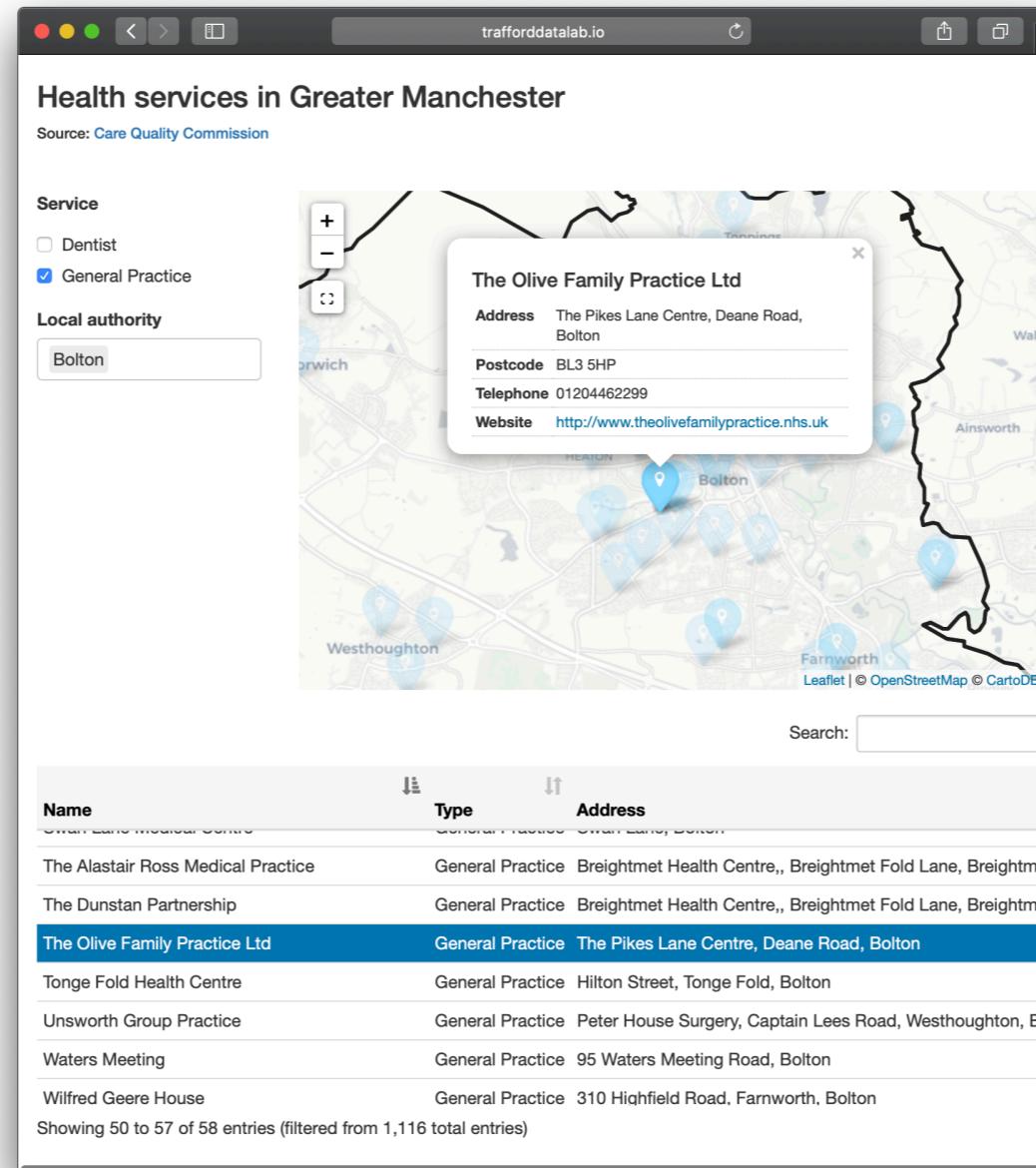
ggplot2 graphics companion

- Setup
- Change over time
- Correlation
- Scatterplot
- Bubble chart
- Deviation
- Distribution
- Flow
- Magnitude
- Part-to-whole
- Ranking
- Spatial
- Useful resources

tutorials



Have a play



</> https://github.com/traffordDataLab/talks/tree/master/2019-03-22_GMPHIN/play

Questions?