# Tom's simple guide to writing a script

Thomas Battram (he/him)

Credit: Dr Gwen Fernandes on helping prepare slides
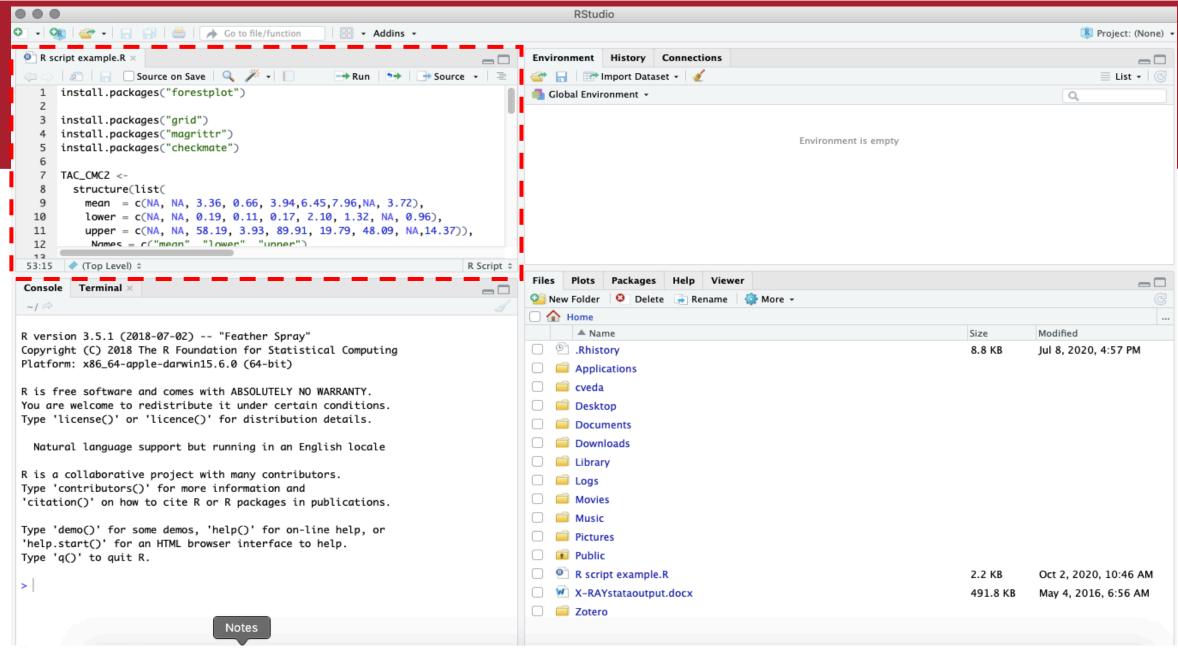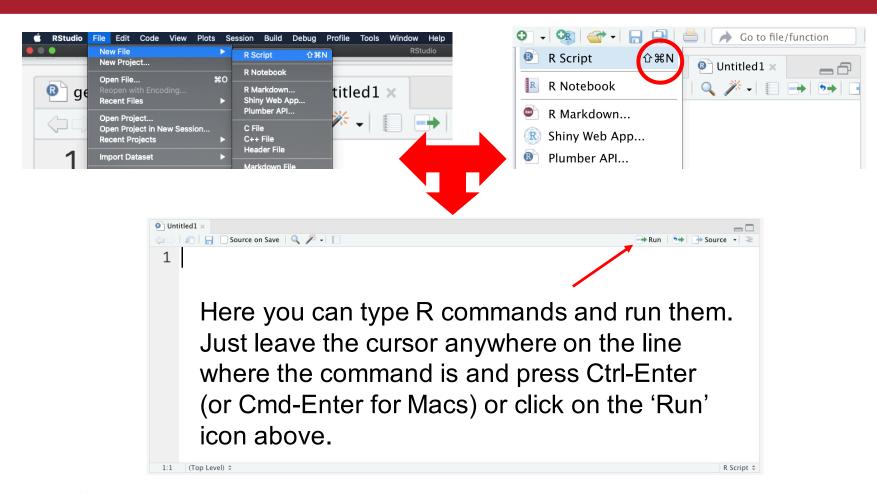
University of BRISTOL

bristol.ac.uk

# Outline

- Making a new R script
- Commenting your scripts
- What to think about before starting to script
- How the top of your script should look
- Making your code easier to read
- Saving your scripts
- Bad/good script comparison

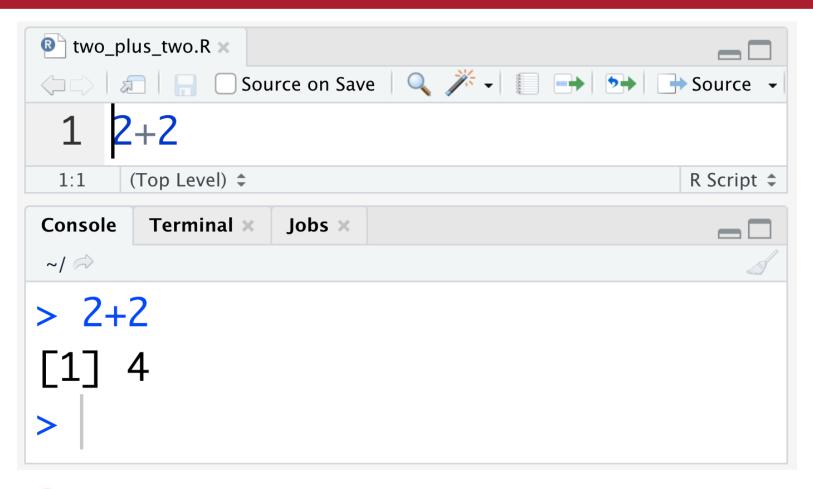bristol.ac.uk

# Making a new R script

– The usual R Studio screen has four windows:

> ➢ CONSOLE
>
> ➢ WORKPLACE AND HISTORY
>
> ➢ FILES, PLOTS, PACKAGES AND HELP
>
> ➢ **R SCRIPT AND DATA VIEW** (this is where you keep a record of your work. For Stata users, this would be like your do-file, for SPSS users it is like the syntax and for SAS users, the SAS program)

# Making a new R script



Here you can type R commands and run them. Just leave the cursor anywhere on the line where the command is and press Ctrl-Enter (or Cmd-Enter for Macs) or click on the 'Run' icon above.

# Making a new R script

# Commenting

- Code with a # before it does not get run
- This is useful for making your scripts much easier to read!
- Comment on WHY and WHAT (to start with)
  - ➢ Start with many comments! Also, use comments to split up the script to make it clearer



bristol.ac.uk

# RStudio sections

# 3 steps before starting

1. Why am I writing this script?

2. What do I need to write this script?

3. How am I going to write this script?

# Step 1: Why am I writing this script?

- Think about it and write it down!
  - e.g. clean a dataset OR assess association between x and y
- Give the script a good name, write a descriptive title + add a couple of lines that describe the purpose of the script

```
generate_pack_years.R

1  # ------------------------------------
2  # Generating pack years in ALSPAC
3  # ------------------------------------
4
5  # This script extracts smoking variables from
6  # the mothers within ALSPAC and uses this to
7  # generate pack years.
8  # Authors: Thomas Battram, Gwen Fernandes.
9  # Date: 2020/02/15
```

bristol.ac.uk

# Step 2: What do I need to write this script?

- Datasets!
- Using only base R can make things difficult… Packages!
- Packages are made by others and are there to make your life easier

For example, reading data into R can be tricky depending on how the data is stored, but packages can make this easier!

bristol.ac.uk

# Reading data into R

Depending on what form the data is in, you have to use different functions to read in the data. Data you get may be in:

- excel spreadsheets
- comma seperated value (csv) files
- tab separated value (tsv) files
- spss files
- stata files
- images (e.g. .png files)

```
11  # load packages
12  library(haven)
13
14  # read in data
15  df <- read_dta("my_data.dta")
```

These need different functions to read them in, some of which are only available with certain packages.

# Step 3: How am I going to write this script?

– Linked to why you are writing it
  and what you need to write it!

– Write out each step

```
17  # structure of the script:
18  # 1. Extract the smoking variables
19  # 2. Exclude individuals with withdrawn
20  #      consent and too much missing data
21  # 3. Generate pack years variable
22  # 4. Check for outliers
23  # 5. Write out a table with identifiers
24  #      and pack years variables
```

bristol.ac.uk

# Top of the script

```
generate_pack_years.R ×
      Source on Save                                                          Run          Source

 1  # ----------------------------------------
 2  # Generating pack years in ALSPAC
 3  # ----------------------------------------
 4
 5  # This script extracts smoking variables from
 6  # the mothers within ALSPAC and uses this
 7  # to generate pack years.
 8  # Authors: Thomas Battram, Gwen Fernandes.
 9  # Date: 2020/02/15
10
11  # load packages
12  library(haven)
13
14  # read in data
15  df <- read_dta("my_data.dta")
16
17  # structure of the script:
18  # 1. Extract the smoking variables
19  # 2. Exclude individuals with withdrawn consent and too much missing data
20  # 3. Generate pack years variable
21  # 4. Check for outliers
22  # 5. Write out a table with identifiers and pack years variables

6:17    (Untitled)                                                                    R Script
```

bristol.ac.uk

# Make your code easy to read

1. Use a consistent style when writing code

```
a_var <- c(1, 2, 3)
a.var <- c(1, 2, 3)
aVar <- c(1, 2, 3)
```

2. Use spaces appropriately

```
x <-c(1,2,51,124,4124)
x <- c(1, 2, 51, 124, 4124)
x    <-              c(1    ,2    ,    51,124      , 4124)
```

3. Use indents appropriately
   a) No long lines of code!

```
my_plot <- plot(x = my_data$x, y = mydata$y, xlab = "
```

```
my_plot <- plot(x = my_data$x,
                y = my_data$y,
                xlab = "My X-axis label",
                ylab = "My Y-axis label",
                main = "My plot title")
```

bristol.ac.uk

# Saving scripts

–It's important to save scripts as you go and you can always come back to them and send them to other people.

–Give the script a good name and save it regularly!



CTRL (or CMD) + S

# Summary

3 questions to think about before writing a script:

    1.Why am I writing this script?

    2.What do I need to write this script?

    3.How am I going to write this script?

    –Comment your code a lot (using #)

    –Make your code easy to read

    –Save your scripts regularly

# Bad Script

1. No title
2. No description
3. No sections
4. No comments at all!
5. Spacing could potentially be improved

```
1    setwd("a_directory")
2    data_cveda <- read.delim("data_cveda.txt")
3    data_cveda
4
5    names(data_cveda)
6    data_cveda$ageband
7
8    data_cveda$child<- data_cveda$ageband==1
9    summary
10
11   mean(data_cveda$p_sdq_tot, na.rm=TRUE)
12   sd(data_cveda$p_sdq_tot,na.rm=TRUE)
13   median(data_cveda$p_sdq_tot, na.rm=TRUE)
14   IQR(data_cveda$p_sdq_tot,     na.rm=TRUE)
15   mad(data_cveda$p_sdq_tot,na.rm=TRUE)
16   min(data_cveda$p_sdq_tot,na.rm=TRUE)
17   max(data_cveda$p_sdq_tot, na.rm=TRUE)
18   range(data_cveda$p_sdq_tot, na.rm=TRUE)
19   summary(data_cveda$p_sdq_tot, na.rm=TRUE)
20
21   d<- data_cveda[data_cveda$ageband==1,]
22   mean(d$p_sdq_tot, na.rm=TRUE)
23   sd(d$p_sdq_tot, na.rm=TRUE)
24   median(d$p_sdq_tot, na.rm=TRUE)
25   IQR(d$p_sdq_tot,na.rm=TRUE)
26   mad(d$p_sdq_tot, na.rm=TRUE)
27   min(d$p_sdq_tot, na.rm=TRUE)
28   max(d$p_sdq_tot, na.rm=TRUE)
29   range(d$p_sdq_tot, na.rm=TRUE)
30   summary(d$p_sdq_tot, na.rm=TRUE)
31
32   cor(data_cveda$p_sdq_emotion, data_cveda$p_sdq_conduct, use= "pairwise.complete.obs")
33
34   table(data_cveda$depanx_1stdeg)
35   table(data_cveda$depanx_1stdeg, data_cveda$sex)
36
```

# Good Script

```r
1    # ---------------------------------------
2    # cVEDA: Association between parental punishment and behaviour
3    # ---------------------------------------
4
5    # This script assesses the association between parental corporal punishment
6    # and total sdq (Strength and Difficulty Questionnaire) score within the cveda cohort
7    # Date: 2020-10-21
8
9    # Set directory and import data and save as data_cveda
10   setwd("a_directory")
11
12   data_cveda <- read.delim("data_cveda.txt")
13
14   # ---------------------------------------
15   # Examine the data and add a variable for children
16   # ---------------------------------------
17
18   # list the variables in data_cveda
19   names(data_cveda)
20   ### The names can be used to refer to the relevant variables by using the $ symbol.
21   data_cveda$ageband
22
23   ### Add which individuals in the dataset are children or in ageband=1
24   # Code below gives you the number of children included in the study.
25   # An additional T/F column is added to the data_cveda data frame ###
26   data_cveda$child <- data_cveda$ageband == 1
27
28   summary(data_cveda$child)
29   # 1918 children in the dataset
30
31   # Examine the statistical qualities of the data
32   mean(data_cveda$p_sdq_tot, na.rm = TRUE)
33   sd(data_cveda$p_sdq_tot, na.rm = TRUE) # standard deviation
34   median(data_cveda$p_sdq_tot, na.rm = TRUE)
35   IQR(data_cveda$p_sdq_tot, na.rm = TRUE) # interquartile range
36   mad(data_cveda$p_sdq_tot, na.rm = TRUE) # mean absolute deviation
37   min(data_cveda$p_sdq_tot, na.rm = TRUE) # minimum value
38   max(data_cveda$p_sdq_tot, na.rm = TRUE) # maximum value
39   range(data_cveda$p_sdq_tot, na.rm = TRUE)
```

bristol.ac.uk

# Good Script

```
53    summary(d_child$p_sdq_tot, na.rm = TRUE) # 0, 8, 12, 11.97, 16, 32, NA (56)
54
55    # ------------------------------------
56    # Run the analyses!
57    # ------------------------------------
58
59    # correlation between SDQ emotion problems and SDQ conduct problems
60    cor(data_cveda$p_sdq_emotion, data_cveda$p_sdq_conduct, use = "pairwise.complete.obs")
61
62    # table of numbers with anxiety/depression in a first degree relative
63    table(data_cveda$depanx_1stdeg)
64    ###The following will produce a three-way cross-tabulation between anxiety/depression in a first degree relative with gender. ###
65    table(data_cveda$depanx_1stdeg, data_cveda$sex)
```

# Any Questions?