

Data Wrangling 101

for BCS 206

Fall 2019

Contents

1	Goals for next two weeks	1
2	Preliminaries	2
2.1	Version control	2
2.2	Reproducibility and literate coding	2
3	Data wrangling	2
3.1	An example data set	2
3.2	Dplyr's verbs	3
3.3	Maggritr's pipes	3
3.4	Putting it together: Wrangling through pipes	3
4	Exercises	5
5	Combining data wrangling and visualization: an example from the Haefner group	5
5.1	Design	5
5.2	Loading data	7
5.3	Figure 2 from Herce Castañón et al.	8
5.3.1	Panel A	8
5.3.2	Panel B	9
5.3.3	Panel A and B together	10
6	Session info	11

1 Goals for next two weeks

- Thinking about workflow in R:
 - Version control
 - R Markdown
- Data wrangling: Turning the data into the form you need (*dplyr*)
- Data visualization:
 - General principles
 - How to plot in R (*ggplot*, *plotly*)

We only have a relatively short time, so we will focus on learning what tools are available and on *examples* of use (rather than an in-depth tutorial). There are great online tutorials and cheatsheets that contain further information.

2 Preliminaries

2.1 Version control

RStudio makes version control, data backup, and data sharing easy (e.g., via Github.com). To use it, download and install git on your computer. Get a free github.com or bitbucket.com account. You only have to do this once.

Then, for each project, create a new project in RStudio and link it to the remote repository (select “Create project” > “Version control”). You will have to enter a URL for the remote repository, which you get, for example, at github.com under the repository’s main page by clicking the “Clone or download button”.

For step by step instructions, see:

- Setting up RStudio for version control
- RStudio help on version control
- Reverting a file to an earlier version

2.2 Reproducibility and literate coding

R and RStudio support reproducibility oriented literate coding via Sweave and Knitr: lab books, presentations, and papers can weave/knit together data, code, and text. The document you share contains the code needed to create its outputs (figures, tables, etc.). This is achieved by combining latex or R markdown with R code (or, for that matter, code from other programming languages). For an excellent video-based introduction, see this tutorial on R markdown. *This document is R markdown compiled with RStudio’s knitr.

3 Data wrangling

The *R* libraries *dplyr* provide us with efficient ways to transform (‘wrangle’) our data tables. The library *magrittr* let’s us concatenate these operations in transparent and easy to read code.

3.1 An example data set

We will illustrate the use of *dplyr* with the following data from an experiment with a 2AFC task in three within-subject conditions (A, B, C), for which we have extracted correctness (1 = correct; 0 = incorrect) and reaction times (RT):

```
summary(d)
```

```
## condition      trial      subject      correct
## A:2688   Min.    : 1.00    1      : 192   Min.    :0.0000
## B:2688   1st Qu.:16.75    2      : 192   1st Qu.:0.0000
## C:2688   Median  :32.50    3      : 192   Median :1.0000
##          Mean    :32.50    4      : 192   Mean   :0.5905
##          3rd Qu.:48.25    5      : 192   3rd Qu.:1.0000
##          Max.    :64.00    6      : 192   Max.   :1.0000
##                                     (Other):6912
##
##          RT
## Min.    : 201.1
## 1st Qu.: 405.7
## Median : 573.8
```

```
## Mean    : 771.3
## 3rd Qu.:1074.4
## Max.    :2877.3
##
```

```
glimpse(d)
```

```
## Observations: 8,064
## Variables: 5
## $ condition <fct> A, A, A, A, A, A, A, A, A, A, A, A, A, A, A, A, A, A...
## $ trial     <int> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1...
## $ subject   <fct> 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 1...
## $ correct   <int> 1, 1, 0, 0, 0, 1, 1, 0, 0, 1, 0, 0, 0, 1, 1, 0, 1, 1...
## $ RT        <dbl> 572.054, 518.748, 358.922, 424.082, 402.441, 350.779...
```

3.2 Dplyr's verbs

Dplyr has 'verbs' like filter, select, summarize, mutate, transmute, etc. to let use conduct operations on our data, and reshape the data frame into the format we need. We can use dplyr, for example, to calculate the proportion correct answers in our experiment by using *summarise*.

```
summarise(d, meanCorrect = mean(correct))
```

```
## # A tibble: 1 x 1
##   meanCorrect
##         <dbl>
## 1         0.591
```

Or just for condition A:

```
d.A = filter(d, condition == "A")
summarise(d.A, meanCorrect = mean(correct))
```

```
## # A tibble: 1 x 1
##   meanCorrect
##         <dbl>
## 1         0.461
```

3.3 Maggritr's pipes

Here we will use only of the 'pipes' magrittr provides:

- `x %>% f`: takes `x` and hands it to the function `f` on the right, as `f`'s first argument
- `x %<>% f1 %>% f2 %>%` etc.: takes `x` hands it to `f1`, takes the output of `f1` and hands it to `f2`, etc. And since the first pipe was `%<>%` (rather than just `%>%`), the final result will be written back into `x`.

3.4 Putting it together: Wrangling through pipes

Remember how we got the mean proportion correct for just Condition A?

```
d.A = filter(d, condition == "A")
summarise(d.A, meanCorrect = mean(correct))
```



Figure 1: Magritt's pipe



Figure 2: Magrittr's pipe

```
## # A tibble: 1 x 1
##   meanCorrect
##       <dbl>
## 1         0.461
```

This is inelegant and hard to read. Pipes let us make this more transparent:

```
d %>%
  filter(condition == "A") %>%
  summarise(meanCorrect = mean(correct))
```

```
## # A tibble: 1 x 1
##   meanCorrect
##       <dbl>
## 1         0.461
```

And this advantage becomes even clearer, the more operations we concatenate. For example, `group_by` is an elegant operator that tells the pipes to conduct all subsequent operations for each of the groups (and then put all the separate outcomes back together into a single data frame). So if we want the proportion correct for all groups:

```
d %>%
  group_by(condition) %>%
  summarise(meanCorrect = mean(correct))
```

```
## # A tibble: 3 x 2
##   condition meanCorrect
##   <fct>         <dbl>
## 1 A           0.461
## 2 B           0.537
## 3 C           0.774
```

4 Exercises

How can we:

- View the entire data set? (*View*)
- Calculate the by-subject averages for all three conditions? (*group_by, summarise*)
- Calculate the by-subject standard deviations around those averages? (*group_by, summarise*)
- Attach this information (the averages and SDs) to each row of the present data.frame? (*group_by, mutate*)
- Determine whether RTs were on average faster for correct, as compared to incorrect, trials?
- Add a column for log-transformed RTs to the data set?
- Remove the old column for raw RTs? (*select*)
- Sort the data by log-transformed reaction times? (*arrange*)

Say we further have an additional data frame with information about our subjects:

```
## Source: local data frame [42 x 3]
## Groups: <by row>
##
## # A tibble: 42 x 3
##   subject gender  age
##   <fct>   <chr> <dbl>
## 1 1      male    21
## 2 2      male    19
## 3 3      female  18
## 4 4      male    20
## 5 5      female  19
## 6 6      male    20
## 7 7      female  19
## 8 8      female  20
## 9 9      male    20
## 10 10     male    19
## # ... with 32 more rows
```

- How can we join the information from the two data sources together? (*left_join*)

5 Combining data wrangling and visualization: an example from the Haefner group

This group seeks to replicate Herce Castañón et al. (2019).

5.1 Design

The design of the present study crossed two levels of contrast (Low = 15%, High = 60%), 3 levels of variance (0, 4, 10), and how the trials in the block were cued (L = left, R = right, N = uncued), for a total of $2 \times 3 \times 3 = 18$ within-subject conditions.

Fig. 2

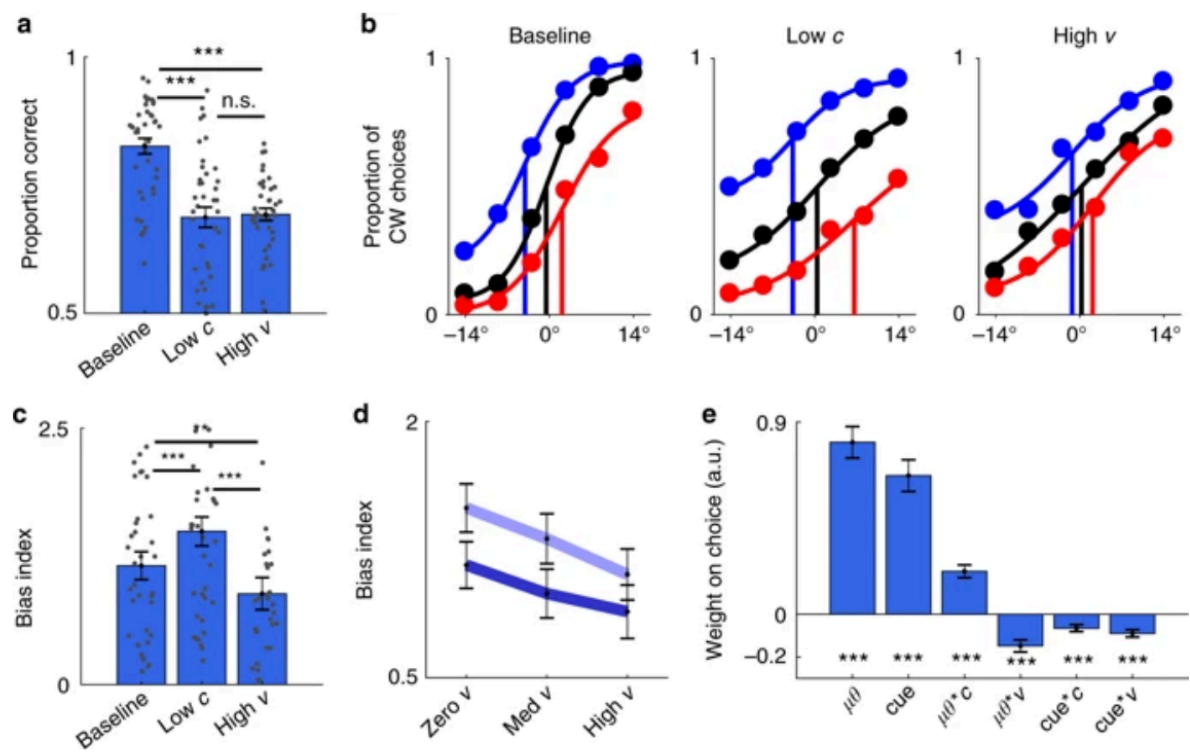


Figure 3: Figure 2 from Herce Castañón et al. (2019)

5.2 Loading data

The data are stored in a MatLab (.mat) file. The file contains one matrix with fields: participant, exp(eriment), stimuli and response. Within each field, there is further information. The important information seems to be in the response field. Some of the important parts include:

- responseRight: the response of the subject (0 for CCW, 1 for CW, w.r.t horizontal)
- correct: what the correct answer is (0 for CCW, 1 for CW, w.r.t horizontal)
- accuracy: whether subject got the correct answer (1) or not (0)
- reaction time: time in seconds the subject took to answer
- confidence: whether the subject was confident in their answer (1) or not (-1)
- cue: whether the cue on that trials is left (-1), right (1), or no cue (0)
- contrast: the contrast of the gabor patch on that trial
- variance: variability in the orientation of gratings of gabor patches on that trial
- isCuedBlock: whether a block (of trials) will have cues (1) or no cues (0)

```
# Load a matlab file and extract the "data" matrix out of it
d.haefner = readMat("./data/Haefner/uncertaintyV1-subject18-1-EarlyQuit.mat")
d.haefner = d.haefner[["data"]][, , 1][["response"]][, , 1]
d.haefner[["trueOrientations"]] <- NULL

# Look at what we've imported.
# NB: str() gives you the structure of an R object
str(d.haefner)

d.haefner %<>%
  map(.f = function(x) c(x)) %>%
  as_tibble()

# The data we have are preliminary pilot data from one of the
# experimenters, and that run did contain all trials. We omit
# all the trials with missing information.
d.haefner %<>%
  na.omit()

# Add the definition of the three conditions of interest in the
# original paper
d.haefner %<>%
  mutate(
    condition = case_when(
      variance == min(variance) & contrast == max(contrast) ~ "baseline",
      variance == max(variance) & contrast == max(contrast) ~ "high variance",
      variance == min(variance) & contrast == min(contrast) ~ "low contrast",
      T ~ ""
    )
  )
```

Now that we've imported the data into an R data frame (or *tibble*), let's have a look at it. First, we can get a general idea of the data by using `str()` (for structure) or `print()`:

```
## # A tibble: 864 x 12
##   randSeed responseRight correct accuracy reactionTime confidence
##   <dbl>         <dbl>   <dbl>   <dbl>         <dbl>         <dbl>
## 1  2.20e8             0       1       0           0.708          -1
## 2  2.20e8             1       0       0           0.609          -1
## 3  2.20e8             0       1       0           1.73           0
```

```
## 4 2.20e8 1 0 0 0.684 0
## 5 2.20e8 0 1 0 0.550 -1
## 6 2.20e8 1 1 1 0.565 -1
## 7 2.20e8 0 1 0 0.492 -1
## 8 2.20e8 0 0 1 0.994 -1
## 9 2.20e8 0 1 0 0.872 0
## 10 2.20e8 0 1 0 0.782 1
## # ... with 854 more rows, and 6 more variables: isCuedBlock <dbl>,
## # cue <dbl>, orientationMean <dbl>, contrast <dbl>, variance <dbl>,
## # condition <chr>
```

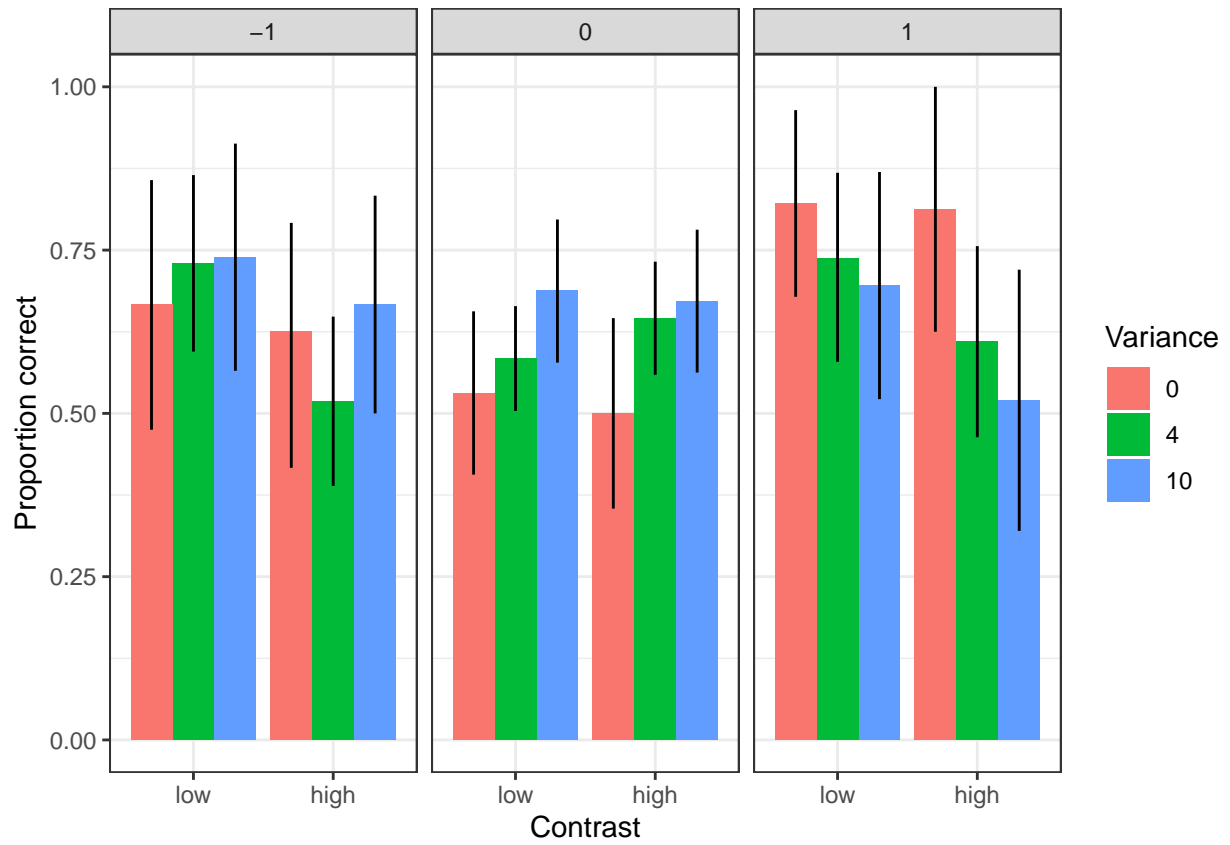
To instead get a summary of the data:

```
##      randSeed      responseRight      correct      accuracy
## Min. :220286057 Min. :0.000 Min. :0.0000 Min. :0.0000
## 1st Qu.:220377520 1st Qu.:0.000 1st Qu.:0.0000 1st Qu.:0.0000
## Median :220475950 Median :1.000 Median :1.0000 Median :1.0000
## Mean :220480775 Mean :0.559 Mean :0.5081 Mean :0.6319
## 3rd Qu.:220587872 3rd Qu.:1.000 3rd Qu.:1.0000 3rd Qu.:1.0000
## Max. :220674325 Max. :1.000 Max. :1.0000 Max. :1.0000
##      reactionTime      confidence      isCuedBlock      cue
## Min. :0.1983 Min. : -1.0000 Min. :0.0000 Min. : -1.00000
## 1st Qu.:0.4322 1st Qu.: -1.0000 1st Qu.:0.0000 1st Qu.: 0.00000
## Median :0.5677 Median : -1.0000 Median :0.0000 Median : 0.00000
## Mean :0.6970 Mean : -0.5289 Mean :0.4167 Mean : -0.02083
## 3rd Qu.:0.8143 3rd Qu.: 0.0000 3rd Qu.:1.0000 3rd Qu.: 0.00000
## Max. :2.9974 Max. : 1.0000 Max. :1.0000 Max. : 1.00000
##      orientationMean      contrast      variance      condition
## Min. : -26.67646 Min. :0.1500 Min. : 0.00 Length:864
## 1st Qu.: -5.95720 1st Qu.:0.1500 1st Qu.: 4.00 Class :character
## Median : 0.15682 Median :0.1500 Median : 4.00 Mode :character
## Mean : 0.08092 Mean :0.3734 Mean : 4.66
## 3rd Qu.: 5.95273 3rd Qu.:0.6000 3rd Qu.:10.00
## Max. : 25.98183 Max. :0.6000 Max. :10.00
```

5.3 Figure 2 from Herce Castañón et al.

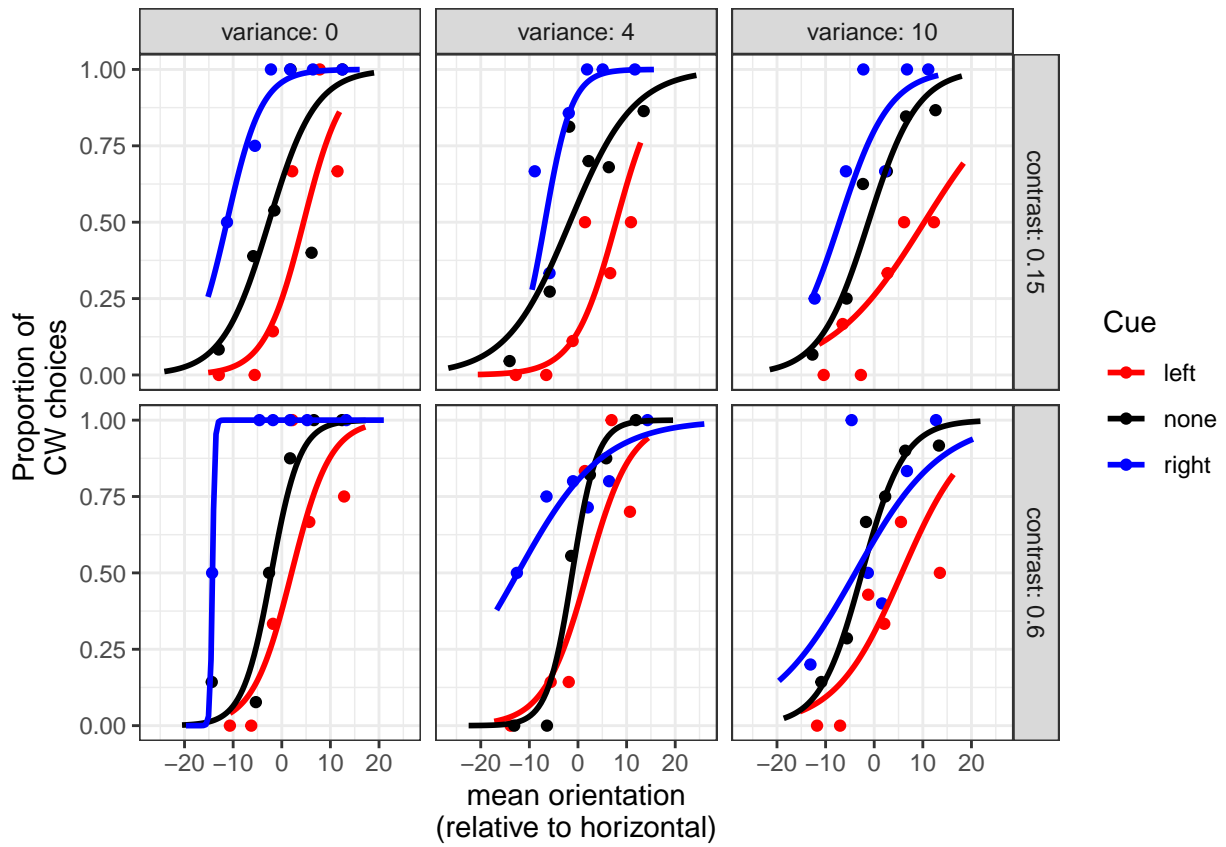
5.3.1 Panel A

We begin by plotting the proportion of correct choices for *all* conditions:

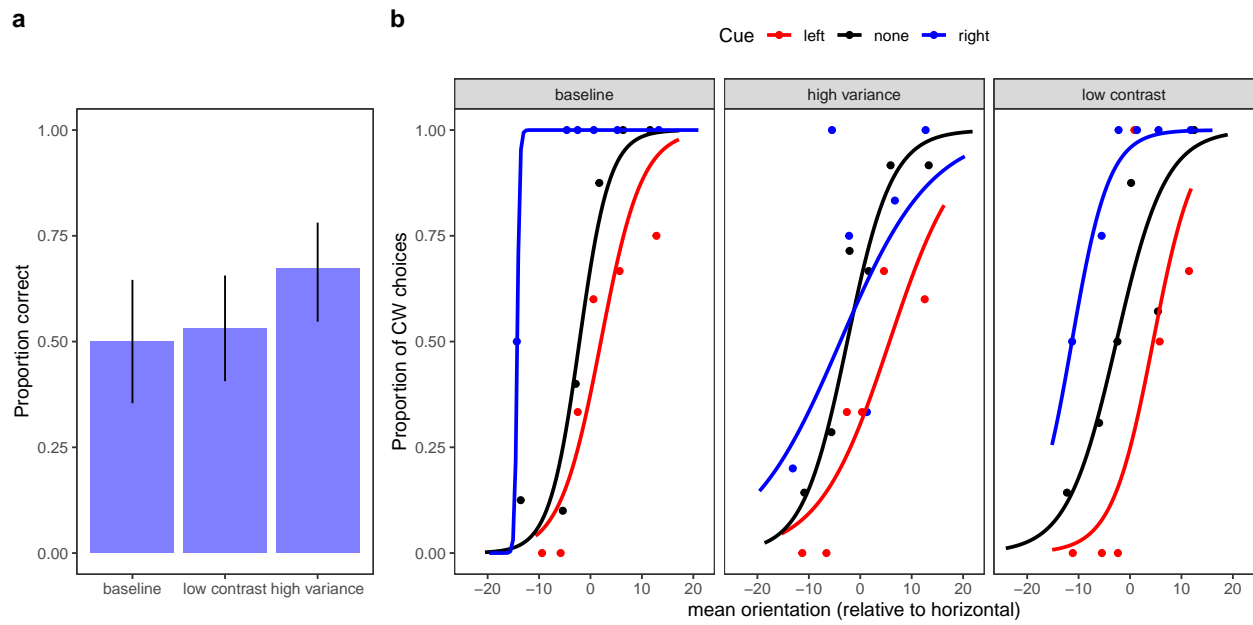


5.3.2 Panel B

We begin by plotting the proportion of CW choices for *all* conditions:



5.3.3 Panel A and B together



6 Session info

```
## - Session info -----
## setting value
## version R version 3.6.0 (2019-04-26)
## os macOS High Sierra 10.13.6
## system x86_64, darwin15.6.0
## ui X11
## language (EN)
## collate en_US.UTF-8
## ctype en_US.UTF-8
## tz America/New_York
## date 2019-11-04
##
## - Packages -----
## package * version date lib source
## acepack 1.4.1 2016-10-29 [1] CRAN (R 3.6.0)
## assertthat 0.2.1 2019-03-21 [1] CRAN (R 3.6.0)
## backports 1.1.5 2019-10-02 [1] CRAN (R 3.6.0)
## base64enc 0.1-3 2015-07-28 [1] CRAN (R 3.6.0)
## broom 0.5.2 2019-04-07 [1] CRAN (R 3.6.0)
## callr 3.3.2 2019-09-22 [1] CRAN (R 3.6.0)
## cellranger 1.1.0 2016-07-27 [1] CRAN (R 3.6.0)
## checkmate 1.9.4 2019-07-04 [1] CRAN (R 3.6.0)
## cli 1.1.0 2019-03-19 [1] CRAN (R 3.6.0)
## cluster 2.1.0 2019-06-19 [1] CRAN (R 3.6.0)
## colorspace 1.4-1 2019-03-18 [1] CRAN (R 3.6.0)
## cowplot * 1.0.0 2019-07-11 [1] CRAN (R 3.6.0)
## crayon 1.3.4 2017-09-16 [1] CRAN (R 3.6.0)
## data.table 1.12.6 2019-10-18 [1] CRAN (R 3.6.0)
## desc 1.2.0 2018-05-01 [1] CRAN (R 3.6.0)
## devtools 2.2.1 2019-09-24 [1] CRAN (R 3.6.0)
## digest 0.6.22 2019-10-21 [1] CRAN (R 3.6.0)
## dplyr * 0.8.3 2019-07-04 [1] CRAN (R 3.6.0)
## ellipsis 0.3.0 2019-09-20 [1] CRAN (R 3.6.0)
## evaluate 0.14 2019-05-28 [1] CRAN (R 3.6.0)
## fansi 0.4.0 2018-10-05 [1] CRAN (R 3.6.0)
## forcats * 0.4.0 2019-02-17 [1] CRAN (R 3.6.0)
## foreign 0.8-72 2019-08-02 [1] CRAN (R 3.6.0)
## Formula 1.2-3 2018-05-03 [1] CRAN (R 3.6.0)
## fs 1.3.1 2019-05-06 [1] CRAN (R 3.6.0)
## generics 0.0.2 2018-11-29 [1] CRAN (R 3.6.0)
## ggplot2 * 3.2.1 2019-08-10 [1] CRAN (R 3.6.0)
## glue 1.3.1 2019-03-12 [1] CRAN (R 3.6.0)
## gridExtra 2.3 2017-09-09 [1] CRAN (R 3.6.0)
## gtable 0.3.0 2019-03-25 [1] CRAN (R 3.6.0)
## haven 2.1.1 2019-07-04 [1] CRAN (R 3.6.0)
## Hmisc 4.2-0 2019-01-26 [1] CRAN (R 3.6.0)
## hms 0.5.2 2019-10-30 [1] CRAN (R 3.6.0)
## htmlTable 1.13.2 2019-09-22 [1] CRAN (R 3.6.0)
## htmltools 0.4.0 2019-10-04 [1] CRAN (R 3.6.0)
## htmlwidgets 1.5.1 2019-10-08 [1] CRAN (R 3.6.0)
## httr 1.4.1 2019-08-05 [1] CRAN (R 3.6.0)
## jsonlite 1.6 2018-12-07 [1] CRAN (R 3.6.0)
```

```

## knitr                1.25      2019-09-18 [1] CRAN (R 3.6.0)
## labeling             0.3       2014-08-23 [1] CRAN (R 3.6.0)
## lattice              0.20-38   2018-11-04 [1] CRAN (R 3.6.0)
## latticeExtra         0.6-28    2016-02-09 [1] CRAN (R 3.6.0)
## lazyeval             0.2.2     2019-03-15 [1] CRAN (R 3.6.0)
## lifecycle            0.1.0     2019-08-01 [1] CRAN (R 3.6.0)
## lubridate            1.7.4     2018-04-11 [1] CRAN (R 3.6.0)
## magrittr             * 1.5      2014-11-22 [1] CRAN (R 3.6.0)
## Matrix               1.2-17    2019-03-22 [1] CRAN (R 3.6.0)
## memoise              1.1.0     2017-04-21 [1] CRAN (R 3.6.0)
## modelr               0.1.5     2019-08-08 [1] CRAN (R 3.6.0)
## munsell              0.5.0     2018-06-12 [1] CRAN (R 3.6.0)
## nlme                 3.1-141   2019-08-01 [1] CRAN (R 3.6.0)
## nnet                 7.3-12    2016-02-02 [1] CRAN (R 3.6.0)
## pillar               1.4.2     2019-06-29 [1] CRAN (R 3.6.0)
## pkgbuild             1.0.6     2019-10-09 [1] CRAN (R 3.6.0)
## pkgconfig            2.0.3     2019-09-22 [1] CRAN (R 3.6.0)
## pkgload              1.0.2     2018-10-29 [1] CRAN (R 3.6.0)
## plotly               * 4.9.0     2019-04-10 [1] CRAN (R 3.6.0)
## plyr                 1.8.4     2016-06-08 [1] CRAN (R 3.6.0)
## prettyunits          1.0.2     2015-07-13 [1] CRAN (R 3.6.0)
## processx             3.4.1     2019-07-18 [1] CRAN (R 3.6.0)
## ps                   1.3.0     2018-12-21 [1] CRAN (R 3.6.0)
## purrr               * 0.3.3     2019-10-18 [1] CRAN (R 3.6.0)
## R.matlab             * 3.6.2     2018-09-27 [1] CRAN (R 3.6.0)
## R.methodsS3          1.7.1     2016-02-16 [1] CRAN (R 3.6.0)
## R.oo                 1.22.0    2018-04-22 [1] CRAN (R 3.6.0)
## R.utils              2.9.0     2019-06-13 [1] CRAN (R 3.6.0)
## R6                   2.4.0     2019-02-14 [1] CRAN (R 3.6.0)
## RColorBrewer         1.1-2     2014-12-07 [1] CRAN (R 3.6.0)
## Rcpp                 1.0.2     2019-07-25 [1] CRAN (R 3.6.0)
## readr               * 1.3.1     2018-12-21 [1] CRAN (R 3.6.0)
## readxl              1.3.1     2019-03-13 [1] CRAN (R 3.6.0)
## remotes              2.1.0     2019-06-24 [1] CRAN (R 3.6.0)
## reshape2            1.4.3     2017-12-11 [1] CRAN (R 3.6.0)
## rlang                0.4.1     2019-10-24 [1] CRAN (R 3.6.0)
## rmarkdown            1.16      2019-10-01 [1] CRAN (R 3.6.0)
## rpart                4.1-15    2019-04-12 [1] CRAN (R 3.6.0)
## rprojroot            1.3-2     2018-01-03 [1] CRAN (R 3.6.0)
## rstudioapi           0.10      2019-03-19 [1] CRAN (R 3.6.0)
## rvest                0.3.4     2019-05-15 [1] CRAN (R 3.6.0)
## scales              1.0.0     2018-08-09 [1] CRAN (R 3.6.0)
## sessioninfo          1.1.1     2018-11-05 [1] CRAN (R 3.6.0)
## stringi              1.4.3     2019-03-12 [1] CRAN (R 3.6.0)
## stringr             * 1.4.0     2019-02-10 [1] CRAN (R 3.6.0)
## survival             2.44-1.1  2019-04-01 [1] CRAN (R 3.6.0)
## testthat             2.2.1     2019-07-25 [1] CRAN (R 3.6.0)
## tibble              * 2.1.3     2019-06-06 [1] CRAN (R 3.6.0)
## tidyr               * 1.0.0     2019-09-11 [1] CRAN (R 3.6.0)
## tidyselect           0.2.5     2018-10-11 [1] CRAN (R 3.6.0)
## tidyverse           * 1.2.1     2017-11-14 [1] CRAN (R 3.6.0)
## usethis              1.5.1     2019-07-04 [1] CRAN (R 3.6.0)
## utf8                 1.1.4     2018-05-24 [1] CRAN (R 3.6.0)
## vctrs                0.2.0     2019-07-05 [1] CRAN (R 3.6.0)

```

```
## viridisLite      0.3.0      2018-02-01 [1] CRAN (R 3.6.0)
## withr            2.1.2      2018-03-15 [1] CRAN (R 3.6.0)
## xfun             0.10       2019-10-01 [1] CRAN (R 3.6.0)
## xml2             1.2.2      2019-08-09 [1] CRAN (R 3.6.0)
## yaml             2.2.0      2018-07-25 [1] CRAN (R 3.6.0)
## zeallot          0.1.0      2018-01-28 [1] CRAN (R 3.6.0)
##
## [1] /Library/Frameworks/R.framework/Versions/3.6/Resources/library
```