

Reproducible Research using R - A template for analysing behavioural experiments

Tuomas Eerola, Durham University

11/22/2018

Contents

1	Why reproducible research?	1
2	Why R (and not Matlab or SPSS)?	2
3	Analysis Template	2
3.1	Suggested folder structure	2
3.2	Using the Template	3
4	Help for statistics with R	11
4.1	Online tutorials	11
4.2	Other Online Resources	11
5	References	11

19/1/2015, update 22/11/2018

1 Why reproducible research?

1. To comply with the increased demands for transparency and open access data, originally formulated within the computer sciences but later spread to biosciences and currently taking hold in social sciences (see Asendorpf et al., 2013; Tomasello & Call, 2011). We should aim to do all our statistical analyses transparently and in a reproducible fashion.
2. To collaborate more easily and effectively. This also helps to spot mistakes and encourage learning and trying out new things (Sandve et al., 2013).
3. To communicate your research more effectively by writing clear analysis paths that are able to produce the key statistics, figures and tables effortlessly. This also helps to gain visibility to research if shared fully (see Piwowar et al., 2007).

2 Why R (and not Matlab or SPSS)?

Here I have chosen *R* to be the chosen tool for reproducible analyses, although of course statistical software could also be used. However, there are several good arguments to support R as the best choice. And I have long personal experience of SPSS and Matlab, both powerful but hampered by various design issues, but R has several advantages over these:

1. R is the most accessible software. R is free, open source, available for all operating systems. Matlab is great for certain type of work, but expensive, fussy about the operating systems, not to mention SPSS in this regard.
2. R is completely programming driven (thus fully transparent). Matlab is equally so, but since it is essentially a *MATrix LABoratory*, it is very good for numerical analyses, but R is a little more versatile for strings and data structures more commonly used in statistics. SPSS also has the syntax option, but it is much more cryptic and unwieldy than R and Matlab. Clear syntax driven operation makes the analyses easily human readable, which is important for collaborations.
3. R has excellent coverage of statistical modelling tools. Thousands of R packages exist for any state-of-the-art statistical technique (bayesian, structural equation modelling, rare regression analytics, all machine-learning algorithms with effective implementations.
4. R is rational and even pedagogical in many of its functionalities (it warns about calculating means for categorical variables, is much more explicit about the outputs, and `data.frames`, etc.)
5. Good support for producing reports ([R Markdown](#), [knitr](#), [Sweave](#)) or interactive websites ([Shiny](#)).

3 Analysis Template

I have prepared an analysis template, which contains examples of the whole process of data analysis; from loading to preprocessing data and analysing and reporting the results. I suggest a certain folder structure to keep the different parts of the processes tightly in different folders. I have been influenced by existing templates¹ and style guides², but it is basically the cleaned up version of the structures that I have for each project.

3.1 Suggested folder structure

A project should have a dedicated folder with a descriptive name it. Within the folder, there is a master file called the `contents.R` which contains the necessary commands to load the data, pre-process it, analyse, and produce figures and tables. For clarity, it is good idea to keep things organised in particular special subfolders.

¹[ProjectTemplate](#)

²[Style Guide](#)

- config** Configuration file (`libraries.R`), which define custom functions and loads the R libraries needed in the project.
- data** Stores the original data, preferably in read-only format. Can be Excel or ascii files, or folders of separate files exported from experiment data collection interface. Be wary of different encodings (UTF-8, Western, UTF-16, etc.) which is the usual cause of problems when reading in data.
- munge** Munger refers to “data munging”, which means cleaning and transforming the data to a suitable format for the analysis. If you have to recode (e.g., invert the scale for a question that has been asked with a reverse wording in comparison to other items) or combine items into indices of instruments, this is the place to do it.
- scr** Scripts used in the analysis. This is the main folder for the analysis. Examples of diagnostics, statistical testing, plotting, generating tables are given.
- figures** All figures and graphs produced by analysis scripts (preferably in pdf format)
- reports** Reports summarising the analysis, either manually or automatically. The latter can be done with [R markdown](#).

3.2 Using the Template

Once you have the template including the data and folder structure as well as *R* installed (or *RStudio*), it should be straightforward to proceed to using the template. The following example loads the first part of the Exp. 1 data in the Sweet Sorrow project containing background survey responses and self-reports given in the experiment from 42 participants (this data is a subset of the study that came out as - Eerola, T., Vuoskoski, J. K., & Kautiainen, H. (2016). Being Moved by Unfamiliar Sad Music Is Associated with High Empathy. *Frontiers in Psychology*, 7, 1176. <http://journal.frontiersin.org/article/10.3389/fpsyg.2016.01176>

3.2.1 Initialise the analysis

Start R and open up the `contents.R` file using your preferred editor. Check that the directory after the first command `setwd` is pointing the location of your analysis directory and run the first lines of the code:

```
## INITIALISE: SET PATH, CLEAR MEMORY AND LOAD LIBRARIES
# set working directory to the project
setwd('~/.Documents/custom/reproducible_data_analysis/data_analysis/')
rm(list=ls(all=TRUE))           # cleans the R memory, just in case
source('config/libraries.R')    # loads the extra R libraries needed
```

If you get errors at this stage with new installation of R, they might refer to the special libraries that were loaded or installed in `libraries.R`. This script should install the required libraries for you such as `ggplot2`, but there might be issues with your particular setup.

3.2.2 Load, preprocess and diagnose the data

Next, it is time to load the data with two scripts, the first one `read_data_survey.R` is simply reading an excel file stored in data folder whereas the `read_data_behavioural.R` reads separate ascii files exported from Matlab (and stored in a sub-folder under the data).

Note about existing data formats. If your data comes from *Qualtrics*, *Survey Monkey*, *Online Surveys*, or some other existing service, it is likely that there are existing R scripts designed for reading the raw data exported from the service. Similarly, if you use existing experiment presentation software such as *Psychopy*, *OpenSesame*, *Psychtoolbox*, similar scripts might be available or the data is readily readable by `read.scv` function in R.

```
## READ DATA (Exp. 1 survey)
source('scr/read_data_survey.R')          # OUTPUT IS IN VARIABLE v
```

```
## N x Variables:42 132
```

```
source('scr/read_data_behavioural.R')      # OUTPUT IS IN VARIABLE B
```

```
## N x Variables:42 177
```

This should get variables `v` and `B`, which are data frames containing all the data. Next the raw responses are munged, that is, preprocessed in terms of the instruments (some items are reversed and then several items are combined). Missing values are imputed (with means), and when the items are combined, the reliability of the instruments is printed (Cronbach alpha). Finally the survey and self-report responses are combined into a single dataframe called `df`.

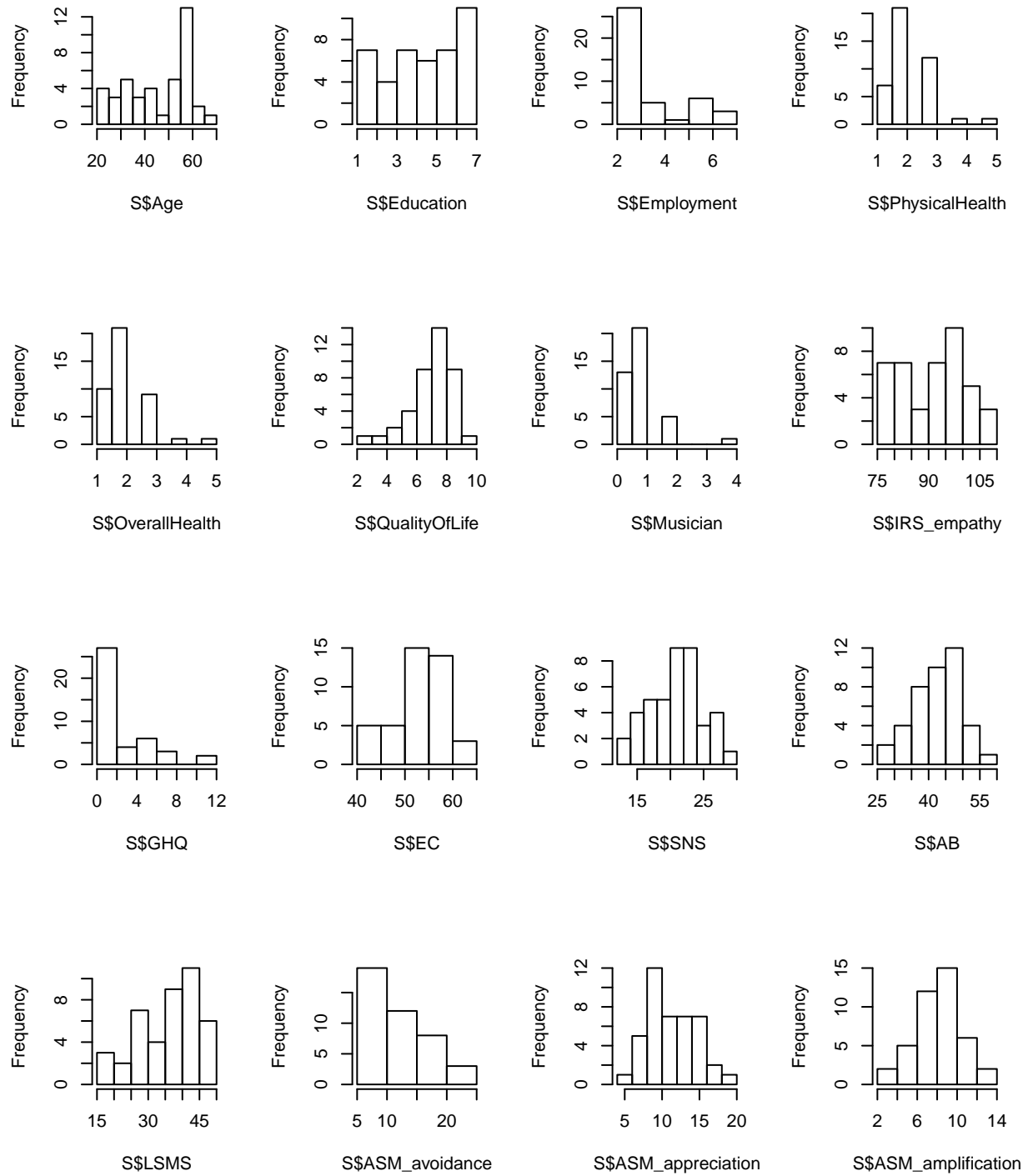
```
## MUNGE DATA (preprocess, recode, etc.)
source('munge/recode_instruments.R')       # Produces S data frame (Survey)
source('munge/merge_data.R')              # Combines behavioural & survey data
```

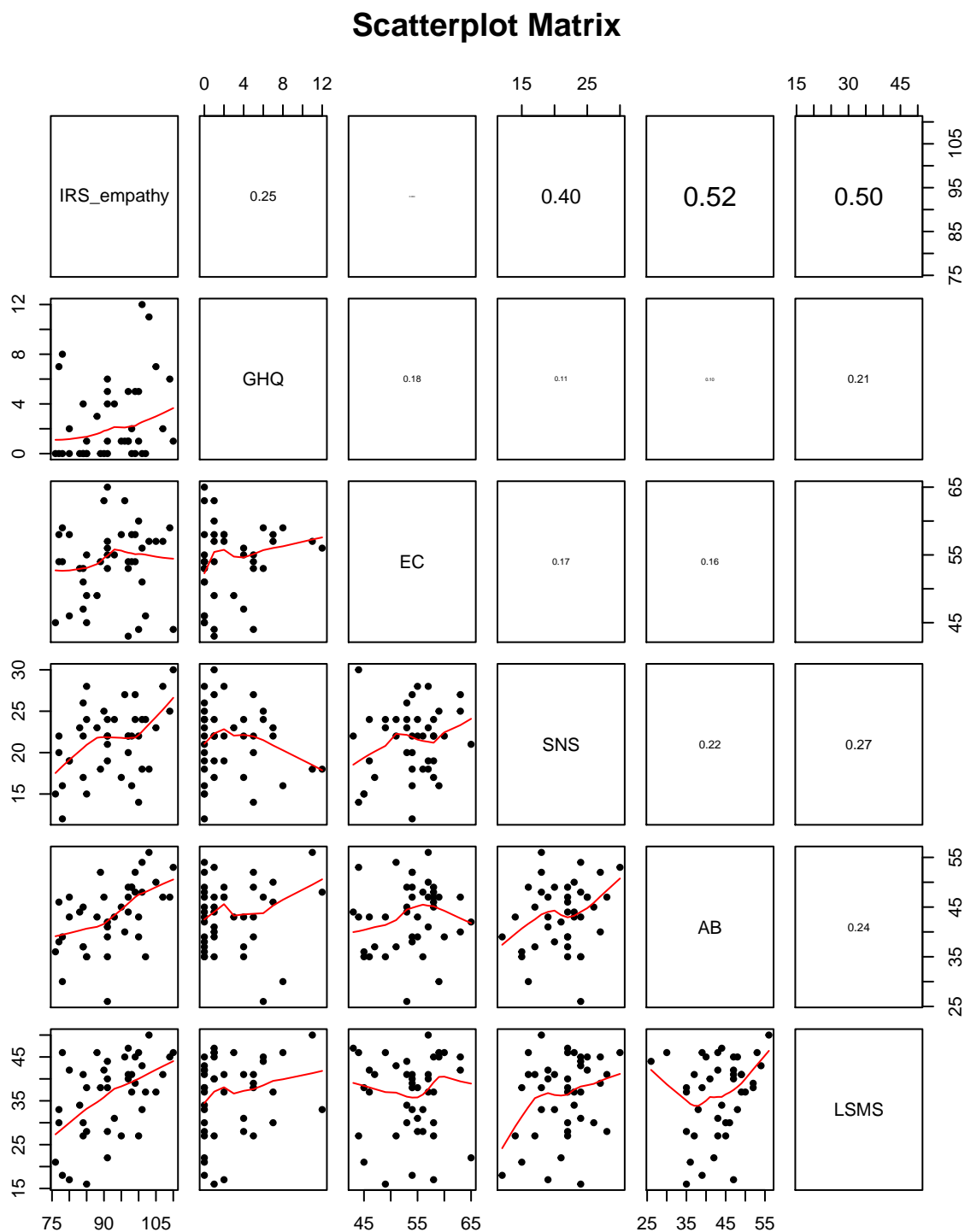
After the munging, it is prudent to check some aspects of the data, `N`, age, gender in order to remind of the dataset properties. It is also useful plot the distributions of the collected data in order to learn whether one needs to use certain operations (transformation or non-parametric statistics) in the subsequent analyses. Correlations between selected instruments is also shown.

```
## DIAGNOSE data
source('scr/diagnose_simple.R')            # describe the data (N, etc.)
```

```
## N = 42
## females = 26 males= 16
## Age range from 20 to 67
## M = 45.73171 SD = 13.96607
## Mean Age by Gender (female, male)
## 47.08 43.4
## Mean GHQ by Gender (female, male)
## 2.92 1.88
```

```
##
## Correlations between the Attitudes Towards Sad Music (ASM) factors
##      avoid. autob. reviv.  appr. inters.  ampl.
## avoid.   1.000 -0.203 -0.750 -0.275  -0.242 -0.013
## autob.  -0.203  1.000  0.347  0.334   0.272  0.253
## reviv.  -0.750  0.347  1.000  0.541   0.375  0.251
## appr.   -0.275  0.334  0.541  1.000   0.418  0.315
## inters. -0.242  0.272  0.375  0.418   1.000  0.124
## ampl.   -0.013  0.253  0.251  0.315   0.124  1.000
source('scr/diagnostic_plots.R')      # display histograms
```





If everything seems to be fine, it is time to proceed into the actual analysis.

3.2.3 Analyse the data

Finally we get to test the planned hypotheses of the experiment. Here the first analysis script builds a linear model for explaining the appreciation of sad music (`ASM_appreciation`)

with other background variables (emotional contagion, nostalgia-proneness, quality of life, musicianship, etc.) in a shape of linear regression. It also gives a simple example of exploring whether this same variable is different across gender and education. These are just examples which do not actually deliver deep answers to the questions.

```
source('scr/linear_models.R') # simple regression and ANOVA examples

## R.adj.squared = 0.4707871
## Analysis of Variance Table
##
## Response: ASM_appreciation
##
```

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
LSMS	1	138.832	138.832	24.0859	2.797e-05 ***
EC	1	23.697	23.697	4.1112	0.05127 .
SNS	1	26.951	26.951	4.6758	0.03842 *
QualityOfLife	1	36.474	36.474	6.3279	0.01728 *
Musician	1	4.343	4.343	0.7534	0.39206
GHQ	1	2.263	2.263	0.3925	0.53555
IRS_empathy	1	0.784	0.784	0.1360	0.71484
IRS_fantasy	1	12.748	12.748	2.2116	0.14708
Residuals	31	178.684	5.764		

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## ANOVA
##
```

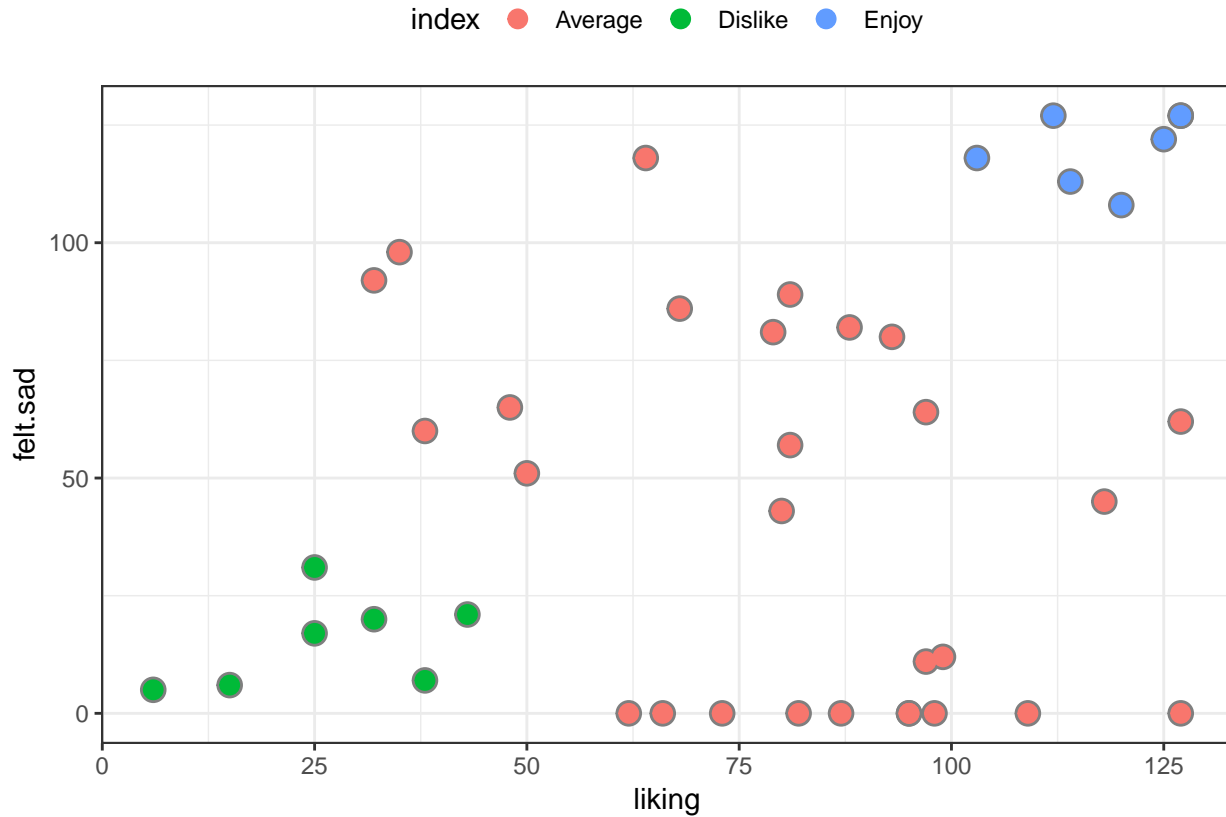
	Df	Sum Sq	Mean Sq	F value	Pr(>F)
Gender	1	0.4	0.423	0.037	0.849
Education	1	2.8	2.825	0.244	0.624
Gender:Education	1	0.6	0.572	0.050	0.825
Residuals	38	439.2	11.557		

Another idea is to attempt to predict those who have been highly influenced by the experiment in terms of giving high sadness ratings and high liking for these as well, whom we can call these “sad music lovers”. The opposite group, who disliked the music and do not report any sadness could be called “sad music haters”. We first define these by quantiles of the self-reports, then try to predict the members into three groups (lovers, neutrals, and haters) from the background variables using logistic regression.

```
source('scr/categorical_models.R') # predict groupings
```

Let’s visualise the values of the liking and feeling sad for the three groups. For clarity, it is a good idea to put the plots in separate scripts.

```
source('scr/figure1.R', print.eval=TRUE) # create Figure 1 for manuscript
```

One can also produce tables in the same way using a simple script. Here is an example of how to report the means and standard deviations of the background variables grouped according to the simple median split of the Liking for Sad Music (LSMS) variable. P-values from ANOVAs are thrown in as a bonus.

```
source('scr/table1.R')           # create Table 1 for manuscript
```

	Low (M)	High (M)	Low (SD)	High (SD)	Pval
## LSMS	29.409091	43.50	7.088952	2.910959	0.000
## Age	46.285714	45.15	14.269348	13.985989	0.798
## GHQ	2.272727	2.80	3.239368	3.122078	0.595
## IRS_fantasy	22.954545	25.25	4.961540	4.766937	0.135
## IRS_empathy	88.636364	95.75	8.796890	8.807144	0.012
## IRS_concern	25.363636	27.90	3.958661	2.693071	0.021
## IRS_perspective	22.772727	24.70	4.460748	4.414092	0.168
## IRS_distress	17.500000	17.95	4.983305	4.160908	0.754
## EC	53.590909	54.05	5.179254	5.898037	0.790
## SNS	20.500000	22.40	4.079566	4.057482	0.139
## AB	42.047619	44.60	5.499784	7.625442	0.225
## ASM_avoidance	15.136364	9.00	4.570042	2.772041	0.000
## ASM_revival	12.136364	17.45	4.632139	2.438183	0.000
## ASM_amplification	8.454545	8.85	2.344746	2.230766	0.580
## ASM_autobiographical	16.818182	18.00	3.540734	3.060788	0.256
## ASM_appreciation	10.454545	13.05	2.987709	3.119970	0.009
## ASM_appreciation	10.454545	13.05	2.987709	3.119970	0.009

Happy exploring. The intention is carrying out the analysis this way is to get a clear sense of the process and deliver outputs of the analyses that are easy to bring to the manuscript, and which also should be transparent for the other readers (supervisors and collaborators).

3.2.4 Combine report and analysis (optional)

It is also possible to combine the reporting of the analysis and the actual analysis to make the process even more transparent. An example of this can be found in `report.Rmd`, which basically runs the steps in the example template in sequence within a particular syntax (R `md`, using `knitr`, and this also creates a pdf or html report to the same folder (take a look at the `report.pdf`), which can contain all sorts of written arguments, comments and so on. In order to use this option, the `report.Rmd` should be on the project folder (not within reports subfolder) and your machine should have additional libraries installed (at least *LaTeX* and *knitr*, see [knitr in a nutshell](#)).

Update in December 2016: RStudio now supports **R Notebooks**, which combine the scripts and reports in a very nice fashion. See `report_NB.Rmd` for an example.

4 Help for statistics with R

4.1 Online tutorials

An Introduction to R The official guidance from *The Comprehensive R Archive Network* (CRAN). May not be always the most compelling introduction but exhaustive at least.

Quick-R Really good source of R examples for almost all operations (manipulation, representation, functions, syntax, stats, figures, etc.).

R tutorials Another fairly clear collection of tutorials.

RStudio online learning pages R Studio is fancy and great visual GUI on top of the R for all platforms and they have released very useful documentations, tutorials, demos, etc.

Advanced R Author of the best packages, Hadley Wickham, has created this resource (book and online version).

4.1.1 Statistics Handbooks with complete R scripts (online)

Practical Regression and Anova using R A handbook of the basic statistical operations written by Julian J. J. Faraway.

Data Analysis and Graphics Using R - An Example-Based Approach Handbook in 3rd printing, written by John Maindonald and John Braun. This source contains exercises, slides, the scripts for all graphs of the book, etc.

4.2 Other Online Resources

R blogger Multipurpose source for news and latest issues in R.

Collection of Resources at CRAN Large collection of different resources (e.g. R for Matlab-minded, Fitting Distributions with R, Reference Cards, Data-mining with R, and so on).

R Documentation Searchable online documentation

StackOverflow Forum of questions and answers about computer programming, including R. Contains over 40,000 questions related to R.

5 References

- Asendorpf, J. B., Conner, M., De Fruyt, F., De Houwer, J., Denissen, J. J., Fiedler, K., Fiedler, S., Funder, D. C., Kliegl, R., Nosek, B. A., & others (2013). Recommendations for increasing replicability in psychology. *European Journal of Personality*, 27(2), 108–119.
- Piwowar, H. A., Day, R. S., & Fridsma, D. B. (2007). Sharing detailed research data is associated with increased citation rate. *PloS one*, 2(3), e308.

- Sandve, G. K., Nekrutenko, A., Taylor, J., & Hovig, E. (2013). Ten simple rules for reproducible computational research. PLoS computational biology, 9(10), e1003285.
- Tomasello, M. & Call, J. (2011). Methodological challenges in the study of primate cognition. Science, 334(6060), 1227–1228.

This document is available in GitHub: <https://github.com/tuomaseerola/ReproR>