# ShinyItemAnalysis: Test and Item Analysis

**Acknowledgements**

# Preface

Educational and psychological testing is present in many areas of our everyday life, including assessing academic achievement, certifying qualifications and proficiency, or assessing one's fatigue, depression or pain. In many situations, testing is a key moment of people's life with long-reaching consequences, such as in university admissions or hiring of new employees. For these reasons, assessments that are used to measure ability, knowledge or other latent traits need to produce valid, reliable and fair scores.

While many methodological books exits that present the methodology for test validation, often practical examples of item and test analysis are missing, or are presented in commercial software which may be costly and thus unavailable. As an alternative, freely available statistical software `R` has been present for many years and many `R` packages have been developed to cover general psychometric concepts or specific psychometric topics. However for those who are new to `R`, it may be hard to overcome the initial burden of `R` code-based environment.

This book introduces selected topics in psychometrics and covers test and item validation with practical examples in `ShinyItemAnalysis` (Martinková, Drabinová, Leder, & Houdek, 2018). `ShinyItemAnalysis` is an `R` package and a freely available online application which provides interactive online environment to support teaching of psychometric concepts and test development with `R`, and to enforce routine validation of educational and psychological measurement worldwide.

The book is prepared as manual for `ShinyItemAnalysis` and explains methodology as well as practical features of each of its sections. It can thus serve those who use `ShinyItemAnalysis` - teachers who assess knowledge of their students, educators who develop new assessments, psychologists and researchers who use established or newly developed instruments in their projects, or even university stakeholders who want to introduce routine validation in their admission tests or classroom assessments. Inside the interactive environment of `ShinyItemAnalysis`, we introduce the reader also to examples in `R`programming language. Individual sections include Selected R code, as well as Exercises. The book can thus be very well used also in graduate courses of measurement and psychometrics and can serve as a gentle introduction to these topics in (or without) `R`.

Recent news can be found at www.ShinyItemAnalysis.org. The `ShinyItemAnalysis` application and `R` package were created with the aim to strengthen understanding of psychometric concepts and to support teaching of these concepts, to empower routine analysis of tests and also to present novel psychometric research. We hope the book you are reading will help fulfill these goals.

# Contents

# Chapter 1

# Introduction

## 1.1 Psychometrics, Measurement, Test development

Psychometrics is a field of study concerned with the theory and technique of psychological, educational and behavioral measurement. It is concerned with objective measurement (testing, assessment) of skills, knowledge, abilities, educational achievement, attitudes, personality traits and other.

As outlined in the Standards for educational and psychological testing AERA, APA and NCME 2014, assessments that are used to measure students' ability or knowledge need to produce valid, reliable and fair scores. To achieve these standards, many aspects of test development need to be taken care of (Haladyna & Downing, 2011). While core aspects of measurement are the same no matter the type of measurement, some specific topics may arise in different areas (see (Brennan, 2006) for overview on topic of Educational measurement).

Psychometric analysis and routine validation of tests is usually present in development of standardized tests, especially those used as admission tests to higher education (SAT, ACT, TOEFL, MCAT, BMAT, etc.), international large scale assessments (PISA, TIMSS, PIRLS), or annual testing of students performed in some countries and states (e.g. ...). Testing companies developing and administering these tests (e.g. College Board, ETS, etc.) often have departments or units taking care of item and test properties.

Test analysis is nowadays being more present also in regions and scientific areas where psychometrics does not have a long tradition. Complex test analyses can be found in development and validation of conceptual assessments (see e.g McFarland et al., 2017) - tests of students' conceptual understanding of key topics in certain fields.

## 1.2 Software

This book introduces selected topics in psychometrics with practical examples in `ShinyItemAnalysis` (Martinková et al., 2018), which provides gentle introduction to psychometric analyses. It uses powers of the many psychometric `R` packages in user-friendly interface. Many interactive features are present to support understanding of presented concepts.

## 1.3 Book overview

Several topics which we find most important parts of test analysis, "base stones", are covered in individual chapters. **Chapter 2** helps the reader to get started with the software.

**Chapter 3** offers introduction to measurement data.

**Chapter 4** provides introduction to measurement error and various ways to get proofs of reliability of the test.

**Chapter 5** describes analyses which may help to provide proofs of test validity.

In **Chapter 6**, traditional item analysis is provided, including various item characteristics based on ratios or correlations, mainly describing item difficulty, discrimination power, guessing or response rate. Detailed distractor analysis is presented to provide better understanding to functioning of all offered options in multiple-choice tests for low as well as high performing students.

**Chapter 7** introduces regression models for description of item properties.

**Chapter 8** explains various item response theory models for binary, ordinal as well as nominal data describing tests including dichotomous, partial credit, Likert-scale, multiple-choice and other types of items.

**Chapter 9** covers topic of differential item functioning and presents various methods for detection of DIF.

Further topics not yet covered by the `ShinyItemAnalysis` application are described in **Chapter 10**.

**Chapter 11** describes how reports may be generated with `ShinyItemAnalysis`.

**Appendices** provide detailed guidance about how to install R, LaTeX, etc.

In summary, the `ShinyItemAnalysis` software as well as this book cover the basic topics in psychometrics and test analysis to provide a solid base stone. Moreover, by introducing these concepts in R, we aim to open the door to much wider and rich methodology for quantitative analysis in education, psychology and sociology to advance quantitative methodology in the behavioral sciences.

# Chapter 2

# Getting started

## 2.1 Getting started with ShinyItemAnalysis

`ShinyItemAnalysis` online application is available at

https://shiny.cs.cas.cz/ShinyItemAnalysis

Other mirrors are specified at http://www.ShinyItemAnalysis.org. As we discuss below, it is also possible to run the application locally.

Intro page (Figure 2.1) includes general information about the application. Various tools are included in separate tabs with logical ordering into separate sections.
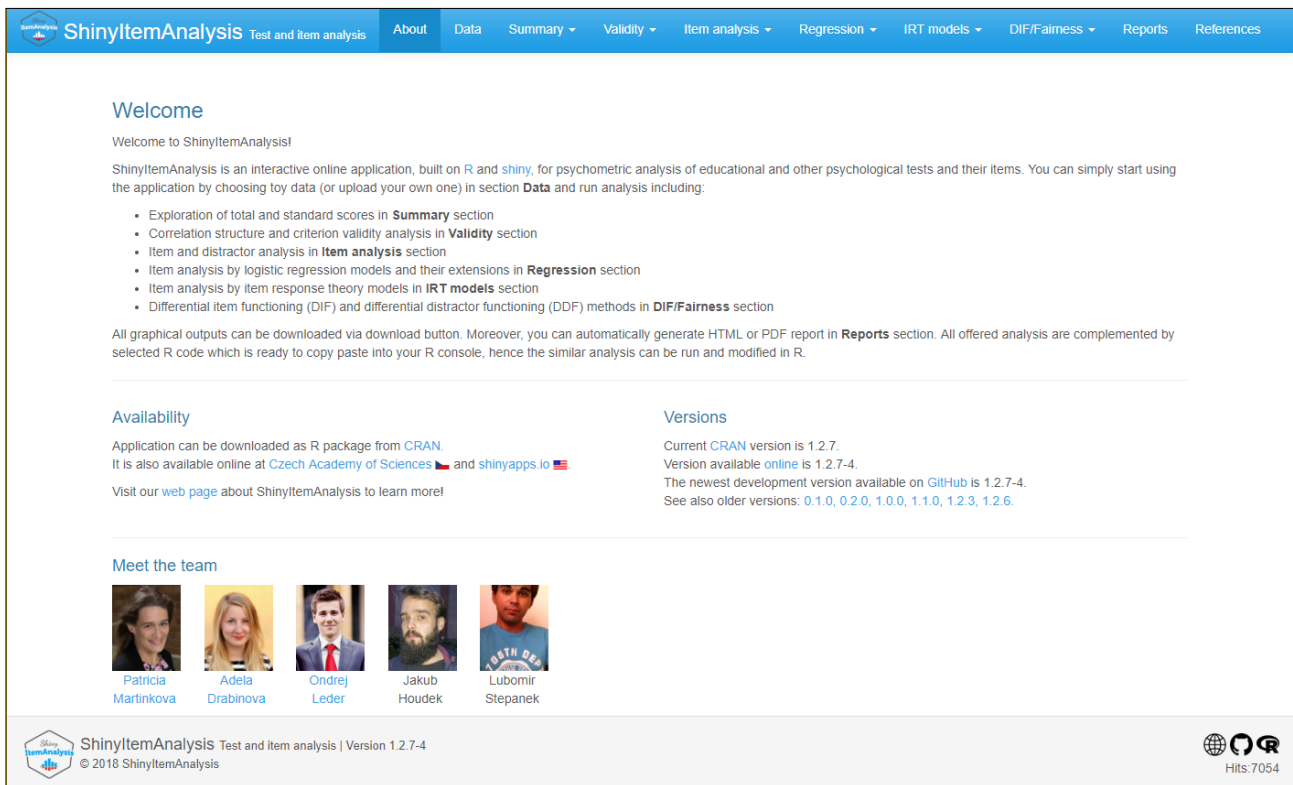


Figure 2.1: Intro page.

## 2.2 Getting started with R

While you may read this book and try most of the exercises without even opening `R`, to get most of the book, we recommend to also download and use `R`. Two main benefits you get from using `R` are as follows: First, you may then use the application locally, which may become faster and more efficient than the online version. Second,

besides running the online application locally in your `R` , you may instead use `R` console to try `R` examples provided inside `ShinyItemAnalysis` and modify them for your purposes.

R can be installed from

https://cran.r-project.org/

Detailed instructions can be found in appendix. [Consider mentioning RStudio.] Basic introductory of `R` can be given e.g. by Paradis (2002).

Once you have your `R` installed, it is easy to install the last stable version of `ShinyItemAnalysis` package from **CRAN** by writing the following code

```
install.packages("ShinyItemAnalysis")
```

The newest development version can be instead downloaded from **GitHub** (with `devtools` package) using the following lines

```
devtools::install_github("patriciamar/ShinyItemAnalysis")
```

Once the `ShinyItemAnalysis` package is installed it can be loaded and the online application can be run locally:

```
library(ShinyItemAnalysis)
startShinyItemAnalysis()
```

Main function `startShinyItemAnalysis()` launches an interactive application as described above.

`ShinyItemAnalysis` uses several `R` packages to provide wide palette of psychometric tools to analyze data (see Table 2.1). Overview of many other psychometric libraries is provided on Psychometric CRAN Task (Mair, 2018).

One can also use `R` console to try selected `R` code provided in `ShinyItemAnalysis`, or to modify the code as needed.

Table 2.1: `R` packages used for developing `ShinyItemAnalysis`.

| R package | Citation | Title |
| --- | --- | --- |
| corrplot | (Wei & Simko, 2017) | Visualization of a correlation matrix |
| CTT | (Willse, 2018) | Classical test theory functions |
| data.table | (Dowle & Srinivasan, 2017) | Extension of `data.frame` |
| deltaPlotR | (Magis & Facon, 2014) | Identification of dichotomous differential item functioning using Angoff's delta plot method |
| difNLR | (Drabinová, Martinková, & Zvára, 2018) | DIF and DDF detection by non-linear regression models |
| difR | (Magis, Beland, Tuerlinckx, & De Boeck, 2010) | Collection of methods to detect dichotomous differential item functioning |
| DT | (Xie, 2018) | A wrapper of the `JavaScript` library 'datatables' |
| ggplot2 | (Wickham, 2016) | Create elegant data visualisations using the grammar of graphics |
| gridExtra | (Auguie, 2017) | Miscellaneous functions for "grid" graphics |
| knitr | (Xie, 2015) | A general-purpose package for dynamic report generation in `R` |
| lattice | (Sarkar, 2008) | Trellis graphics for `R` |
| latticeExtra | (Sarkar & Andrews, 2016) | Extra graphical utilities based on lattice |
| lme4 | (Bates, Mächler, Bolker, & Walker, 2015) | Linear mixed-effects models using `Eigen` and S4 |
| ltm | (Rizopoulos, 2006) | Latent trait models under IRT |
| MASS | (Venables & Ripley, 2002) | Support functions and datasets for Venables and Ripley's MASS |
| Matrix | (Bates & Maechler, 2017) | Sparse and dense matrix classes and methods |
| mirt | (Chalmers, 2012) | Multidimensional item response theory |
| moments | (Komsta & Novomestky, 2015) | Moments, cumulants, skewness, kurtosis and related tests |
| msm | (Jackson, 2011) | Multi-state Markov and hidden Markov models in continuous time |
| multilevel | (Bliese, 2016) | Multilevel functions |
| nlme | (Pinheiro, Bates, DebRoy, Sarkar, & R Core Team, 2017) | Linear and nonlinear mixed effects models |
| nnet | (Venables & Ripley, 2002) | Feed-forward neural networks and multinomial log-linear models |
| plotly | (Sievert et al., 2017) | Create interactive web graphics via 'plotly.js' |
| polycor | (Fox, 2016) | Polychoric and polyserial correlations |
| psych | (Revelle, 2018) | Procedures for psychological, psychometric, and personality research |
| psychometric | (Fletcher, 2010) | Applied psychometric theory |
| RColorBrewer | (Neuwirth, 2014) | Colorbrewer palettes |
| reshape2 | (Wickham, 2007) | Flexibly reshape data: A reboot of the reshape package |
| rmarkdown | (Allaire et al., 2017) | Dynamic documents for `R` |
| shiny | (Chang, Cheng, Allaire, Xie, & McPherson, 2017) | Web application framework for `R` |
| shinyBS | (Bailey, 2015) | Twitter bootstrap components for shiny |
| shinydashboard | (Chang & Borges Ribeiro, 2018) | Create dashboards with shiny |
| shinyjs | (Attali, 2018) | Easily improve the user experience of your shiny apps in seconds |
| stringr | (Wickham, 2018) | Simple, consistent wrappers for common string operations |
| WrightMap | (Irribarra & Freund, 2014) | IRT item-person map with 'conquest' integration |
| xtable | (Dahl, 2016) | Export tables to LaTeX or `HTML` |

# Chapter 3

# Measurement data

Throughout the book, we will be working with several measurement datasets. `ShinyItemAnalysis` also allows the users to upload their own datasets. These mostly contain responses of students (in rows) to test items (columns). Responses may be binary (i.e., true/false, or 1/0), ordinal (e.g., on Likert scale $1-2-3-4-5$) or nominal (e.g., $A-B-C-D$ for multiple-choice items). [consider mentioning mix of item formats]

Some further respondent covariates may be present, e.g. group membership: gender, ethnicity, etc. Besides, some criterion variable may be present, such as repondents' IQ, their future study success, study Grade Point Average (GPA), etc.

## 3.1 Toy data

In `ShinyItemAnalysis`, five training datasets may be uploaded using the **Select dataset** button in section **Data**:

**GMAT** is a simulated dataset from `ShinyItemAnalysis` R package. The dataset represents responses of 2,000 subjects (1,000 males, 1,000 females) to multiple-choice test of 20 items. The answers were generated using parameters of real Graduate Management Admission Test (GMAT) (Kingston, Leary, & Wightman, 1985). The distribution of total scores is the same for both groups. However, first two items were manipulated to function differently for the two groups. See Martinková et al. (2017) for further discussion. GMAT dataset also containts simulated continuous criterion variable.

Similarly, **GMAT2** is a simulated dataset based on parameters of real GMAT (Kingston et al., 1985) from `difNLR` R package (Drabinová et al., 2018). The dataset represents responses of 1,000 subjects (500 males, 500 females) to multiple-choice test of 20 items. Also in this dataset, the first two items were simulated to function differently in uniform and non-uniform way respectively.

**Medical 100** is a real dataset of admission test to medical school from `ShinyItemAnalysis` R package. The data set represents responses of 2,392 subjects (750 males, 1,633 females and 9 subjects without gender specification) to multiple-choice test of 100 items. Medical 100 contains criterion variable - indicator whether student remained in the study after one year or not.

**MSAT-B** is a subset of real Medical School Admission Test in Biology (MSAT-B) in Czech Republic from `difNLR` R package (Drabinová et al., 2018). The dataset represents responses of 1,407 subjects (484 males, 923 females) to selection of 20 multiple-choice items. First item was previously detected as functioning differently for the two genders. For more details on this dataset, see Drabinová and Martinková (2017).

**HCI** (McFarland et al., 2017) is a real dataset of Homeostasis Concept Inventory (HCI) offered by R package `ShinyItemAnalysis`. The dataset represents responses of 651 subjects (405 males, 246 females) to multiple-choice test of 20 items. HCI contains criterion variable - indicator whether student plans to major in the life sciences.

## 3.2 Data upload

Own data may be uploaded as csv files and previewed in the **Data** section. Main data file should contain responses of individual respondents (rows) to given items (columns). Data need to be either binary or nominal (e.g., in A-B-C-D format). Ordinal data (such as Likert scale) are currently treated as nominal.

Individual items need to be separated by comma, semicolon or tab, this is specified in check box "Separator". Around each reponse value, double-quote or quote may be present (this is often typical in nominal data and it may be specified using check box "Quote"). Header may contain item names, no row names should be included. In all data sets header should be either included or excluded, this is specified by check box "Header". Columns of dataset are by default renamed to Item and number of particular column, however, keeping original names

may be forced using check box "Keep items names". Missing values in scored dataset are by default evaluated as 0, however treating them as missing may be forced using check box "Keep missing values".

For nominal data, it is necessary to upload key of correct answers. Group vector may also be included, this is a binary vector, where 0 represents reference group and 1 represents focal group. Its length need to be the same as number of individual respondents in main dataset. Missing values are not supported for group membership vector and should be removed. Finally, criterion variable can be included, this is either discrete or continuous vector (e.g. future study success or future GPA in case of admission tests). Again, its length needs to be the same as number of individual respondents in the main dataset.

To replicate examples involving HCI dataset (McFarland et al., 2017), csv files for upload are provided on [include link]



Figure 3.1: Page to select or upload data.

## 3.3 Data summary and exploration

Data inspection is a crucial first step in any data analysis. Summary statistics should always be checked before proceeding to further analyses. Summary tab offers basic summaries of data including counts in nominal and binary categories in items, counts for groups and basic statistics for criterion variable (Figure 3.2).

Further checks of any suspicious data may be done in Data exploration tab (Figure 3.3).

## 3.4 Total scores

Total score also known as raw score or sum score is a total number of correct answers. In what follows we label total score of person $p$ as $X_p$. Let $\bar{X} = \frac{1}{n} \sum_{p=1}^{n} X_p$ be a sample mean of total scores $X_p$ and $s^2 = \frac{1}{n-1} \sum_{p=1}^{n} \left( X_p - \bar{X} \right)^2$ their sample variance. Z-score or also standardized score is a linear transformation of total score with a mean of 0 and with variance of 1, that is Z-score for person $p$ is given by

$$Z_p = \frac{X_p - \bar{X}}{s}.$$

T-score is transformed Z-score with a mean of 50 and standard deviation of 10, that is

$$T_p = 10Z_p + 50.$$

Section **Summary** offers summary statistics and histogram of the total scores. The summary table (Table 3.1) offers their basic characteristics such as minimum and maximum, mean, median, standard deviation, skewness
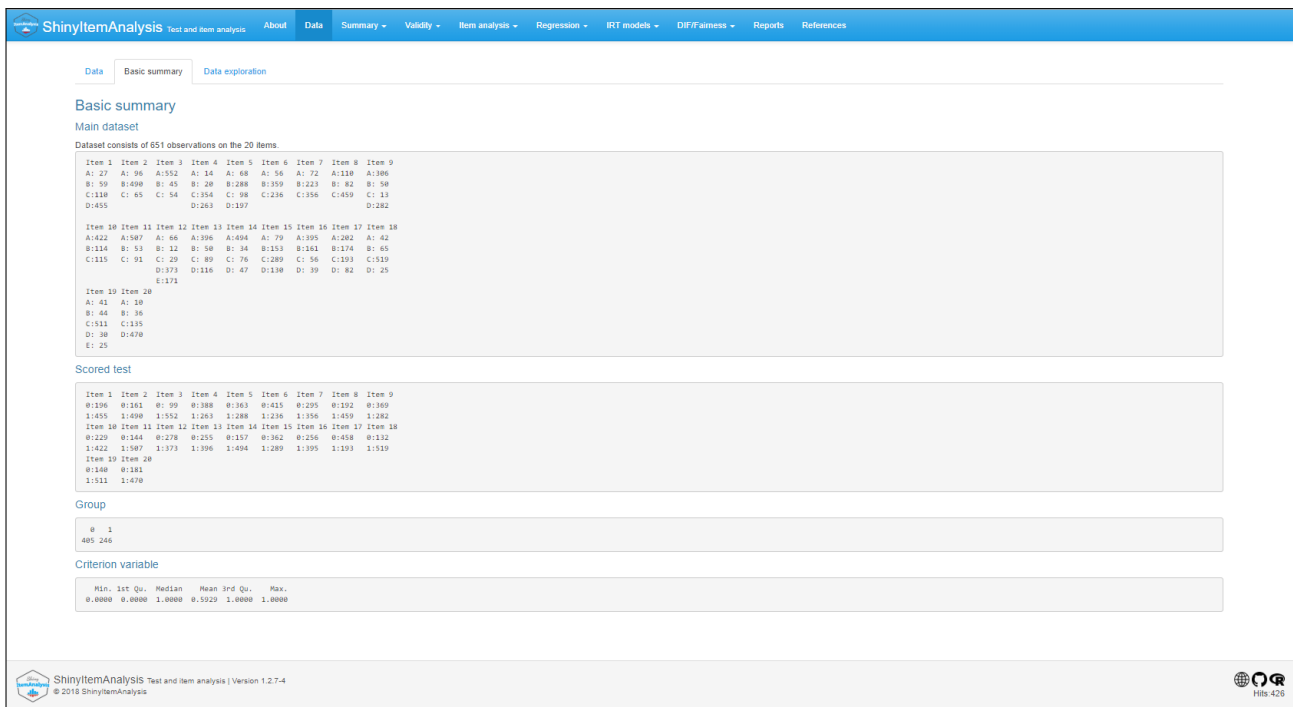
Figure 3.2: Page of basic data summary.

Table 3.1: Summary table of total scores for HCI dataset

| Min | Max | Mean | Median | SD | Skewness | Kurtosis |
|------|-------|-------|--------|------|----------|----------|
| 3.00 | 20.00 | 12.21 | 12.00 | 3.64 | $-0.20$ | 2.35 |

and kurtosis. The kurtosis here is estimated by sample kurtosis

$$g_2 = \frac{m_4}{s^4} = \frac{\frac{1}{n}\sum_{p=1}^{n}\left(X_p - \bar{X}\right)^4}{\left[\frac{1}{n}\sum_{p=1}^{n}\left(X_p - \bar{X}\right)^2\right]^2}.$$

The skewness is estimated by sample skewness

$$b_1 = \frac{m_3}{s^3} = \frac{\frac{1}{n}\sum_{p=1}^{n}\left(X_p - \bar{X}\right)^3}{\left[\frac{1}{n-1}\sum_{p=1}^{n}\left(X_p - \bar{X}\right)^2\right]^{3/2}}.$$

The kurtosis for normally distributed scores is near the value of 3 and the skewness is near the value of 0.

Besides the summary statistics the histogram of total scores is provided to describe the distribution of total scores (Figure 3.4). The cut-score may be specified to better visualize the distribution of total scores in two groups (e.g. of those who passed a test and those who did not). For selected cut-score, blue part of histogram shows respondents with total score above the cut-score, grey column shows respondents with total score equal to the cut-score and red part of histogram shows respondents below the cut-score. Bell-shaped histograms are typical for normally distributed data. On the other hand, two-peaked histograms may signalize that the data is actually composed out of two different subgroups.

Total scores with various standard scores (Z-scores, T-scores) are summarized, together with percentile and success rate for each level of total score in one table (Table 3.2).

From numbers in Table 3.2, we can for example read that students with 16 points were in 87th percentil with 80% success rate.
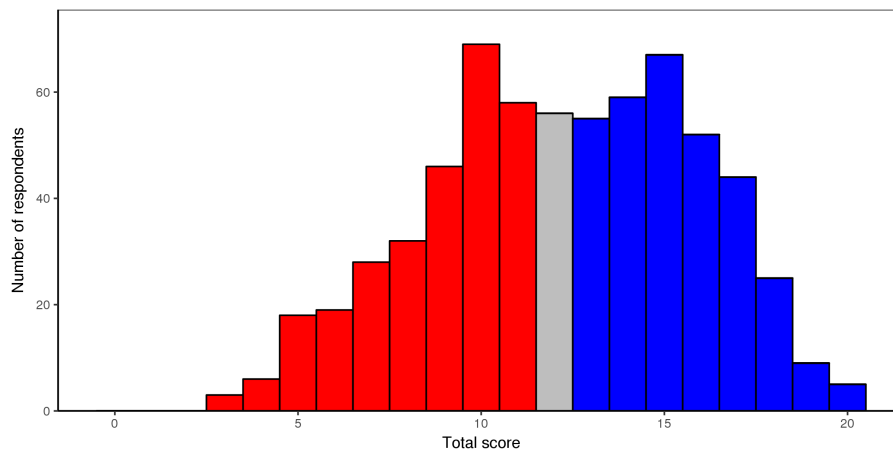
Figure 3.3: Data exploration tab.



Figure 3.4: Histogram of total scores for HCI data.

## 3.5 Selected R code

```r
library(difNLR)
library(ggplot2)
library(moments)

# loading data
data(GMAT)
data <- GMAT[, 1:20]

# total score calculation
score <- apply(data, 1, sum)

# summary of total score
c(min(score), max(score), mean(score), median(score), sd(score),
  skewness(score), kurtosis(score))

# colors by cut-score
```

Table 3.2: Standard scores for HCI dataset.

| Total score | Percentile | Success rate | Z-score | T-score |
|---|---|---|---|---|
| 3.00 | 0.00 | 15.00 | $-2.53$ | 24.69 |
| 4.00 | 0.01 | 20.00 | $-2.26$ | 27.44 |
| 5.00 | 0.04 | 25.00 | $-1.98$ | 30.19 |
| 6.00 | 0.07 | 30.00 | $-1.71$ | 32.93 |
| 7.00 | 0.11 | 35.00 | $-1.43$ | 35.68 |
| 8.00 | 0.16 | 40.00 | $-1.16$ | 38.43 |
| 9.00 | 0.23 | 45.00 | $-0.88$ | 41.18 |
| 10.00 | 0.34 | 50.00 | $-0.61$ | 43.92 |
| 11.00 | 0.43 | 55.00 | $-0.33$ | 46.67 |
| 12.00 | 0.51 | 60.00 | $-0.06$ | 49.42 |
| 13.00 | 0.60 | 65.00 | 0.22 | 52.17 |
| 14.00 | 0.69 | 70.00 | 0.49 | 54.91 |
| 15.00 | 0.79 | 75.00 | 0.77 | 57.66 |
| 16.00 | 0.87 | 80.00 | 1.04 | 60.41 |
| 17.00 | 0.94 | 85.00 | 1.32 | 63.16 |
| 18.00 | 0.98 | 90.00 | 1.59 | 65.90 |
| 19.00 | 0.99 | 95.00 | 1.87 | 68.65 |
| 20.00 | 1.00 | 100.00 | 2.14 | 71.40 |

```r
cut <- median(score) # cut-score
color <- c(rep("red", cut - min(score)),
           "gray",
           rep("blue", max(score) - cut))
df <- data.frame(score)

# histogram
ggplot(df, aes(score)) +
  geom_histogram(binwidth = 1, fill = color, col = "black") +
  xlab("Total score") +
  ylab("Number of respondents") +
  theme_app()

# scores calculations
score <- apply(data, 1, sum) # Total score
tosc <- sort(unique(score)) # Levels of total score
perc <- cumsum(prop.table(table(score))) # Percentiles
sura <- 100 * (tosc / max(score)) # Success rate
zsco <- sort(unique(scale(score))) # Z-score
tsco <- 50 + 10 * zsco # T-score
```

## 3.6   Exercises

**Ex. 3.1**   Run `ShinyItemAnalysis` and try basic data exploration. Using default dataset, answer following questions.

- What is its name?
- How many observations does dataset consist of?
- How many observations do come from focal and reference group?
- What are the maximum and minimum values of criterion variable?

**Ex. 3.2**   Upload data into `ShinyItemAnalysis` and explore them.

- What is mean and standard deviation of total scores?
- Calculate Z-score for student with total score 10. Provide whole calculation.

- Calculate T-score for student with total score 10. Provide whole calculation.
- How many points did student with 90th percentile receive?

**Ex. 3.3** Create short `R` script including following tasks.

- Upload data from previous section.
- Calculate total scores for uploaded dataset, their mean, median and their standard deviation.
- Draw histogram of total scores. Values smaller than median should be red, values larger than median should be blue, median should be gray.
- Calculate Z-score for uploaded dataset.
- Calculate T-score for uploaded dataset.

# Chapter 11

# Reports generation

To support routine usage of psychometric methods in test development, `ShinyItemAnalysis` offers possibility to upload data for analysis as csv files, and to generate PDF or HTML reports. Sample PDF report and csv files used for its generation are provided in Supplemental Materials.

Report generation uses `rmarkdown` templates and `knitr` for compiling (see Figure 11.1). LaTeX is used for creating PDF reports. Latest version of LaTeX with properly set paths is needed to generate PDF reports locally.



Figure 11.1: Report generation workflow.

Report page setting allows to specify dataset name, to include name of person who generated the report, to select from available methods and to customize settings (see Figures 11.2 and 11.3). **Generate report** button starts analyses needed for report generation. Subsequently, **Download report** button initializes compiling the text, tables and figures into PDF or HTML file.



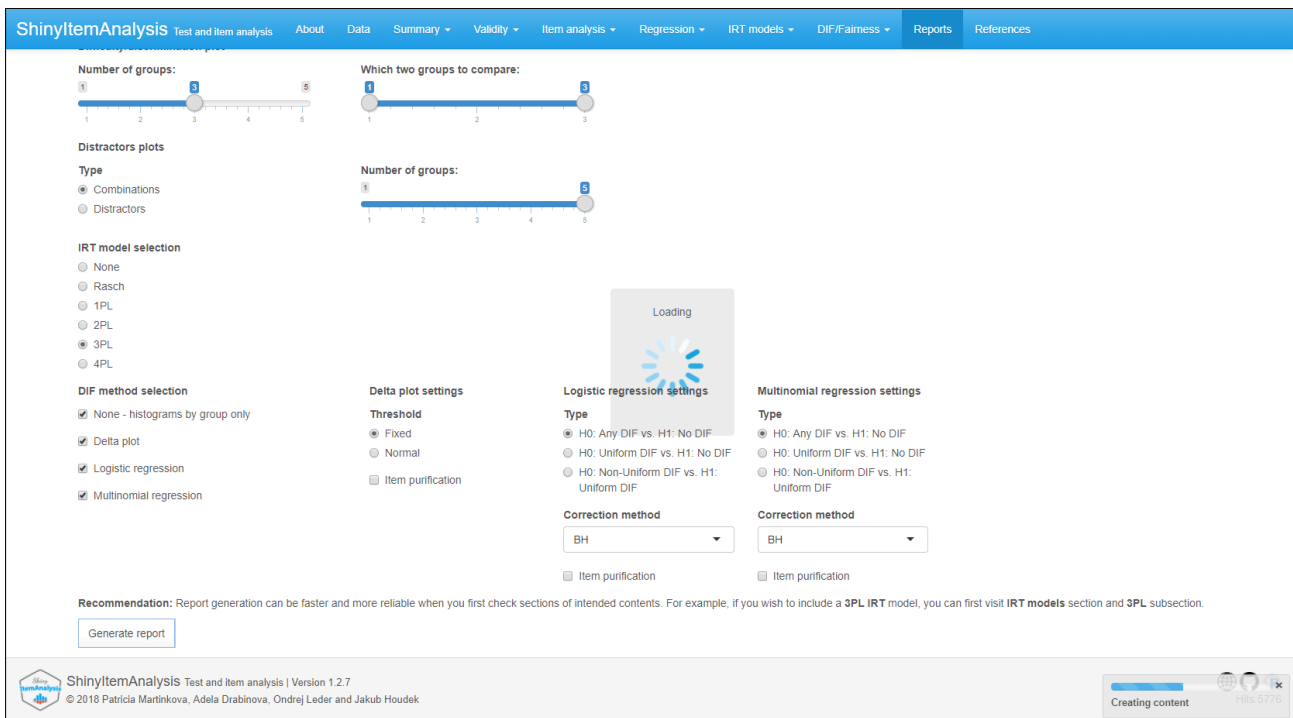Figure 11.2: Report setting of HCI data analysis.

Figure 11.3: Report setting of HCI data analysis.

Sample pages of PDF report on HCI dataset are displayed in Figure 11.4. Reports provide quick overview of test characteristics and may be a helpful material for test developers, item writers and institutional stakeholders.
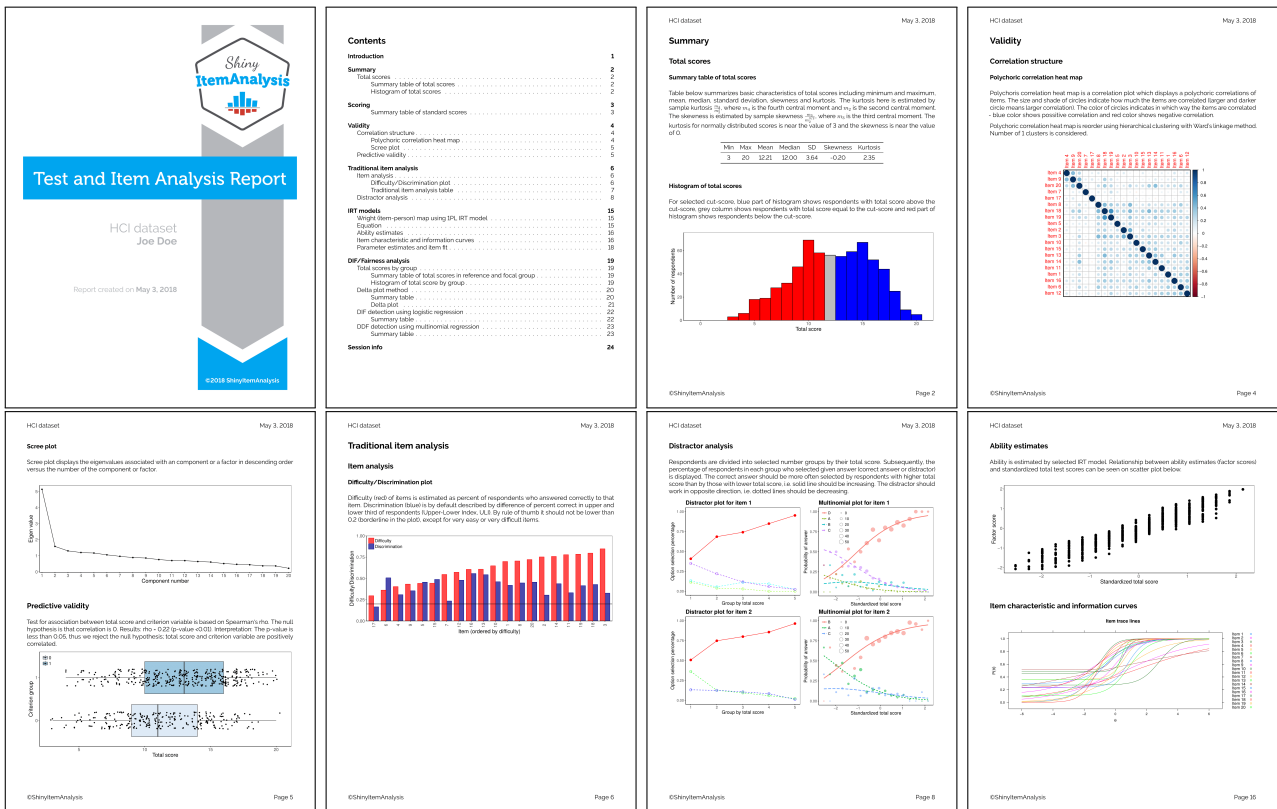
Figure 11.4: Selected pages of report on HCI data.

# Appendices

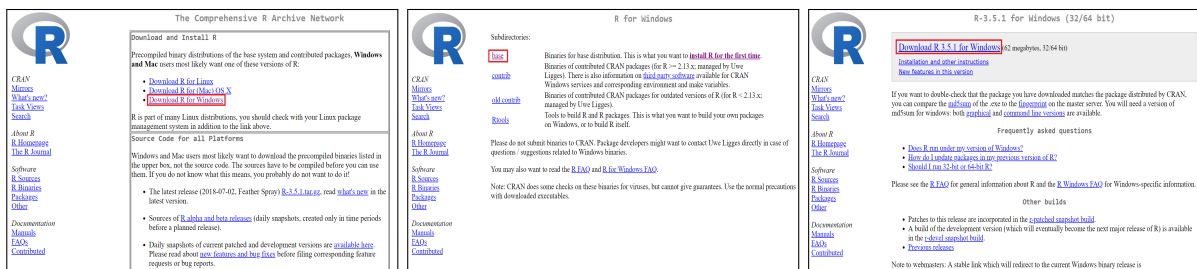## A1    Installation of `R` and RStudio

Here we provide detailed instruction for installation of `R` in Windows and Mac OS X. We also recommend you to install RStudio.

### A1.1    Windows

1. Go to

   https://cran.r-project.org/

   and click on **Download `R` for Windows**, then **base** and finally **Download `R` 3.x.x for Windows**.
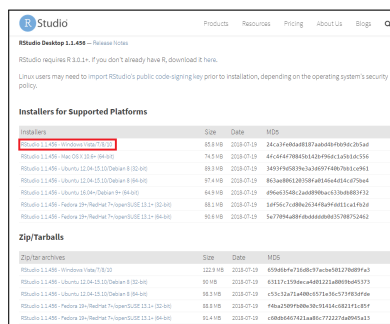
   

   This starts downloading of installer R-3.X.X-win.exe.
2. Open downloaded installer and follow instructions to install `R`. Leave all default settings in the installation options.
3. Go to

   https://www.rstudio.com/products/rstudio/download/

   and click on **RStudio 1.1.XXX - Windows Vista/7/8/10** on the bottom of the page.

   

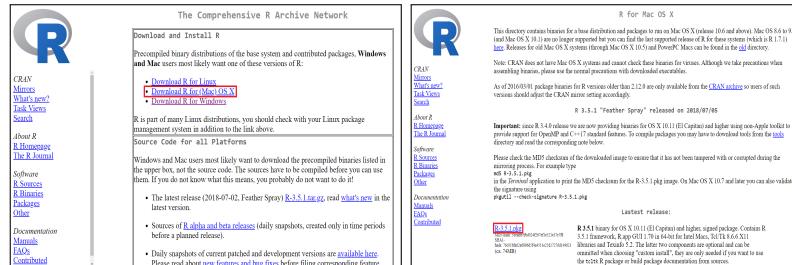   This starts downloading of installer RStudio-1.1.XXX.exe.
4. Open downloaded installer and follow instructions to install RStudio. Leave all default settings in the installation options.

## A1.2 Mac OS X

1. Go to

https://cran.r-project.org/

and click on **Download R for Mac OS X** and then **Download R-3.x.x.pkg**.

2. Install R. Leave all default settings in the installation options.

3. Go to

https://www.rstudio.com/products/rstudio/download/

and click on **RStudio 1.1.XXX - Mac OS X 10.6+ (64-bit)** on the bottom of the page.

4. Install RStudio by dragging the application icon to your Applications folder.

## A2 Installation of `ShinyItemAnalysis`

1. Open RStudio (or R) and install and load `ShinyItemAnalysis` with typing following commands into console:

```
install.packages("ShinyItemAnalysis")
library(ShinyItemAnalysis)
```

2. In case that some dependency packages have not been installed automatically, you can use command

```
install.packages("MISSING-PACKAGE")
```

where MISSING-PACKAGE is replaced by name of not installed package. To install packages, you can use also clickable environment of RStudio:

3. Now, `ShinyItemAnalysis` is ready to run. To launch the application, type into console

```
startShinyItemAnalysis()
```



Click on **Open in browser** button to open application in your favourite browser.

# A3   Installation of T<sub>E</sub>Xdistribution

Here we provide links to detailed tutorials for MiKTex installation. MiKTeX is T<sub>E</sub>Xdistribution and can be freely downloaded:

https://miktex.org/download

Please, follow tutorials for your choice of operation system at this webpage.

# References

Akaike, H. (1974). A new look at the statistical model identification. *IEEE Transactions on Automatic Control*, *19*(6), 716–723.

Allaire, J., Xie, Y., McPherson, J., Luraschi, J., Ushey, K., Atkins, A., . . . Chang, W. (2017). rmarkdown: Dynamic documents for r [Computer software manual]. Retrieved from https://CRAN.R-project.org/package=rmarkdown (R package version 1.8)

American Educational Research Association (AERA), American Psychological Association (APA), National Council on Measurement in Education (NCME). (2014). *Standards for educational and psychological testing*. American Educational Research Association.

Ames, A. J., & Penfield, R. D. (2015). An ncme instructional module on item-fit statistics for item response theory models. *Educational Measurement: Issues and Practice*, *34*(3), 39–48.

Andrich, D. (1978). A rating formulation for ordered response categories. *Psychometrika*, *43*(4), 561–573.

Angoff, W. H., & Ford, S. F. (1973). Item-race interaction on a test of scholastic aptitude. *Journal of Educational Measurement*, *10*(2), 95–105.

Attali, D. (2018). shinyjs: Easily improve the user experience of your shiny apps in seconds [Computer software manual]. Retrieved from https://CRAN.R-project.org/package=shinyjs (R package version 1.0)

Auguie, B. (2017). gridextra: Miscellaneous functions for "grid" graphics [Computer software manual]. Retrieved from https://CRAN.R-project.org/package=gridExtra (R package version 2.3)

Bailey, E. (2015). shinybs: Twitter bootstrap components for shiny [Computer software manual]. Retrieved from https://CRAN.R-project.org/package=shinyBS (R package version 0.61)

Bates, D., Mächler, M., Bolker, B., & Walker, S. (2015). Fitting linear mixed-effects models using lme4. *Journal of Statistical Software*, *67*(1), 1–48. doi: 10.18637/jss.v067.i01

Bates, D., & Maechler, M. (2017). Matrix: Sparse and dense matrix classes and methods [Computer software manual]. Retrieved from https://CRAN.R-project.org/package=Matrix (R package version 1.2-12)

Bliese, P. (2016). multilevel: Multilevel functions [Computer software manual]. Retrieved from https://CRAN.R-project.org/package=multilevel (R package version 2.6)

Bock, R. D. (1972). Estimating item parameters and latent ability when responses are scored in two or more nominal categories. *Psychometrika*, *37*(1), 29–51.

Brennan, R. L. (2006). *Educational measurement*. Praeger.

Brown, W. (1910). Some experimental results in the correlation of mental abilities. *British Journal of Psychology, 1904-1920*, *3*(3), 296–322.

Chalmers, R. P. (2012). mirt: A multidimensional item response theory package for the R environment. *Journal of Statistical Software*, *48*(6), 1–29. doi: 10.18637/jss.v048.i06

Chang, W., & Borges Ribeiro, B. (2018). shinydashboard: Create dashboards with 'shiny' [Computer software manual]. Retrieved from https://CRAN.R-project.org/package=shinydashboard (R package version 0.7.0)

Chang, W., Cheng, J., Allaire, J., Xie, Y., & McPherson, J. (2017). shiny: Web application framework for r [Computer software manual]. Retrieved from https://CRAN.R-project.org/package=shiny (R package version 1.0.5)

Cronbach, L. J. (1951). Coefficient alpha and the internal structure of tests. *Psychometrika*, *16*(3), 297–334.

Dahl, D. B. (2016). xtable: Export tables to latex or html [Computer software manual]. Retrieved from https://CRAN.R-project.org/package=xtable (R package version 1.8-2)

Dowle, M., & Srinivasan, A. (2017). data.table: Extension of 'data.frame' [Computer software manual]. Retrieved from https://CRAN.R-project.org/package=data.table (R package version 1.10.4-3)

Drabinová, A., & Martinková, P. (2017). Detection of differential item functioning with nonlinear regression: A non-IRT approach accounting for guessing. *Journal of Educational Measurement*, *54*(4), 498–517.

Drabinová, A., & Martinková, P. (2018). difnlr: Generalized logistic regression models for dif and ddf detection. *R Journal*. (Submitted)

Drabinová, A., Martinková, P., & Zvára, K. (2018). difnlr: Dif and ddf detection by non-linear regression models. [Computer software manual]. Retrieved from https://CRAN.R-project.org/package=difNLR (R package version 1.2.2)

Ebel, R. L. (1954). Procedures for the analysis of classroom tests. *Educational and Psychological Measurement*, *14*(2), 352–364.

Fletcher, T. D. (2010). psychometric: Applied psychometric theory [Computer software manual]. Retrieved from https://CRAN.R-project.org/package=psychometric (R package version 2.2)

Fox, J. (2016). polycor: Polychoric and polyserial correlations [Computer software manual]. Retrieved from https://CRAN.R-project.org/package=polycor (R package version 0.7-9)

Haladyna, T. M., & Downing, S. M. (2011). *Handbook of test development*. Routledge.

Irribarra, D. T., & Freund, R. (2014). Wright map: Irt item-person map with conquest integration [Computer software manual]. Retrieved from http://github.com/david-ti/wrightmap

Jackson, C. H. (2011). Multi-state models for panel data: The msm package for R. *Journal of Statistical Software*, *38*(8), 1–29. Retrieved from http://www.jstatsoft.org/v38/i08/

Kingston, N., Leary, L., & Wightman, L. (1985). An exploratory study of the applicability of item response theory methods to the graduate management admission test. *ETS Research Report Series*, *1985*(2), 1–56.

Komsta, L., & Novomestky, F. (2015). moments: Moments, cumulants, skewness, kurtosis and related tests [Computer software manual]. Retrieved from https://CRAN.R-project.org/package=moments (R package version 0.14)

Lord, F. M. (1980). *Applications of item response theory to practical testing problems*. Routledge.

Magis, D., Beland, S., Tuerlinckx, F., & De Boeck, P. (2010). A general framework and an r package for the detection of dichotomous differential item functioning. *Behavior Research Methods*, *42*, 847–862.

Magis, D., & Facon, B. (2014). deltaPlotR: An R package for differential item functioning analysis with angoff's delta plot. *Journal of Statistical Software, Code Snippets*, *59*(1), 1–19. Retrieved from http://www.jstatsoft.org/v59/c01/

Mair, P. (2018). *CRAN task view: Psychometric models and methods*. Retrieved 2018-08-16, from https://CRAN.R-project.org/view=Psychometrics

Mantel, N., & Haenszel, W. (1959). Statistical aspects of the analysis of data from retrospective studies of disease. *Journal of the national cancer institute*, *22*(4), 719–748.

Martinková, P., Drabinová, A., Leder, O., & Houdek, J. (2018). ShinyItemAnalysis: Test and item analysis via shiny [Computer software manual]. Retrieved from https://CRAN.R-project.org/package=ShinyItemAnalysis (R package version 1.2.6)

Martinková, P., Drabinová, A., Liaw, Y.-L., Sanders, E. A., McFarland, J. L., & Price, R. M. (2017). Checking equity: Why differential item functioning analysis should be a routine part of developing conceptual assessments. *CBE-Life Sciences Education*, *16*(2), rm2.

Martinková, P., Štěpánek, L., Drabinová, A., Houdek, J., Vejražka, M., & Štuka, v. (2017). Semi-real-time analyses of item characteristics for medical school admission tests. In *Computer science and information systems (fedcsis), 2017 federated conference on* (pp. 189–194).

Masters, G. N. (1982). A rasch model for partial credit scoring. *Psychometrika*, *47*(2), 149–174.

McFarland, J. L., Price, R. M., Wenderoth, M. P., Martinková, P., Cliff, W., Michael, J., . . . Wright, A. (2017). Development and validation of the homeostasis concept inventory. *CBE-Life Sciences Education*, *16*(2), ar35.

Muraki, E. (1992). A generalized partial credit model: Application of an em algorithm. *ETS Research Report Series*, *1992*(1).

Neuwirth, E. (2014). Rcolorbrewer: Colorbrewer palettes [Computer software manual]. Retrieved from https://CRAN.R-project.org/package=RColorBrewer (R package version 1.1-2)

Nunnally, J. C., & Bernstein, I. (1994). *Psychometric theory* (2nd ed.). McGraw-Hill New York.

Paradis, E. (2002). *R for beginners*. Montpellier (F): University of Montpellier. Retrieved from https://cran.r-project.org/doc/contrib/Paradis-rdebuts_en.pdf

Pinheiro, J., Bates, D., DebRoy, S., Sarkar, D., & R Core Team. (2017). nlme: Linear and nonlinear mixed effects models [Computer software manual]. Retrieved from https://CRAN.R-project.org/package=nlme (R package version 3.1-131)

Raju, N. S. (1988). The area between two item characteristic curves. *Psychometrika*, *53*(4), 495–502.

Raju, N. S. (1990). Determining the significance of estimated signed and unsigned areas between two item response functions. *Applied Psychological Measurement*, *14*(2), 197–207.

Rasch, G. (1960). *Studies in mathematical psychology: I. probabilistic models for some intelligence and attainment tests*. Nielsen & Lydiche.

Revelle, W. (1979). Hierarchical cluster analysis and the internal structure of tests. *Multivariate Behavioral Research*, *14*(1), 57–74.

Revelle, W. (2018). psych: Procedures for psychological, psychometric, and personality research [Computer software manual]. Evanston, Illinois. Retrieved from https://CRAN.R-project.org/package=psych (R package version 1.8.3)

Rizopoulos, D. (2006). ltm: An r package for latent variable modelling and item response theory analyses. *Journal of Statistical Software*, *17*(5), 1–25. Retrieved from http://www.jstatsoft.org/v17/i05/

Rosseel, Y. (2012). lavaan: An R package for structural equation modeling. *Journal of Statistical Software*, *48*(2), 1–36. Retrieved from http://www.jstatsoft.org/v48/i02/

Rust, J., & Golombok, S. (2014). *Modern psycometrics* (3rd ed.). Routledge.

Samejima, F. (1970). Estimation of latent ability using a response pattern of graded scores. *Psychometrika*, *35*(1), 139–139.

Sarkar, D. (2008). *Lattice: Multivariate data visualization with r.* New York: Springer. Retrieved from http://lmdvr.r-forge.r-project.org (ISBN 978-0-387-75968-5)

Sarkar, D., & Andrews, F. (2016). latticeextra: Extra graphical utilities based on lattice [Computer software manual]. Retrieved from https://CRAN.R-project.org/package=latticeExtra (R package version 0.6-28)

Schwarz, G., et al. (1978). Estimating the dimension of a model. *The Annals of Statistics*, *6*(2), 461–464.

Sievert, C., Parmer, C., Hocking, T., Chamberlain, S., Ram, K., Corvellec, M., & Despouy, P. (2017). plotly: Create interactive web graphics via 'plotly.js' [Computer software manual]. Retrieved from https://CRAN.R-project.org/package=plotly (R package version 4.7.1)

Spearman, C. (1910). Correlation calculated from faulty data. *British Journal of Psychology, 1904-1920*, *3*(3), 271–295.

Swaminathan, H., & Rogers, H. J. (1990). Detecting differential item functioning using logistic regression procedures. *Journal of Educational measurement*, *27*(4), 361–370.

van der Linden, W. J. (2017). *Handbook of item response theory.* CRC Press.

Venables, W. N., & Ripley, B. D. (2002). *Modern applied statistics with s* (Fourth ed.). New York: Springer. Retrieved from http://www.stats.ox.ac.uk/pub/MASS4 (ISBN 0-387-95457-0)

Wei, T., & Simko, V. (2017). R package "corrplot": Visualization of a correlation matrix [Computer software manual]. Retrieved from https://github.com/taiyun/corrplot ((Version 0.84))

Wickham, H. (2007). Reshaping data with the reshape package. *Journal of Statistical Software*, *21*(12), 1–20. Retrieved from http://www.jstatsoft.org/v21/i12/

Wickham, H. (2016). *ggplot2: Elegant graphics for data analysis.* Springer-Verlag New York. Retrieved from http://ggplot2.org

Wickham, H. (2018). stringr: Simple, consistent wrappers for common string operations [Computer software manual]. Retrieved from https://CRAN.R-project.org/package=stringr (R package version 1.3.0)

Willse, J. T. (2018). Ctt: Classical test theory functions [Computer software manual]. Retrieved from https://CRAN.R-project.org/package=CTT (R package version 2.3.2)

Xie, Y. (2015). *Dynamic documents with R and knitr* (2nd ed.). Boca Raton, Florida: Chapman and Hall/CRC. Retrieved from https://yihui.name/knitr/ (ISBN 978-1498716963)

Xie, Y. (2018). Dt: A wrapper of the javascript library 'datatables' [Computer software manual]. Retrieved from https://CRAN.R-project.org/package=DT (R package version 0.4)

Zinbarg, R. E., Revelle, W., Yovel, I., & Li, W. (2005). Cronbach's $\alpha$, revelle's $\beta$, and mcdonald's $\omega$ h: Their relations with each other and two alternative conceptualizations of reliability. *psychometrika*, *70*(1), 123–133.

# Acronyms

**DDF** Differential Distractor Functioning.

**GMAT** Graduate Management Admission Test.
**GPA** Grade Point Average.
**GPCM** Generalized Partial Credit Model.
**GRM** Graded Response Model.

**HCI** Homeostasis Concept Inventory.

**ICC** Item Characteristic Curve.
**IIC** Item Information Curve.

**MSAT-B** Medical School Admission Test in Biology.

**NRM** Nominal Response Model.

**PCM** Partial Credit Model.

**RSM** Rating Scale Model.

**ULI** Upper-Lower Index.

# Index