

GEORG-AUGUST UNIVERSITÄT GÖTTINGEN

STATISTICAL PROGRAMMING WITH R

Gotta Read 'Em All: An RStudio Add-In to visually read different file-formats into R

Author:

Stanislaus STADLMANN,
Student ID: 21144637

Supervisor

Paul WIEMANN, M.Sc.

Submitted on August 22, 2016



GEORG-AUGUST-UNIVERSITÄT
GÖTTINGEN

Contents

1	Motivation	1
2	Underlying Frameworks	2
2.1	The Shiny Framework	2
2.2	Shiny Gadgets and RStudio Add-Ins	4
3	Implementation	4
3.1	The Main Function: GREYA()	5
3.2	The File Reading Function: GREYA_read()	7
3.3	Helper Functions	9
4	Usage	9
4.1	Requirements	9
4.2	Installation	9
4.3	Starting the Add-In	10
4.4	Selecting the Dataset	10
4.5	The preview window	11
4.6	Adjusting reading conditions	11
	References	13

List of Figures

1	Screenshot of an example Shiny Application.	3
2	Selecting a file	10
3	Looking at the preview window	11
4	A ready-to use R command created by GREASE	11

Listings

1	Example code for a Shiny Application[2].	3
2	Structure of the <code>GREASE()</code> function	5
3	Contents of reactive function for “Done”-event	6
4	Reading procedure of <code>GREASE_read()</code>	8
5	Return Command of the <code>GREASE_read()</code> function	9
6	Installation of GREASE	10

1 Motivation

R is a statistical software with an almost uncountable number of functions for different statistical methods, procedures and graphs. The Comprehensive R Archive Network (CRAN), the most popular network for adding new features to R via so called “packages”, has more than 8000 different packages ready to be downloaded¹. Each of them provide a variety of functions to solve different tasks. For example, a package called “Vector Generalized Linear and Additive Models” (or short VGAM) can be easily installed via the R command `install.packages("VGAM")` and, once loaded in R, provides functions to estimate a variety of different regression models.

These aforementioned packages with the underlying functions make R the popular statistical package it is today. But most of these additional features require data to be of any use, most prominently regression model estimation functions. Datasets do exist shipped with R but for most academic purposes, external data will be required to generate new insights.

There are multiple ways of reading data into R depending on the data type (e.g. .csv, .xlsx) and also the data size (big data, small data). Reading a .csv file, for example, can be achieved via the built-in R function `base::read.table("filename.csv")`. Big .csv files can be read very quickly with the function `data.table::fread("bigfile.csv")` from the data.table package [4, Dowle et al 2015]. Other packages provide even more functionality, e.g. for dealing with strings or a smaller number of required arguments inside of a function. Most of those packages are also available on CRAN.

The availability of packages to read different filetypes in numerous ways is very helpful for the advanced R user, because there is almost no filetype that cannot be read via an R function. But it also poses a problem: if there are so many ways that a user can read a file into R, how will she/he remember all the necessary packages and functions, and their function arguments? This problem is often encountered by new R users who want to use R’s extended functionalities but fail at importing data into their working environment.

An answer to this problem is provided by Thomas Leeper’s R package called “rio” [5, Leeper et al 2016], which tries to minimize redundancy by wrapping

¹There were exactly 8895 packages on CRAN at August 4, 2016.

reading functions into one import `rio::import()` and one export function `rio::export()`.

The R package introduced with this paper takes it one step further. Built on the Shiny Framework and implemented as an RStudio Add-In, “Gotta Read ’Em All” provides a General User Interface (GUI) for reading all different file-formats into R.

The general process is the following: in the beginning, the user selects a file on her/his computer. After some adjustments (which are done interactively), the proper function to read the file is pasted into the console, with an object name that can be specified by the user. In between, the user can always head to the preview to see what the parsed file would look like with the current options.

Using this Add-In, the user can read data into R without remembering any code, but still obtains the correct R code to re-parse the data at a later point.

2 Underlying Frameworks

2.1 The Shiny Framework

The Shiny Framework² is in itself an R package designed to create interactive visualisations with R functions and HTML code. The author describes the package as “combining the computational power of R with the interactivity of the modern web” [6, RStudio 2016]. The goal is to create applications with clickable interfaces quickly showcasing different scenarios. This is done via reactive R functions, which are run everytime a user interacts with the GUI.

Figure 1 shows an example Shiny Application, consisting of some user interface (UI) control elements (a select button, and tick boxes) and a graph. The UI is reactive; whenever the user ticks a box or selects a different value in the first box, the graph changes its look. This is built upon reactive functions, which run every time a value inside them changes. The values that are allowed to change are therefore bound to the UI elements.

Shiny Apps consist of two elements: The UI and a “server” side. The UI includes functions which wrap HTML to build the viewable part of the App, for

²<http://shiny.rstudio.com>

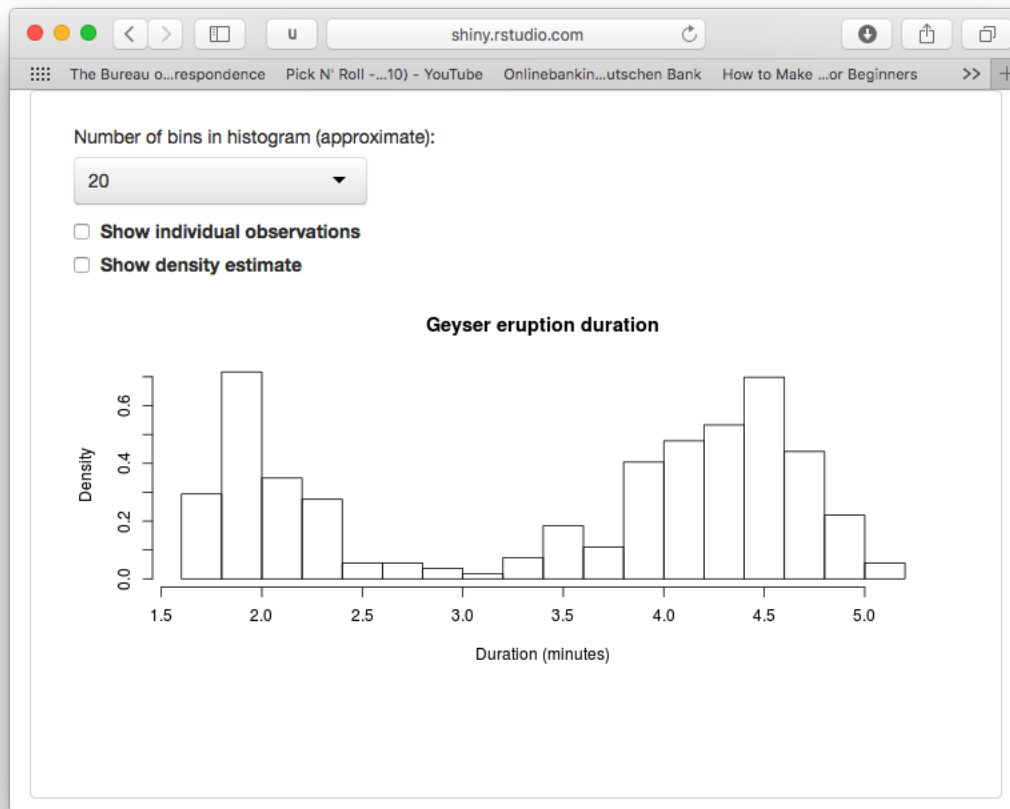


Figure 1: Screenshot of an example Shiny Application.

example buttons and placement of graphs. The server side consists of functions that specify the reactive R functions which create dynamic output displayed on the UI side, and when the functions should be run. Both sides then are able to automatically communicate with each other to create a smooth interactive experience for the user. Code for an example Shiny app is attached below:

```

1 # Define UI
2 ui <- bootstrapPage(
3   numericInput("n", "Number of obs", 100),
4   plotOutput("plot")
5 )
6 # Define Server
7 server <- function(input, output) {
8   output$plot <- renderPlot({ hist(runif(input$n)) })
9 }
10 app <- shinyApp(ui, server)
11 # Run App
12 runApp(app)

```

Code-Chunk 1: Example code for a Shiny Application[2].

It is now possible to see in Code-Chunk 1 that both the UI and server elements are R objects, while the server also resembles an R function. `runApp()` then opens up the Application.

2.2 Shiny Gadgets and RStudio Add-Ins

Shiny Applications are useful for displaying interactive visualisations, but they are made for displaying results to the end user. “Shiny Gadgets”³, an extension of the Shiny framework, are supposed to be part of the programming or analysing process. They are built on the same framework that was introduced with Shiny, but serve the purpose of making programming challenges a little easier. For example, a Shiny Gadget could be used to provide a UI for downloading certain data from complex websites.

“RStudio Add-Ins”⁴ are Shiny Gadgets that are built right into RStudio, an Integrated Developer Environment (IDE) for R. Calling the Shiny Gadget is made easier, as the RStudio user only has to press two buttons. Furthermore, Add-Ins have extended access to RStudio itself via a package called “rstudioapi” [8, Wickham and Allaire 2016]. For example, Add-Ins are able to paste a string into the console and can modify the currently opened R script.

The combination of the Shiny framework and RStudio add-ins creates the ideal setting for an Add-In that helps the user parse any data format into R. The Shiny Framework provides interactiveness and the RStudio connection makes it easier to call the Application out of an IDE.

3 Implementation

An R-Studio Add-In has to be installed via the R package ecosystem, so GREa is also wrapped up in a package called GREa. Calling the Add-In is done via the main function `GREa::GREa()`. Also, there exist a couple of helper functions, which were necessary to reduce redundant R code. The following functions are implemented:

- `GREa()`

³<http://shiny.rstudio.com/articles/gadgets.html>

⁴<https://rstudio.github.io/rstudioaddins/>

- `GREA_read()`
- `wd_check()`
- `fileChoose()`

As stated above, `GREA::GREA()` is the function that starts the Add-In. `GREA::GREA_read()` is the most important helper function, which converts any filetype into the `data.frame` R class inside the global environment. For checking if a file is inside the current R working directory, `GREA::wd_check()` was written. This function was created from my personal frustration with large filepaths. `GREA::fileChoose()` has the same utility as `base::file.choose()` (choosing a file interactively), except that it returns `NULL` when cancelled.

3.1 The Main Function: `GREA()`

As explained in Chapter 2.2, Add-Ins are built on the Shiny Framework. `GREA()` is therefore also a wrapper for a Shiny Application. In this section, the general structure of the function will be demonstrated.

```

1 GREA <- function() {
2   ui <- miniPage(
3     # User Interface Functions #
4     # ...
5   )
6   server <- shinyServer(function(input, output, session) {
7     # Reactive Server Functions #
8     # ...
9     observeEvent(input$done, {
10      # Functions that paste the code into the console
11    })
12  })
13  # Functions that start the Add-In
14  app <- shinyApp(ui = ui, server = server)
15  viewer <- dialogViewer(dialogName = "GREA",
16                        height = 350, width = 500)
17  runGadget(app, viewer = viewer, stopOnCancel = FALSE)
18 }
```

Code-Chunk 2: Structure of the `GREA()` function

Code-Chunk 2 shows the structure of GREA’s main function. Similarly to Shiny Applications, we can see that the `ui`(lines 2-5) and `server`(lines 6-12) objects are specified first. For creating the `ui` object, the `miniUI::miniPage()` function from the `miniUI` package [3, Cheng 2016] is used, which is specifically targeted at small user interfaces in the style of smartphone apps. The `server` object, also resembling an R function, consists of reactive functions that are called every time the user interacts with the UI. This includes the selection of a file to be read in on the user’s computer, or the adjustment of reading conditions (e.g. a change in the comma separator in text-delimited files).

Next, both objects are assembled to a new object called `app`(line 14). In lines 15 and 16, an object named `viewer` is generated via `shiny::dialogViewer()`. This step ensures that GREA is started as an integrated window in RStudio with the right height and width. In the end, the function `runGadget()`(line 17) is called to start the Add-In. In essence, this function has the same functionality as the `shiny::runApp()` function, with the addition of better automatic handling in the event of the user cancelling the app.

The piece of code that sets this Add-In apart from Shiny applications and Shiny Gadgets is resembled by lines 9-11 of Code-Chunk 2, a function which is run when the user presses the “done” button.

```

1 observeEvent(input$done, {
2   # Paste Code into Console
3   if (nzchar(fileloc()) && nzchar(input$name_dataset) && !
4     is.null(dataset())) {
5     # Get code that was used to read dataset
6     expr <- attributes(dataset())$GREACommand
7     # Assemble code
8     code <- paste0(input$name_dataset, " <- ", expr)
9     # Paste into Console
10    rstudioapi::insertText(text = code, id = "#console")
11  }
12  # ... and then stop the app
13  stopApp()
14 })

```

Code-Chunk 3: Contents of reactive function for “Done”-event

Code-Chunk 3 shows the detailed contents of lines 9-11 in Code-Chunk 2. After

the user has specified all necessary variables interactively, she/he presses the “done” button. This triggers the event called `input$done`, which leads to the above function being called.

Line 3 checks if three conditions are met:

1. A file location is correctly specified.
2. A dataset name is correctly specified.
3. Reading the dataset with the options provided by the user is successful.

These conditions make sure that the user doesn’t obtain code for reading a file which will yield an error. If and only if these conditions are met, the code to read the data is assembled (lines 5-7) and then pasted into the R console (line 9). The procedure to paste the code into the console makes use of the `insertText()` function from the “rstudioapi” package, previously mentioned in Chapter 2.2. After this procedure, the user obtains the code to read the specified file and only has to execute the command to attach the new data to R’s global environment.

3.2 The File Reading Function: `GREA_read()`

Chapter 3.1 explained the implementation of the main function, `GREA::GREA()`. Though this function is the cornerstone of the overall Add-In, it relies heavily on a function that reads filetypes into R, `GREA::GREA_read()`. Writing this function was challenging, as it had to fulfill the following requirements at once:

1. Read a filetype and convert it to an R object.
2. Keep the code after a successful reading and attach it to the R object.
3. Omit rarely used arguments in the code that should be pasted into the console (e.g. `NA` values).

Requirement 1 is the goal that many R reading functions want to achieve. Requirement 2 is necessary, because after successful interactive parsing, the right code should be pasted into the active R console. Requirement 3 was included, because only the reading conditions that are not default values should be included in the code that the user obtains in the end.

To fulfill Requirement 1, `GREA_read()` automatically detects the filetype of the

selected file. Then, depending on which filetype was detected, it utilizes the following functions from other packages to read data:

- `base::read.table()` for text-delimited files (e.g. .csv, .txt, etc),
- `R.matlab::readMat()` for .mat files (MATLAB) [1, Bengtsson 2016],
- and `rio::import()` for all other filetypes.

After detecting the filetype and selecting the right function for reading, a “call” is assembled. In general, an R call is an executed or non-executed command. In this case, the command is created depending on the arguments that the user specified. To show the procedure, an example code snippet is provided below.

```
1 else if (any(filetype == c("xls", "xlsx"))) {  
2     expr <- quote(rio::import())  
3     expr[c("file", "which")] <- list(filelocation,  
4         sheetIndex)  
5     if (!missing(na.values))  
6         expr[c("na")] <- na.values  
7     if (skip > 0)  
8         expr[c("skip")] <- skip  
9 }
```

Code-Chunk 4: Reading procedure of `GREA_read()`

In Line 1 of Code-Chunk 4, the condition for the selected file being of an Excel type (.xls, .xlsx) is examined. Once the condition is met, Line 2 utilizes the `base::quote()` function, which creates a non-executed call of the first argument, in this case `rio::import()`. In Line 3, the arguments `file` (filelocation) and `which` (which of the Excel sheets should be parsed) are specified using the arguments of `GREA::GREA_read()`. Lines 4-7 fulfill Requirement 3 which was specified at the beginning of this chapter: For rare function arguments, these arguments should only be included when they are explicitly specified. In this case, arguments `na` and `skip` are only then added to the call (object `expr`) when they deviate from default values (`na.values` is not a missing argument and `skip` is greater than 0).

After the call is assembled correctly, it should not only be executed (Requirement 1), but it should also be attached to the object that was created during the call (Requirement 2). This is achieved by the following code at the end of the function:

```
1 | return(structure(eval(expr), GREACommand = deparse(expr)))
```

Code-Chunk 5: Return Command of the `GREA_read()` function

Code-Chunk 5 shows the execution of said assembled call. The `eval()` command first executes it and the `structure()` function then attaches the non-executed call (`expr`) to the object that was created, thus fulfilling Requirements 1 and 2. Afterwards, the enhanced object is returned (including its attribute).

Whenever a file is now read into the R environment via `GREA_read()`, the code that was used can be accessed.

3.3 Helper Functions

The GREA package features two helper functions: `fileChoose()` and `wd_check()`. As mentioned before, `fileChoose()` is just an enhanced version of `base::file.choose()` that returns `NULL` if the interactive file selection is cancelled. `wd_check` has two tasks. First, it transforms Windows filepaths to a path that can be read by any Operating System. Second, the function checks if a file selected via `fileChoose()` lies inside of the current R working environment, thus being able to shorten the filepath while reading files (that is where it got its name from). The returned object is just another filepath, which is only a subpath if the file lies inside the current R working directory.

4 Usage

4.1 Requirements

RStudio Add-Ins require the newest release of RStudio⁵ and R⁶.

4.2 Installation

The Add-In was uploaded to GitHub (an open-source code sharing platform), so it can be easily installed. Installation is done via the following code:

⁵RStudio can be downloaded here: <https://www.rstudio.com>.

⁶R can be downloaded here for all major platforms: <https://www.r-project.org>.

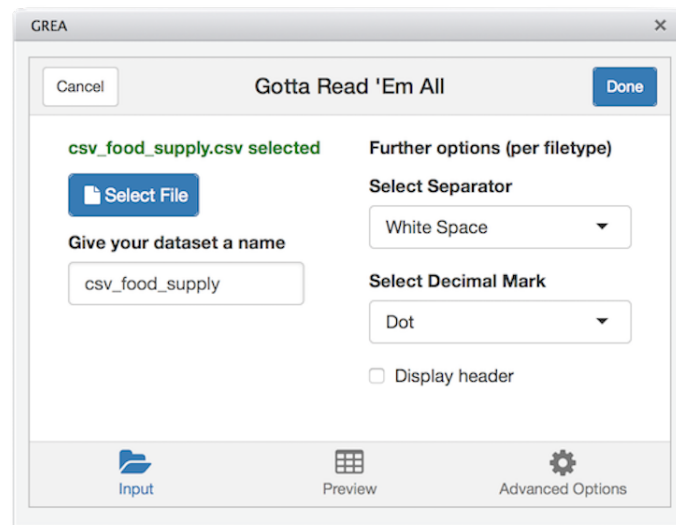


Figure 2: Selecting a file

```

1 | if (!require(devtools))
2 |   install.packages("devtools")
3 | devtools::install_github("Stan125/GREA")

```

Code-Chunk 6: Installation of GREY

Lines 1-2 of Code-Chunk 6 check if the devtools package is installed and install it if not. After executing Line 3, GREY is installed on the user's computer.

4.3 Starting the Add-In

To call the Add-In, the user has to click on the Add-In Tab in RStudio and select "Gotta Read Em All". The Add-In itself then pops up.

4.4 Selecting the Dataset

Once the Add-In is started up, the user has to press the "Select File" button to select a file on your computer, as seen in Figure 2. Then, he/she can type in a name for the name of the dataset in R. A suggestion is made via the selected file's filename. Once the file is loaded into the Add-In, additional options for parsing the file on the right are displayed, depending on the filetype. After the user has made his/her adjustments, he/she can click on the previews tab.

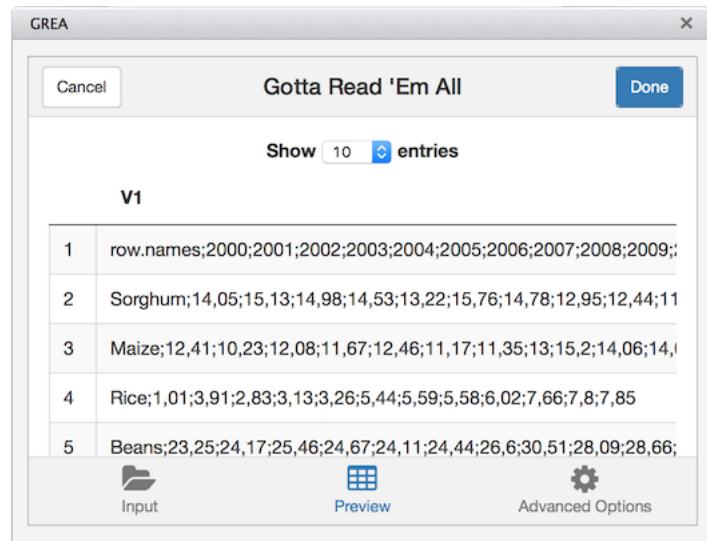


Figure 3: Looking at the preview window

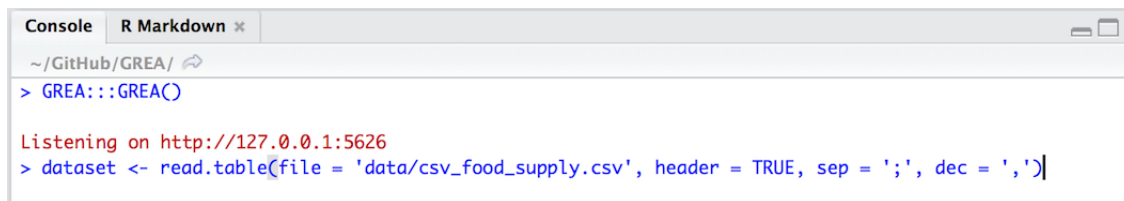


Figure 4: A ready-to use R command created by GREYA

4.5 The preview window

The previews tab (Figure 3) shows a preview of what the dataframe would look like if the user parsed it with the current settings. If something looks odd (e.g. column names fell into the first row of the dataset), the first tab has to be selected again. We can see that in our case, the column and decimal separators are wrongly specified. If everything is right, it is still recommended to head to the first tab.

4.6 Adjusting reading conditions

If the preview of the dataframe seemed off, the user now has the chance to adjust some parameters (e.g. Sheet Index for Excel files, or separator for .csv files). Additional optional parameters can be found in the “Advanced Options” tab. When a name is specified for the newly acquired dataset, the user presses

“done”. Afterwards, the function to read the dataset is pasted into the current R console (as seen in Figure 4), so the user can store it for future use.

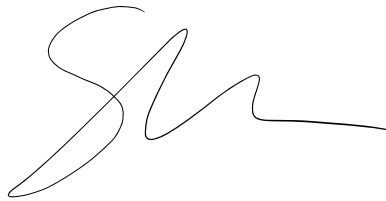
References

- [1] H. Bengtsson (2016). *R.matlab: Read and Write MAT Files and Call MATLAB from Within R*. R package version 3.6.0. <https://CRAN.R-project.org/package=R.matlab>.
- [2] W. Chang et al (2016). *shiny: Web Application Framework for R*. R package version 0.13.2. <https://CRAN.R-project.org/package=shiny>.
- [3] J. Cheng (2016). *miniUI: Shiny UI Widgets for Small Screens*. R package version 0.1.1. <https://CRAN.R-project.org/package=miniUI>.
- [4] M. Dowle et al (2015). *data.table: Extension of Data.frame*. R package version 1.9.6. <https://CRAN.R-project.org/package=data.table>.
- [5] T. Leeper et al (2016). *rio: A Swiss-army knife for data file I/O*. R package version 0.4.8. <https://CRAN.R-project.org/package=rio>.
- [6] RStudio (2016). 'Powerfully interactive'. *Shiny byRStudio*. Available: <http://shiny.rstudio.com> [Accessed 16 August 2016].
- [7] The Comprehensive R Archive Network (2016). 'Contributed Packages'. *CRAN*. Available: <https://cran.r-project.org/index.html> [Accessed 4 August 2016].
- [8] H. Wickham and J. Allaire (2016). *rstudioapi: Safely Access the RStudio API*. R package version 0.6. <https://CRAN.R-project.org/package=rstudioapi>.

Declaration of Own Work

I, Stanislaus Stadlmann, confirm that the work for the following paper with the title: “Gotta Read ’Em All: An RStudio Add-In to visually read different file-formats into R” was solely undertaken by myself and that no help was provided from other sources as those allowed. All sections of the paper that use quotes or describe an argument or concept developed by another author have been referenced, including all secondary literature used, to show that this material has been adopted to support my thesis.

Göttingen, August 22, 2016.

A handwritten signature in black ink, consisting of a large, stylized 'S' followed by a series of loops and a horizontal line extending to the right.