

RobStat™ Package Vignette

November 10, 2021

Ricardo A. Maronna, R. Douglas Martin, Victor J. Yohai, Matias Salibian-Barrera

Abstract

The RobStat™ R package is a companion to the book *Robust Statistics: Theory and Methods (with R)* by the authors of this document. The purpose of this vignette is to guide you on how to use the R scripts and related data sets in the package, which replicate examples in the book. A short Section 2 on convergence of robust estimation algorithms is also included

Contents

1 The RobStat™ Scripts and Data Sets	3
1.1 Accessing the Example R Scripts in Your Installed RobStat™	4
1.2 Loading Data Sets Provided in RobStat™	4
1.3 Installing Other Packages Needed for Data or Functions Access	5
1.4 Loading Data Sets Provided in Another Package	5
1.5 Using Functions from Another Installed Package without Loading the Package	5
1.6 Differences Between the Example R Scripts Results and the Book Examples	6
1.7 Running the Example Scripts	6
2 Convergence of Algorithms for Computing Robust Estimates	6

1 The RobStat™ Scripts and Data Sets

The first column of Table 1 below lists the numbering of the examples in the book, and the second column contains the names in the book of the corresponding R scripts. The third and fourth columns provide the names of data sets used by the scripts. Some of the R scripts in the table make use of packages other than RobStat™, for either accessing data, or for using certain R functions, and such packages are listed in the fifth column of the table. In order to run those scripts, you need to have already installed those package(s) from CRAN (<https://cran.r-project.org/>).

EXAMPLE	NAME	RobStat™ DATA	OTHER DATA	OTHER PACKAGES REQUIRED
4.1	shock.R	shock		quantreg
4.2	oats.R	oats		
5.1	mineral.R	mineral		quantreg
5.2	wood.R		wood	robustbase
5.3	step.R		made up data	
5.4	algae.R	algae		
5.5	ExactFit.R		synthetic data	
6.1	biochem.R	biochem		
6.2	wine.R	wine		
6.3	vehicle.R	vehicle		rrcov
6.4	bus.R	bus		
6.5-6.6	wine1.R	wine		GSE
6.7	autism.R		autism	WWGbook, robustvarComp, nlme
7.1	leukemia.R		leuk.dat	robust
7.2	skin.R	skin		
7.3	epilepsy.R		breslow.dat	robust
8.1	ar1.R		synthetic data	robustarima
8.2	ar3.R		synthetic data	robustarima
8.3	identAR2.R		synthetic data	robustarima
8.4	identMA1.R		synthetic data	robustarima
8.5	MA1-A0.R		synthetic data	robustarima
8.6	resex.R	resex		robustarima

Table 1: R Scripts and Data in the RobStat™ R Package

We note that the script `flour.R`, and the data set `flour` used by the script, contained in RobStat™ but not listed in the above table are for Example 1.1, Figure 2.1, and Table 2.4 of the book. Furthermore, the data set `neuralgia` contained in RobStat™ but not listed in the table above, is used only in Problem 7.1 of the book.

1.1 Accessing the Example R Scripts in Your Installed RobStatTM

In case you have not already installed RobStatTM, install it with the command:

```
install.packages("RobStatTM")
```

With RobStatTM installed, load it in your current R session with:

```
library("RobStatTM")
```

In order to access the example R scripts, you need to find your installed RobStatTM “scripts” folder, which is one of several RobStatTM package folders. You can find the location of the scripts folder on your computer by using the function `system.file()` as follows:

```
system.file("scripts", package = "RobStatTM")
```

NOTE: Copy/paste of the above line does not typically work, so you should type it in. The result of using this command will depend upon your computer. For example, in the case of a particular computer running Windows 10, the result is:

```
[1] "C:/Users/Doug/Documents/R/win-library/3.4/RobStatTM/scripts"
```

Then you just need to navigate to the `scripts` folder, where you will see all of the example R scripts, with the same name they are given in the book, i.e., the name in the `NAME` column of Table 1. You should then copy/paste any script, or all the scripts, to some other location on your computer where you want to run them.

1.2 Loading Data Sets Provided in RobStatTM

To load any data set that is in RobStatTM, just use the `data()` function with the data set name. For example, in the case of the `shock` data set:

```
data(shock)
head(shock, 2)

##   n.shocks time
## 1         0  11.4
## 2         1  11.9
```

1.3 Installing Other Packages Needed for Data or Functions Access

You will notice in the table on the previous page that there is a column named OTHER PACKAGES REQUIRED. The other package is required for accessing a data set in the package, or for using an R function in the package, or both. In the case of any example R script that requires another package for one of those reasons, you need to have installed the package before running the script. You could do this on a case by case basis, but it may be easier to just be sure and install all the packages in the OTHER PACKAGES REQUIRED once and for all. You can do this quite easily from RStudio.

1.4 Loading Data Sets Provided in Another Package

To load a data set that is in another package that you have already installed in your R (but have not loaded, and in fact should not load) you use the `data()` function with an optional argument that specifies the package the data is in. For example, in the case of the `wood` data set that is in the `robustbase` package:

```
data(wood, package = "robustbase")
head(wood, 1)

##      x1      x2      x3      x4      x5      y
## 1 0.573 0.1059 0.465 0.538 0.841 0.534
```

Note that the above use of `data()` does not result in loading the named package, and this is the recommended way to load data in another package. The reason for not first loading another package with `library()`, and then using the `data()` function, is that the result of loading a package can result in problematic masking of functions in RobStatTM.

1.5 Using Functions from Another Installed Package without Loading the Package

When a script requires a function from a package other than RobStatTM, the other package needs to have been already installed. However, in order to avoid masking problems mentioned above, the script does not load the library, and instead uses that standard package referencing mechanism “::”. For example, you can use the function `rq` in the `quantreg` package to do an L1 (least absolute deviation regression) fit of `zinc` to `copper` in the `mineral` data set with the code below:

```
minerall1 <- quantreg::rq(zinc ~ copper,
  data = mineral)
```

1.6 Differences Between the Example R Scripts Results and the Book Examples

There are a number of example R scripts where part of the output, either pure numerical values or numerical values in plots, is not exactly the same as in the book examples. Scripts where this is the case include `oats.R` (different p-values, Figure 4.4), `mineral.R` (Figures 5.3, 5.5, 5.6), `wood.R` (Figures 5.10, 5.11), `algae.R` (Figure 5.14, 5.15), `wine.R` (Figure 6.11). This is because the final code for some of the scripts is different from the code used for the book examples, and the scripts code is usually improved in some way. Typically, the conclusions drawn from the figures, concerning which data points are outliers, and same for the book figures and the figures produced by the example scripts.

1.7 Running the Example Scripts

We recommend running all the example scripts, and especially the Table 1 Examples 5.1, 5.2, 5.3, 5.4, 5.5, which reflect our recent recommendation to use `family = "mopt"`, with `efficiency = 0.95` for the robust regression function `lmrobdetMM`.

2 Convergence of Algorithms for Computing Robust Estimates

RobStatTM contains sophisticated optimization algorithms for computing robust estimates, and convergence of the algorithms is determined by certain parameters. For example for the `lmrobM` and `lmrobdetMM` functions, the convergence parameters are set with the functions `lmrobM.control` and `lmrobdet.control`, respectively, which contain a number of parameters that have default settings. For example, `max.it` with default value `max.it = 100`, determines the maximum number of iterations of an iterated weighted least square (IRWLS) algorithm, and `rel.tol` with a default value `rel.tol = 1e-07`, determines a relative change threshold. See the book Section 4.5.2 for the algorithm used by `lmrobM`, including the relative change being measured, where step 4.3 uses for ϵ the value of `rel.tol`.

When a RobStatTM function gives a warning message, for example

M-step did NOT converge. Returning unconverged lm-estimate

it does not mean that the estimate is necessarily bad, it is just that the algorithm has not converged given the default values of `rel.tol` and `max.it`, and the estimate is usually (but not always) quite viable. Our recommendation is that when such a warning occurs, one should adjust the parameter `rel.tol` downward, and/or adjust the parameter `max.it` upward.