

Efficient Management of Data in R (Data Structures!)

Data Science Lecture Series: Advanced R

W. Evan Johnson, Ph.D.
Professor, Division of Infectious Disease
Director, Center for Data Science
Rutgers University – New Jersey Medical School

2023-05-07

Section 1

Importing data

Importing data

The first problem a data scientist will usually face is how to import data into R!

Often they have to import data from either a file, a database, or other sources. One of the most common ways of storing and sharing data for analysis is through electronic spreadsheets.

A spreadsheet stores data in rows and columns. It is basically a file version of a data frame (or a tibble!).

Importing data

A common function for importing data is the `read.table` function:

```
mydata <- read.table("mydata.txt")
```

This is looking for a structured dataset, with the same number of entries in each row, and data that is delimited with a single space between values.

Importing data

The `read.table` function can also read tab-delimited data:

```
mydata <- read.table("mydata.txt", sep="\t")
```

Or comma separated (.csv) formats:

```
mydata <- read.table("mydata.txt", sep=",")
```

(also explore the `read.csv` function)

Importing data

We can also add options to set the first column as a header and select a row for the row labels:

```
mydata <- read.table("mydata.txt",  
                      header=TRUE,  
                      row.names="id")
```

Importing data

Excel files can also be directly imported using `read.xlsx`:

```
library(xlsx)
mydata <- read.xlsx("myexcel.xlsx")
```

And one can also select a specific sheet in the Excel file:

```
mydata <- read.xlsx("myexcel.xlsx",
                    sheetName = "mysheet")
```

Other functions for importing data

Other useful importing tools are `scan`, `readLines`, `readr`, and `readxl`. The latter two we will discuss later.

Section 2

Exporting Data

Exporting data

We have many options for exporting data from R. For data frames, one of the easiest ways to output data is with the `write.table` function:

```
write.table(dat, file = "data_out.txt",  
            quote = FALSE, sep = ",",  
            row.names = TRUE,  
            col.names = TRUE)
```

Exporting data

Another important and useful way of inputting/outputting data is in an Rds object:

```
saveRDS(dat, file = "dat.Rds")  
dat.copy <- readRDS(file = "dat.Rds")
```

Section 3

Introduction to Data Structures

Importance of data structures

A data structure is a particular way of organizing data in a computer so that it can be used effectively. The idea is to reduce the space and time complexities of different tasks.

Data structures in R programming are tools for holding multiple values, variables, and sometimes functions.

Please think very carefully about the way you manage and store your data! This can make your life much easier and make your code and data cleaner and more portable!

Types of data structures in R

R's base data structures are often organized by their dimensionality (1D, 2D, nD) and whether they're homogeneous or heterogeneous (elements of identical or various type). Six of the most common data types are:

- 1 Vectors
- 2 Lists
- 3 Matrices
- 4 Arrays
- 5 Factors
- 6 Data frames (or tibbles)

Section 4

Data Frames

Data Frames

The most common data structure for storing a dataset in R is in a **data frame**. Conceptually, we can think of a data frame as a two dimensional table with rows representing observations and the different variables reported for each observation defining the columns. Data frames are particularly useful for datasets because we can combine different data types into one object.

Data Frames

We can convert matrices into data frames using the function `as.data.frame`:

```
mat <- matrix(1:12, 4, 3)
mat <- as.data.frame(mat)
```

Or just generate it directly using the `data.frame` function:

```
dat <- data.frame(x=1:4, y=5:8, z=9:12)
```

A `data.frame` can be indexed as matrices, `dat[1:2, 2:3]`, and columns can be extracted using the `$` operator.

Section 5

Tibbles

Tibbles

Here is a printed version of the data frame:

```
dat
```

```
##      x y  z
##  1  1 5   9
##  2  2 6  10
##  3  3 7  11
##  4  4 8  12
```

Tibbles

A **tibble** is a modern version of a `data.frame`.

```
library(tidyverse)
dat1 <- tibble(x=1:4, y=5:8, z=9:12)
```

Or convert a `data.frame` to a tibble

```
dat <- data.frame(x=1:4, y=5:8, z=9:12)
dat1 <- as_tibble(dat)
```

Tibbles

Here is a printed version of the tibble:

```
dat1
```

```
## # A tibble: 4 x 3
##       x     y     z
##   <int> <int> <int>
## 1     1     5     9
## 2     2     6    10
## 3     3     7    11
## 4     4     8    12
```

Tibbles

Important characteristics that make tibbles unique:

- 1 Tibbles are primary data structure for the tidyverse
- 2 Tibbles display better and printing is more readable
- 3 Tibbles can be grouped
- 4 Subsets of tibbles are tibbles
- 5 Tibbles can have complex entries—numbers, strings, logicals, lists, functions.
- 6 Tibbles can (almost) enable object-orientated programming in R

Section 6

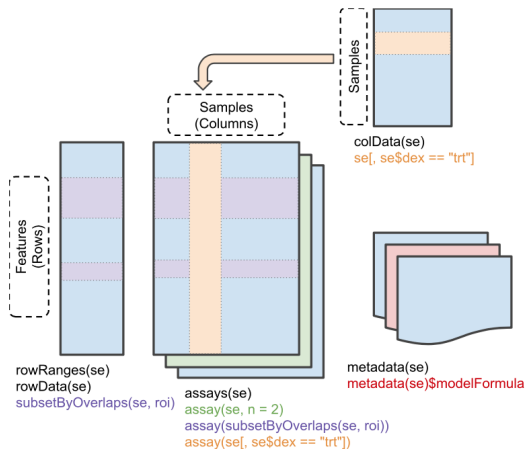
Advanced Data Structures in R

Advanced Data Structures in R

In your Extra Practice, you will explore more advanced R data structures, namely the **S3** and **S4** class objects. These can facilitate object orientated programming.

Advanced Data Structures in R

One example of an S4 class data structure is the **SummarizedExperiment** object.



Summarized Experiment

Session info

```
sessionInfo()
```

```
## R version 4.2.3 (2023-03-15)
## Platform: aarch64-apple-darwin20 (64-bit)
## Running under: macOS Ventura 13.3.1
##
## Matrix products: default
## BLAS:   /Library/Frameworks/R.framework/Versions/4.2-arm64/Resources/lib/libRblas.0.dylib
## LAPACK: /Library/Frameworks/R.framework/Versions/4.2-arm64/Resources/lib/libRlapack.dylib
##
## locale:
## [1] en_US.UTF-8/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8
##
## attached base packages:
## [1] stats      graphics  grDevices  utils      datasets  methods   base
##
## other attached packages:
## [1] lubridate_1.9.2 forcats_1.0.0  stringr_1.5.0  dplyr_1.1.1
## [5] purrr_1.0.1     readr_2.1.4    tidyr_1.3.0    tibble_3.2.1
## [9] ggplot2_3.4.2   tidyverse_2.0.0
##
## loaded via a namespace (and not attached):
## [1] pillar_1.9.0    compiler_4.2.3  tools_4.2.3    digest_0.6.31
## [5] timechange_0.2.0 evaluate_0.20    lifecycle_1.0.3 gtable_0.3.3
## [9] pkgconfig_2.0.3 rlang_1.1.0     cli_3.6.1      rstudioapi_0.14
## [13] yaml_2.3.7      xfun_0.38       fastmap_1.1.1  withr_2.5.0
## [17] knitr_1.42      generics_0.1.3  vctrs_0.6.1    hms_1.1.3
## [21] grid_4.2.3      tidyrselect_1.2.0 glue_1.6.2      R6_2.5.1
## [25] fansi_1.0.4     rmarkdown_2.21  tzdb_0.3.0     magrittr_2.0.3
## [29] scales_1.2.1    htmltools_0.5.5 colorspace_2.1-0 utf8_1.2.3
## [33] stringi_1.7.12  munsell_0.5.0
```