

R DATA INPUT AND OUTPUT

LECTURE 11

Dirk Eddebuettel

STAT 430: Data Science Programming Methods (Fall 2019)

Department of Statistics, University of Illinois

DATA INPUT AND OUTPUT

DATA SCIENTIST



What my friends think I do



What my mom thinks I do



What society thinks I do



What my company thinks I do



What I think I do



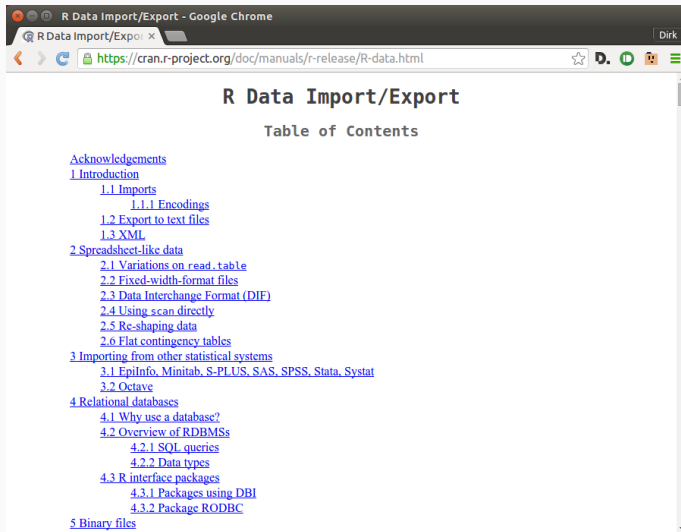
What I actually do

data-science-dojo

Source: <https://twitter.com/datasciencedojo/status/627180472104128512>

Getting Data Into R is part of just about any analysis!

- R excels at this
 - truly broad coverage of file formats
 - as well as ‘backends’ such as databases
 - or different web-based APIs
- Our focus: read/write of csv data
- Mention other formats: json, xml, ...
- Efficient R-specific storage: rds and fst
- Mention protobuf, msgpack, feather, parquet, ...



The screenshot shows a Google Chrome browser window with the title "R Data Import/Export - Google Chrome". The address bar displays the URL "https://cran.r-project.org/doc/manuals/r-release/R-data.html". The main content area is titled "R Data Import/Export" and "Table of Contents". The table of contents lists various sections and sub-sections, all of which are underlined as links. The sections include Acknowledgements, Introduction, Imports (with a sub-section on Encodings), Export to text files, XML, Spreadsheet-like data (with sub-sections on read.table, Fixed-width-format files, Data Interchange Format (DIF), Using scan directly, Re-shaping data, and Flat contingency tables), Importing from other statistical systems (with sub-sections on EpiInfo, Minitab, S-PLUS, SAS, SPSS, Stata, and Systat), Octave, Relational databases (with sub-sections on Why use a database?, Overview of RDBMSs, SQL queries, Data types, R interface packages, Packages using DBI, and Package RODBC), and Binary files.

Acknowledgements
1 Introduction
1.1 Imports
1.1.1 Encodings
1.2 Export to text files
1.3 XML
2 Spreadsheet-like data
2.1 Variations on read.table
2.2 Fixed-width-format files
2.3 Data Interchange Format (DIF)
2.4 Using scan directly
2.5 Re-shaping data
2.6 Flat contingency tables
3 Importing from other statistical systems
3.1 EpiInfo, Minitab, S-PLUS, SAS, SPSS, Stata, Systat
3.2 Octave
4 Relational databases
4.1 Why use a database?
4.2 Overview of RDBMSs
4.2.1 SQL queries
4.2.2 Data types
4.3 R interface packages
4.3.1 Packages using DBI
4.3.2 Package RODBC
5 Binary files

Chapters in R Data Import/Export Manual

- Introduction
- Spreadsheet-like data
- Importing from other statistical systems
- Relational databases
- Binary files
- Image files
- Connections
- Network interfaces
- Reading Excel spreadsheets
- References

Comes with every R installation, but a little dry to read...

Textfiles

- Text files for interchange are common
- File separators have historically been
 - a comma, hence “comma-separated values” or CSV
 - a tab (or other whitespace) in files ending in `.txt`
 - other less common values such as `|` etc
- Possible issues to encounter:
 - other world regions use a comma where we use a dot
 - values must then be quoted and/or alternate separator used
 - “encodings” can also be an issues

Basic use

- The underlying base R function is `read.table()`
 - it has numerous options so see the help page
 - defaults to *any* so-called white space as separator (*i.e.* one or more tabs, spaces, newlines)
 - also defaults to *no header*
 - can handled missing values and different types
- More frequently used: `read.csv()`
 - wrapper around `read.table()`
 - different defaults for separator, header, fill
 - we will discuss both, one can focus on this one

read.table: work horse behind other functions

```
R> args(read.table)
function (file, header = FALSE, sep = "", quote = "\"'", dec = ".",
  numerals = c("allow.loss", "warn.loss", "no.loss"), row.names,
  col.names, as.is = !stringsAsFactors, na.strings = "NA",
  colClasses = NA, nrows = -1, skip = 0, check.names = TRUE,
  fill = !blank.lines.skip, strip.white = FALSE, blank.lines.skip = TRUE,
  comment.char = "#", allowEscapes = FALSE, flush = FALSE,
  stringsAsFactors = default.stringsAsFactors(), fileEncoding = "",
  encoding = "unknown", text, skipNul = FALSE)

NULL
R>
```

See `help(read.table)` for more complete coverage

read.table: Examples

```
# no header, default separator
```

```
read.table("file.txt")
```

```
# header from 1st line, 'pipe' separator
```

```
read.table("file.txt", header=TRUE, sep="|")
```

```
# do not convert text variables to factors
```

```
read.table("file.txt", stringsAsFactors=FALSE)
```

```
# skip 1st column, 2nd column taken as numeric
```

```
read.table("file.txt", colClasses=c("NULL", "numeric"))
```

`read.csv`: comma-separated file reader

```
R> args(read.csv)
```

```
function (file, header = TRUE, sep = ",", quote = "\"",  
          dec = ".", fill = TRUE, comment.char = "", ...)
```

```
NULL
```

```
R>
```

`header`, `row.names`, `col.names`, `stringsAsFactors`, `sep`, ... are all very important, and some are set differently now.

... means other arguments permitted and passed through.

- There are a number of titanic datasets in R packages
- Base R package **datasets** has it too
- Here I wanted a file, so [this page at Stanford](#) had one

```
R> data <- read.csv("titanic.csv") # default Arguments
R> str(data)
# 'data.frame': 887 obs. of  8 variables:
# $ Survived      : int  0 1 1 1 0 0 0 0 1 1 ...
# $ Pclass        : int  3 1 3 1 3 3 1 3 3 2 ...
# $ Name          : Factor w/ 887 levels "Capt. Edward Gifford Crosby",...: 602 823 172 814 733 464 700 3
# $ Sex           : Factor w/ 2 levels "female","male": 2 1 1 1 2 2 2 1 1 ...
# $ Age           : num  22 38 26 35 35 27 54 2 27 14 ...
# $ Siblings.Spouses.Aboard: int  1 1 0 1 0 0 0 3 0 1 ...
# $ Parents.Children.Aboard: int  0 0 0 0 0 0 0 1 2 0 ...
# $ Fare          : num  7.25 71.28 7.92 53.1 8.05 ...
```

- One commonly use option: `stringAsFactors=FALSE`

```
R> data <- read.csv("titanic.csv", stringsAsFactors=FALSE)
R> str(data)
# 'data.frame': 887 obs. of  8 variables:
# $ Survived      : int  0 1 1 1 0 0 0 0 1 1 ...
# $ Pclass        : int  3 1 3 1 3 3 1 3 3 2 ...
# $ Name          : chr  "Mr. Owen Harris Braund" "Mrs. John Bradley (Florence Briggs Thayer) Cumings" "
# $ Sex           : chr  "male" "female" "female" "female" ...
# $ Age           : num  22 38 26 35 35 27 54 2 27 14 ...
# $ Siblings.Spouses.Aboard: int  1 1 0 1 0 0 0 3 0 1 ...
# $ Parents.Children.Aboard: int  0 0 0 0 0 0 0 1 2 0 ...
# $ Fare          : num  7.25 71.28 7.92 53.1 8.05 ...
```

Notice the difference?

Save an individual R object *portably*

Very useful for *efficient* (compressed, binary) and *portable* (from machine to machine) data storage

```
# save obj portably and compressed  
saveRDS(obj, file="somefile.rds")
```

```
# in another session or on another machine  
newobj <- readRDS("somefile.rds")
```

```
R> saveRDS(data, file="titanic.rds") # save to an rds
R> newdata <- readRDS("titanic.rds") # reload
R>
R> identical(newdata, data) # check that bitwise identical
[1] TRUE
R>
R> # compressed and portable storage
R> object.size(gs)
118176 bytes
R> file.info("titanic.rds")$size
[1] 18800
```

fread (in package `data.table`)

- extremely fast as it reads chunks in parallel
- good heuristics for inferring column types
- generally “just works” and is fast
- also reads from commands
- recommended (and `data.table` is next lecture)
- (and there is also `fwrite`)
- (and we will see both later)

read_csv (in package **readr**)

- part of tidyverse, opinionated
- not as fast as **fread**
- not as general (may need columns specified)
- returns as a **tibble** object which or may not be desirable

Add-on Package

- Fastest binary writer / reader
- *Highly* optimized, parallel
- Very fast reads and writers
- Good for large-enough data
- But maybe not worth for small data
- On CRAN as package `fst`
- Website fstpackage.org



Very useful extensions

- In `read.*()` functions, `file=...` generalizes
- We can read from any “Connections” object
 - a URL ie anything “on the web” (or network-reachable)
 - a pipe ie output from another program
 - a compressed file
 - a network socket
- See `help(connections)` and the Data I/O Manual for more.

Example

```
R> # we can read directly via http or https
R> webdata <- read.csv("https://some.site.com/path/data.csv")
R>
R> # we can also read output from other programs directly
R> res <- read.table(pipe("python someOtherScript.py"))
R>
R> # or any other shell program via pipe() ...
```

DATABASES

DBI Abstraction

- DBI: One of the oldest packages for R (and predecessor S(-Plus)).
- Powerful abstraction of db *backend* via a unified interface.
- Connect to a given db using the corresponding driver.
- Once connected, run query, insert, updates, ... as usual.
- To run code on a different engine, just switch one statement.

DBI Backends

- RMySQL
- ROracle
- RPostgreSQL
- RSQLite
- RSQLServer
- MonetDB.R

and more

DBI Example: SQLite

```
R> library("RSQLite")
```

```
Loading required package: DBI
```

```
R> con <- dbConnect(RSQLite::SQLite(), ":memory:")
```

```
R> dbWriteTable(con, "titanic", data)
```

```
[1] TRUE
```

```
R> dim(dbGetQuery(con, "select * from sales limit 10"))
```

```
[1] 10 8
```

```
R>
```


DBI Example: PostgreSQL

```
R> library("RPostgreSQL")  
R> drv <- dbDriver("PostgreSQL")  
R> con <- dbConnect(drv, "username", "password", "dbname")  
R> data <- dbGetQuery(con, "SELECT * from sales")  
R> dbDriver(con)
```

Others

- Lots of other connection packages
- RODBC, RJDBC,
- NoSQL: redis, mongo, monet, ...
- Google bigquery, Apache Cassandra
- Hadoop, Spark, ...
- If something has data, there will be an R package somewhere.

SPREADSHEETS

Reading xls/xlsx

Several CRAN packages offer functionality:

- **readxl** is a more recent and fast reader, not writing
- **openxlsx** is a C++-based bi-directional package
- **XLConnect** allows bi-directional read + write
- (older, now archived) **xlsx** requires Java, reliable, bi-directional, allows creation of full-featured spreadsheets

Writing xls/xlsx

Several CRAN packages offer functionality:

- `writexl` can export to xlsx, fast, no dependency
- `openxlsx` is a C++-based bi-directional package
- `XLConnect`, also bi-directional
- `xlsx` (archived) requires Java, reliable, bi-directional, allows creation of full-featured spreadsheets

OTHER

Choices:

- Several XML readers: `XML`, `xml2`, ...
- Several JSON readers: `RJSONIO`, `jsonlite`, ...
- Other readers for YAML, TOML, ...
- Scientific data formats such as HDF5 and NetCDF ...
- ‘Big data’ formats like Parquet via `arrow` or `miniparquet`
- Plus more domain-specific readers

In general, if a format exists, there may be a reader for it.

The excellent CRAN Task View on [Web Technologies](#) has sections on

- Parsing Data from the Web
- Curl, HTTP, FTP, HTML, XML, SOAP
- Authentication
- Web Frameworks
- JavaScript
- Code Sharing

plus a truly exhaustive list of *Data Sources on the Web Accessible via R* organised by discipline

SUMMARY

R can read or write

- just about anything from text files
- just about any existing database backend
- specialised files (`xls`, `xlsx`, statistics packages)
- via connections from programs, URLs, and more
- higher-performance readers and writers