

R Best Practice

Chad Goymer

2018-08-15

Contents

1	Introduction	5
1.1	End-User Computing	5
1.2	Risk Analysis	6
1.3	Documentation	6
1.4	Testing	6
1.5	Reproducibility	7
1.6	Change control	7
1.7	Access control	8
1.8	Appendix: Solvency II: internal model approval process data review findings	8
2	Software	9
2.1	Core R Applications	9
2.2	Supporting Software	9
2.3	Package Versions	9
3	Writing R Code	11
3.1	Writing Style	11
3.2	Structure	12
3.3	R Markdown	14
4	Recommended Packages	23
4.1	Data Wrangling	23

Chapter 1

Introduction

This document sets out the proposals for best practice when using R. The aim is to provide a consistent and comprehensive approach to using R for analysis and applications. It forms part of the wider guidance on end-user computing. The document is split into the following areas:

1. **Software:** Details the recommended software to install, both the core R applications and supporting software.
2. **Writing:** Gives guidance on how to write R scripts and applications.
3. **Packages:** Lists the recommended packages to use for achieving common tasks.
4. **Development:** Describes the recommended approaches to developing in R.

This document is not supposed to be a comprehensive guide to using R. It gives guidance and sign-posts users to resources with more details.

1.1 End-User Computing

The Bank of England's 2016 report *Solvency II: internal model approval process data review findings* identifies end-user computing as a Solvency II risk. In particular, it states:

Spreadsheets and other user-developed applications are a form of information technology, and all information technology needs to be appropriately controlled.

End-user computing requires a user to think about the following areas when developing in R:

- Risk analysis
- Documentation
- Testing
- Reproducibility
- Change control
- Access control

The result of the risk analysis determines the level at which the subsequent points need to be addressed.

1.2 Risk Analysis

The level of documentation, testing and change control depends on two factors:

1. Complexity
2. Criticality

Complexity: If a piece of work is very simple it may only require a short description, a manual test and saving to a secure location. However, if it is highly complex it is important for others, and your future self, to document it thoroughly, provide repeatable tests and save the scripts in version control.

Criticality: If a piece of work is a simple analysis for curiosity it may not need any controls applied to it. However, if it provides functionality which other systems rely on or informs major business decisions then it is important to ensure it is thoroughly tested and reviewed and changes are not made without careful consideration of how it affects dependents.

In order to assess the level of controls we perform a simple risk analysis by estimating the level of complexity and criticality of the work, using a scale of **high**, **medium** or **low**. If both are low then no further consideration is required. However, If any of the above factors are medium or high then the work is considered material and should be subject to the following controls.

1.3 Documentation

As a minimum the following should be documented with the code. For simple scripts, a commented header in the file is sufficient. For multi-file applications it is recommended that a README file is added.

- **Name:** The name of the project, analysis or application.
- **Description:** A brief description of what the code does (not how).
- **Purpose:** The reason for writing the code.
- **Instructions:** How to use the code - for the benefit of other users.

For **high complexity** work the following should be considered

- Write documentation for each section/function of code separately .
- Provide examples of usage.

For **high criticality** work consider creating an R package then:

- Create a DESCRIPTION file to document the package.
- Use the package **Roxygen2** to document functions.
- Write a vignette explaining its use.

1.4 Testing

As a minimum, manual tests should be performed and the results recorded alongside the code. This may take the following forms:

- Assertions in the code to ensure reasonableness, e.g. result is within a range.
- Comparisons with known results.

For **high complexity** work consider the following:

- Separate tests for each section/function of code.
- Comparisons with simple inputs and expected results

For **high criticality** work consider creating an R package then:

- Use the package `testthat` to write executable tests.
- Execute regression tests whenever changes are made.
- Perform integration tests with dependent systems.

1.5 Reproducibility

Important work should be reproducible by others, and your future self. This means keeping track of inputs, parameters and code used to produce results. The simplest approach is to keep a copy of input data and parameters alongside the code used to produce the saved results. Then it is self-contained.

Also consider using R markdown which captures description, code and results in a single document.

For **high complexity** work consider creating an RStudio project it has the following benefits:

- Sets the current working directory so relative file paths can be used.
- Saves open files so you can start where you left off.

Some other recommendations worth considering are:

- Avoid referencing external files if possible.
- Make dependencies clear.
- Cache versions of input data into a subfolder, if not large.

For **high criticality** work also consider

- Outputting audit logs with time stamp, username and computing environment
- Holding inputs in version control

1.6 Change control

Change control is about managing updates and bug fixes to the code in a way that is tracked, reviewed, and approved. It ensures change, and its implications, is understood and can also be reversed, if necessary. The simplest approach is to produce new files for each version, keeping older versions elsewhere. Each new version should document the reason for the change.

A much better solution is to use version control software to track changes. This makes the process much more robust. It is recommended that Git is used for version control and shared using a hosting application such as GitHub (or TFS if not available).

For **high complexity** and/or **criticality** consider the following:

- Track changes in the hosting application (GitHub or TFS).
- Record description and reason for change and prioritise.
- Create a new branch for each change.
- Merge back into the main branch only after a review.
- Require approval from an appropriate person or group before releasing to production.

1.7 Access control

There should be appropriate controls on who has access to the source code and/or results. The developers should be aware of who has access the work they are producing. The simplest approach is to save files into a location with adequate access controlled by the IT dept. It may also be appropriate to limit who has the ability to change the source code separately.

For **high complexity** work it is important only those with the appropriate level of knowledge can change the code.

For **high criticality** work production code should not be changable directly. A release process should be set up so appropriate approval must be given before code is promoted.

In both cases:

- The list of people with access to the source code and results should be reviewed on a regular basis.
- Appropriate back up and recovery strategies should be implemented to ensure rollback and re-installation is possible.

1.8 Appendix: Solvency II: internal model approval process data review findings

1.8.1 Sub-risk 5: IT environment, technology and tools

Spreadsheets and other user-developed applications are a form of information technology, and all information technology needs to be appropriately controlled.

Finding 9: end-user computing (EUC)

4.36 Spreadsheets and other end-user applications (2012.4.39) remained common in capital and balance sheet modelling. The PRA does not have a view on whether end-user computing (EUC) is appropriate, as it is a form of IT, and all IT needs to be appropriately controlled. Where EUC is material to the internal model data flow, the PRA will be looking for appropriate controls for data quality such as reasonableness checks, input validations, peer reviews, systems environment configuration, logical access management, ongoing change controls (development, build, systems and user acceptance testing) and release management (including implementation and operational testing), disaster recovery, and documentation.

4.37 Automation of spreadsheets reduces the risk of manual error (2012.4.42), but can introduce different problems such as reduced oversight, inadequate transparency about the extent of linking and the proliferation of nested spreadsheets and the attendant issue of ‘broken links’.

4.38 The 2012 report did not engage comprehensively with cyber risk. This is likely to be an area of increasing focus, following alerts and increasing concerns about security as firms move away from localised application and onto networked platforms. As noted in the Bank of England Financial Stability Report,¹ cyber attacks can threaten financial stability by disrupting the provision of critical functions from the financial system to the real economy. The Financial Policy Committee has recommended that resilience testing be a regular part of core firms’ cyber resilience assessment. Insurers providing cover for cyber or business interruption are also indirectly exposed to cyber risk.

Finding 10: IT infrastructure

4.39 Complex IT implementations (2012.4.44) can be challenging to manage without a clear definition of user requirements, design, testing and appropriate controls for effective operation in business as usual. This continued to be an area of risk. One firm took seven years to implement a tactical system, and still has no strategic system for its upstream administration processes.

Chapter 2

Software

When developing R code which may be shared or executed by another person it is important to ensure everyone is using the same versions of software and packages.

2.1 Core R Applications

The minimum software required to write R code is the R application itself. However, it is better to use a separate editor for creating R scripts. The recommended software installation is:

- Microsoft R Open, which is 100% compatible with the standard version of R and provides better performance in some cases. Like the standard version of R, Microsoft's version is free.
- RStudio Desktop, to edit R script and build projects. RStudio can also be used for free.
- RTools, which is required to build R packages. RTools is free to download.

RStudio provides a cheatsheet detailing commonly used functionality of its editor:

2.2 Supporting Software

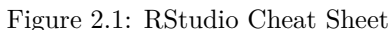
In addition to the above, when writing R code of high complexity or criticality it is recommended that version control is used:

- Git is the market leader and integrates well with RStudio. Git is free
- GitKraken provides a very good graphical interface to Git and access to more features than RStudio. GitKraken is not free, but it does not cost very much.

2.3 Package Versions

A common problem with R is ensuring everyone is using the same version of packages. There are a number of ways to deal with this problem:

1. Another approach is to create a site library, a location on a network drive where everyone has access to. The default library location can then be set to this location so everyone uses the same packages. The disadvantage of this approach is that if someone updates a package it is updated for everyone.



- 1 and 2 can be combined to avoid installing packages for every user.

- The disadvantage of this approach is the required packages are downloaded separately for each project, so should only be used for exceptional cases.

Chapter 3

Writing R Code

3.1 Writing Style

When writing code it is important to follow a common style so it is readable and other people, and your future self, can easily understand and extend it.

3.1.1 Tidyverse

It is recommended, when using R, that we use the “Tidyverse” approach and packages wherever possible. The Tidyverse is a collection of packages designed for data science, as well as a philosophy and style for formatting data and writing functions. The tidyverse.org website details the packages involved and contains articles on how to use them. Some of it will be summarised below, but the website contains the definitive information.

The most in-depth treatment on the Tidyverse can be found in the book “R for Data Science”, which is available online for free at r4ds.had.co.nz.

The tidyverse has an extensive style guide, which we summarise here:

3.1.2 Files

- File names should be meaningful and end in `.R`, or `.Rmd` for R markdown files. Only use letters, numbers and `-` or `_`.
- If your script required packages load them all at once at the very beginning of the file.
- Use comments to explain the “why” not the “what” or “how”.
- Break up files into named sections using commented lines (`# ----`, example below). In RStudio you can collapse and expand sections commented this way.

```
# Load data -----
```

3.1.3 Syntax

- Variable and function names should use only lowercase letters, numbers, and `_`.
- Always indent the code inside curly braces `{}` by two spaces.

3.1.4 Functions

- Use verbs for function names where possible.
- If a function has numerous arguments put each one on a new line.
- A function should do one thing well. If it is doing too much the break it up.
- A function should be easily understandable in isolation. It should not refer to any variables outside the function scope.

3.1.5 Pipes

- Use `%>%` when you find yourself composing three or more functions together, instead of a nested call.
- `%>%` should always have a space before it and a new line after it. After the first step, each line should be indented by two spaces.

3.1.6 Documentation

- Created functions should be documented so others, including the future you, can understand what the function does and how to use it.
- Documentation should be written before the function definition in an roxygen2 style. roxygen uses special comments, starting with `#'`. The first line is the title, and anything else, not prefixed with a keyword forms the description. Keywords start with `@` and the most important ones are `@param` to describe a function parameter and `@return` to describe what the function returns.

Here is a very simple example:

```
#' The length of a string.
#'
#' This function returns the number of characters in the supplied string.
#'
#' @param string input character vector
#'
#' @return integer vector giving number of characters in each element of the
#'   character vector.
#'
#' @export
#'
str_length <- function(string) {
  nchar(string)
}
```

3.2 Structure

Organising files for a particular piece of work becomes more important as the scale and complexity increases. Standard approaches exist to simplify the workflow.

3.2.1 Projects

Use RStudio projects to organise files. This has a number of advantages:

1. Sets the working directory to the project location
2. Reopens the same files when returning to the projects

3.2.2 Folders

When an analysis becomes complex it should be split up into logical parts and stored in subfolders. Store the original data in a folder, unchanged. It is better to “cleanse” input data with an R script as it can be repeated when data changes, and/or the approach changed itself.

R code may also be stored in a separate folder. You may have an R script for cleansing the data and another for performing an analysis. Include an R script at the top level which executes the code in the subfolder in the appropriate order. Use relative paths to the files. If an RStudio project has been created the working directory will be set to the project directory automatically.

Output the results, plots, data, etc, in another folder so it is clear whether data files are results rather than inputs.

An example of a project structure:

- data
 - interesting_data.xlsx
 - reference_data.csv
- R
 - clean_data.R
 - analyse.R
- tests
 - test_cleansed_data.R
 - test_analysis.R
- results
 - cool_plot.png
 - table_of_results.csv
- run_code.R
- README.md

3.2.3 README

Adding a README file is a good way to explain to other, and you future self, what the analysis does and how to use it. The documentation section of the best practice has more detail on what should be included.

It is recommended that markdown is used to write the README. It is a very simple way to specify text formatting in a plain text file and can be converted to many other formats (HTML, docx, PDF) if required. In addition, if the package is stored in GitHub a markdown README is automatically rendered on the repository’s page.

3.2.4 R Packages

When R code has high criticality consider turning it into a package. A package is a way of collecting together related code in a robust way. It has the following advantages:

- Easier to share with others (as a zip file)
- Documentation is compiled into help pages
- All tests can be executed with a single command
- Can implement a development and release process

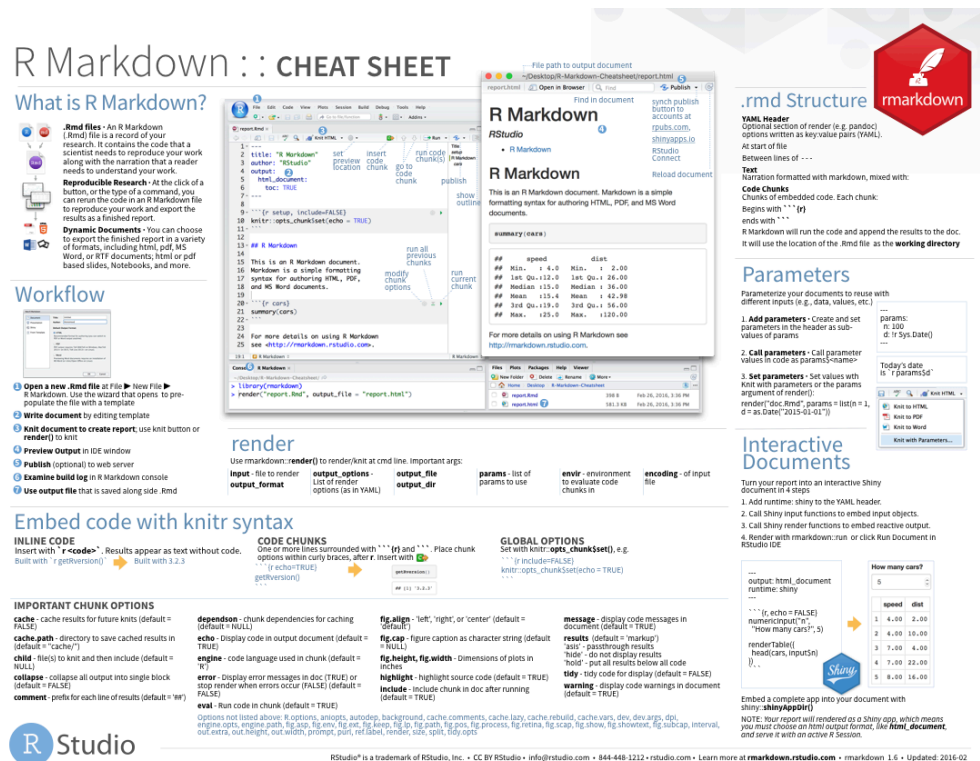


Figure 3.1: R Markdown Cheat Sheet

- Code is broken up into useful functions

Writing a package is very straight forward with the helper packages available today. More information can be found in Package Development.

3.3 R Markdown

R markdown is a way of capturing documentation, code and results and in a single file. The document is written in plain text using a style called markdown. This has a simple syntax for specifying text formatting. R code is added in “chunks” and when the document is rendered the R code is executed and replaced with the results.

R markdown can be used to produce web pages, Word and PDF documents. The provide a robust way of capturing an analysis and the results and can be re-run when the data changes.

RStudio provides a cheatsheet detailing R markdown functionality:

3.3.1 Markdown

Markdown is a lightweight markup language with plain text formatting syntax. It is designed so that it can be converted to HTML and many other formats.

3.3.1.1 Paragraphs

Leave at least one empty line between text to start a new paragraph.

This is the first paragraph.

This is the second paragraph.

This is the first paragraph.

This is the second paragraph.

3.3.1.2 Headers

Header 1

Header 2

Header 3

3.3.1.3 Emphasis

italic ****bold****

italic **__bold__**

italic **bold**

italic **bold**

3.3.1.4 Lists

Unordered List:

- * Item 1
- * Item 2
 - + Item 2a
 - + Item 2b
- Item 1
- Item 2
 - Item 2a
 - Item 2b

Ordered List:

1. Item 1
2. Item 2
3. Item 3
 - a. Item 3a
 - b. Item 3b
1. Item 1
2. Item 2
3. Item 3
 - a. Item 3a
 - b. Item 3b

3.3.1.5 Links

Use a plain http address or add a link to a phrase:

```
http://example.com
```

```
[linked phrase](http://example.com)
```

```
http://example.com
```

```
linked phrase
```

3.3.1.6 Images

Images on the web or local files in the same directory:

```

```

```
![optional caption text](images/octocat.png)
```

3.3.1.7 Reference Style Links and Images

Links

```
A [linked phrase][id].
```

At the bottom of the document:

```
[id]: http://example.com/ "Title"
```

A linked phrase.

At the bottom of the document:

Images

```
![alt text][id]
```

At the bottom of the document:

```
[id]: images/octocat.png "Octocat"
```

At the bottom of the document:

3.3.1.8 Blockquotes

A friend once said:

```
> It's always better to give
> than to receive.
```

A friend once said:

It's always better to give than to receive.



Figure 3.2:



Figure 3.3: optional caption text



Figure 3.4: alt text

3.3.1.9 Plain Code Blocks

Plain code blocks are displayed in a fixed-width font but not evaluated

```
...
This text is displayed verbatim / preformatted
...
```

This text is displayed verbatim / preformatted

3.3.1.10 Inline Code

We defined the ``add`` function to compute the sum of two numbers.

We defined the `add` function to compute the sum of two numbers.

3.3.1.11 LaTeX Equations

Inline equation:

Einstein's famous equation `$E = mc^2$`

Einstein's famous equation $E = mc^2$

Display equation:

```
$$
E = mc^2
$$
```

$$E = mc^2$$

3.3.1.12 Horizontal Rule / Page Break

Three or more asterisks or dashes:

```
*****
```

```
-----
```



3.3.1.13 Tables

First Header	Second Header
Content Cell	Content Cell
Content Cell	Content Cell

First Header	Second Header
Content Cell	Content Cell
Content Cell	Content Cell

3.3.1.14 Manual Line Breaks

End a line with a backslash:

```
Roses are red,\
Violets are blue.
```

```
Roses are red,
Violets are blue.
```

3.3.1.15 Miscellaneous

```
superscript^2^
```

```
~~strikethrough~~
```

```
superscript2
```

```
strikethrough
```

3.3.2 R Markdown

3.3.2.1 R Code Chunks

R code will be evaluated and printed

```
summary(cars$dist)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      2.00   26.00   36.00   42.98   56.00  120.00
```

```
summary(cars$speed)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##       4.0   12.0   15.0   15.4   19.0   25.0
```

Inline R Code

There were 50 cars studied

3.3.3 R Notebooks

Chapter 4

Recommended Packages

4.1 Data Wrangling

Data wrangling is the process of transforming and mapping data from one “raw” data form into another format with the intent of making it more appropriate and valuable for a variety of downstream purposes such as analytics.

– *Wikipedia*

4.1.1 Tidy Data

The Tidyverse of packages are built around the concept of tidy data, first introduced by Jeff Leek in his book *The Elements of Data Analytic Style*. Hadley Wickham summarises the characteristics of tidy data with the following points:

1. Each variable forms a column.
2. Each observation forms a row.
3. Each type of observational unit forms a table.

4.1.2 Cheatsheets

RStudio provide a selection of cheatsheets containing quick reference of R functions useful for common task:

4.1.3 Importing Data

The first step in wrangling is importing the data into R. The methods to do so depend on the source of data:

- Reading files
- Connecting to databases
- Web APIs or pages

4.1.3.1 Reading Files

Files may come in many formats and R has packages to read many of them. The most common for data science are tabular text files, such as CSVs, and Excel spreadsheets. The recommended packages for reading these files are:.

Data Import : : CHEAT SHEET

R's **tidyverse** is built around **tidy data** stored in **tibbles**, which are enhanced data frames.

The front side of this sheet shows how to read text files into R with **readr**.

The reverse side shows how to create tibbles with **tibble** and to layout tidy data with **tidyr**.

OTHER TYPES OF DATA

Try one of the following packages to import other types of files

- **haven** - SPSS, Stata, and SAS files
- **readxl** - excel files (xls and xlsx)
- **DBI** - databases
- **jsonlite** - json
- **xml2** - XML
- **httr** - Web APIs
- **rvest** - HTML (Web Scraping)

Save Data

Save **x**, an R object, to **path**, a file path, as:

Comma delimited file
write_csv(**path**, **na** = "NA", **append** = FALSE, **col_names** = **append**)

File with arbitrary delimiter
write_delim(**x**, **path**, **delim** = " ", **na** = "NA", **append** = FALSE, **col_names** = **append**)

CSV for excel
write_excel_csv(**path**, **na** = "NA", **append** = FALSE, **col_names** = **append**)

String to file
write_file(**x**, **path**, **append** = FALSE)

String vector to file, one element per line
write_lines(**x**, **path**, **na** = "NA", **append** = FALSE)

Object to RDS file
write_rds(**x**, **path**, **compress** = c("none", "gz", "bz2", "xz", ...))

Tab delimited files
write_tsv(**x**, **path**, **na** = "NA", **append** = FALSE, **col_names** = **append**)

Read Tabular Data - These functions share the common arguments:

read_file(**file**, **col_names** = TRUE, **col_types** = NULL, **locale** = default_locale(), **na** = c("", "NA"), **quoted_na** = TRUE, **comment** = "", **trim_ws** = TRUE, **skip** = 0, **n_max** = Inf, **guess_max** = min(1000, **n_max**), **progress** = interactive())

a,b,c 1,2,3 4,5,NA	→	A B C 1 2 3 4 5 NA
a,b,c 1,2,3 4,5,NA	→	A B C 1 2 3 4 5 NA
a b c 1 2 3 4 5 NA	→	A B C 1 2 3 4 5 NA
a b c 1 2 3 4 5 NA	→	A B C 1 2 3 4 5 NA

Comma Delimited Files
read_csv("file.csv")
To make file.csv run:
write_file(**x** = "a,b,c(1,2,3)n4,5,NA", **path** = "file.csv")

Semi-colon Delimited Files
read_csv2("file2.csv")
write_file(**x** = "a;b;c(1,2,3)n4,5,NA", **path** = "file2.csv")

Files with Any Delimiter
read_delim("file.txt", **delim** = "|")
write_file(**x** = "a|b|c(1|2|3)n4|5|NA", **path** = "file.txt")

Fixed Width Files
read_fwf("file.fwf", **col_positions** = c(1, 3, 5))
write_file(**x** = "a b c(1 2 3)n4 5 NA", **path** = "file.fwf")

Tab Delimited Files
read_tsv("file.tsv") Also **read_table**()
write_file(**x** = "a\tb\tc(1\t2\t3)n4\t5\tNA", **path** = "file.tsv")

USEFUL ARGUMENTS

a,b,c 1,2,3 4,5,NA	Example file write_file ("a,b,c(1,2,3)n4,5,NA", "file.csv") f <- "file.csv"	1 2 3 4 5 NA	Skip lines read_csv (f, skip = 1)
A B C 1 2 3 4 5 NA	No header read_csv (f, col_names = FALSE)	A B C 1 2 3	Read in a subset read_csv (f, n_max = 1)
A B C 1 2 3 4 5 NA	Provide header read_csv (f, col_names = c("x", "y", "z"))	A B C NA 2 3 4 5 NA	Missing Values read_csv (f, na = c("1", ""))

Read Non-Tabular Data

Read a file into a single string
read_file(**file**, **locale** = default_locale())

Read each line into its own string
read_lines(**file**, **skip** = 0, **n_max** = -1L, **na** = character(), **locale** = default_locale(), **progress** = interactive())

Read Apache style log files
read_log(**file**, **col_names** = FALSE, **col_types** = NULL, **skip** = 0, **n_max** = -1, **progress** = interactive())

Read a file into a raw vector
read_file_raw(**file**)

Read each line into a raw vector
read_lines_raw(**file**, **skip** = 0, **n_max** = -1L, **progress** = interactive())



Studio

RStudio® is a trademark of RStudio, Inc. • CC BY RStudio • info@rstudio.com • 844-448-1212 • rstudio.com • Learn more with tidyverse.org • readr 1.1.0 • tibble 1.2.12 • tidy 0.6.0 • Updated: 2017-01

Figure 4.1: Data Import

Data Transformation with dplyr : : CHEAT SHEET

dplyr functions work with pipes and expect **tidy data**. In tidy data:

Each variable is in its own column
Each observation, or case, is in its own row
x %>% f(y) becomes **f(x, y)**

Summarise Cases

These apply **summary functions** to columns to create a new table. Summary functions take vectors as input and return one value (see back).

summary function
summarise(**data**, ...) Compute table of summaries. Also **summarise_()**.
summarise(**mtcars**, **avg** = mean(**mpg**))
count(**x**, ..., **wt** = NULL, **sort** = FALSE) Count number of rows in each group defined by the variables in ... Also **tally()**.
count(**iris**, **Species**)

VARIATIONS

summarise_all() - Apply funs to every column.
summarise_at() - Apply funs to specific columns.
summarise_if() - Apply funs to all cols of one type.

Group Cases

Use **group_by()** to create a "grouped" copy of a table. **dplyr** functions will manipulate each "group" separately and then combine the results.

mtcars %>%
group_by(**cy**) %>%
summarise(**avg** = mean(**mpg**))

group_by(**data**, ..., **add** = FALSE) Returns copy of table grouped by ...
g_iris <- **group_by**(**iris**, **Species**)
ungroup(**x**, ...) Returns ungrouped copy of table.
ungroup(**g_iris**)



Studio

RStudio® is a trademark of RStudio, Inc. • CC BY RStudio • info@rstudio.com • 844-448-1212 • rstudio.com • Learn more with browser/ignettes(package = c("dplyr", "tibble")) • dplyr 0.5.0 • tibble 1.2.0 • Updated: 2017-01

Figure 4.2: Data Transformation

Work with strings with stringr : : CHEAT SHEET

The **stringr** package provides a set of internally consistent tools for working with character strings, i.e. sequences of characters surrounded by quotation marks.

Detect Matches

str_detect(string, pattern) Detect the presence of a pattern match in a string. `str_detect(fruit, "a")`

str_which(string, pattern) Find the indexes of strings that contain a pattern match. `str_which(fruit, "a")`

str_count(string, pattern) Count the number of matches in a string. `str_count(fruit, "a")`

str_locate(string, pattern) Locate the positions of pattern matches in a string. Also `str_locate_all`. `str_locate(fruit, "a")`

Subset Strings

str_sub(string, start = 1L, end = -1L) Extract substrings from a character vector. `str_sub(fruit, 1, 3); str_sub(fruit, -2)`

str_subset(string, pattern) Return only the strings that contain a pattern match. `str_subset(fruit, "b")`

str_extract(string, pattern) Return the first pattern match found in each string, as a vector. Also `str_extract_all` to return every pattern match. `str_extract(fruit, "[aeiou]")`

str_match(string, pattern) Return the first pattern match found in each string, as a matrix with a column for each (1) group in pattern. Also `str_match_all`. `str_match(sentences, "(a|the) ([^]+)"`

Manage Lengths

str_length(string) The width of strings (i.e. number of code points, which generally equals the number of characters). `str_length(fruit)`

str_pad(string, width, side = c("left", "right", "both"), pad = " ") Pad strings to constant width. `str_pad(fruit, 17)`

str_trunc(string, width, side = c("right", "left", "center"), ellipsis = "...") Truncate the width of strings, replacing content with ellipsis. `str_trunc(fruit, 3)`

str_trim(string, side = c("both", "left", "right")) Trim whitespace from the start and/or end of a string. `str_trim(fruit)`

Mutate Strings

str_sub(<i>vector</i>, <i>value</i>) Replace substrings by identifying the substrings with `str_sub()` and assigning into the results. `str_sub(fruit, 1, 3) <- "str"`

str_replace(string, pattern, replacement) Replace the first matched pattern in each string. `str_replace(fruit, "a", "x")`

str_replace_all(string, pattern, replacement) Replace all matched patterns in each string. `str_replace_all(fruit, "a", "x")`

str_to_lower(string, locale = "en") Convert strings to lower case. `str_to_lower(sentences)`

str_to_upper(string, locale = "en") Convert strings to upper case. `str_to_upper(sentences)`

str_to_title(string, locale = "en") Convert strings to title case. `str_to_title(sentences)`

Join and Split

str_c(<i>vector</i>, sep = "", collapse = NULL) Join multiple strings into a single string. `str_c(letters, LETTERS)`

str_c(<i>vector</i>, sep = "", collapse = NULL) Collapse a vector of strings into a single string. `str_c(letters, collapse = ",")`

str_dup(string, times) Repeat strings times times. `str_dup(fruit, times = 2)`

str_split_fixed(string, pattern, n) Split a vector of strings into a matrix of substrings (splitting at occurrences of a pattern match). Also `str_split` to return a list of substrings. `str_split_fixed(fruit, "-", n = 2)`

glue::glue(<i>string</i>, ..., sep = "", envir = parent.frame(), open = "{", close = "}") Create a string from strings and (expressions) to evaluate. `glue::glue("Pi is {pi}")`

glue::glue_data(<i>x</i>, ..., sep = "", envir = parent.frame(), open = "{", close = "}") Use a data frame, list, or environment to create a string from strings and (expressions) to evaluate. `glue::glue_data(mtcars, "{row.names(mtcars)} has {hp} hp")`

Order Strings

str_order(x, decreasing = FALSE, na.last = TRUE, locale = "en", numeric = FALSE, ...) Return the vector of indexes that sorts a character vector. `str_order(x)`

str_sort(x, decreasing = FALSE, na.last = TRUE, locale = "en", numeric = FALSE, ...) Sort a character vector. `str_sort(x)`

Helpers

str_conv(string, encoding) Override the encoding of a string. `str_conv(fruit, "ISO-8859-1")`

str_view(string, pattern, match = NA) View HTML rendering of first regex match in each string. `str_view(fruit, "[aeiou]")`

str_view_all(string, pattern, match = NA) View HTML rendering of all regex matches. `str_view_all(fruit, "[aeiou]")`

str_wrap(string, width = 80, indent = 0, exdent = 0) Wrap strings into nicely formatted paragraphs. `str_wrap(sentences, 20)`

1 See bit.ly/ISO639-1 for a complete list of locales.

R Studio

RStudio is a trademark of RStudio, Inc. • CC BY RStudio • info@rstudio.com • 844-448-1212 • rstudio.com • Learn more at stringr.tidyverse.org • Diagrams from @kvaeder • stringr 1.2.0 • Updated: 2017-10

Figure 4.3: Strings

Dates and times with lubridate : : CHEAT SHEET

Date-times

2017-11-28 12:00:00

A **date-time** is a point on the timeline, stored as the number of seconds since 1970-01-01 00:00:00 UTC

`dt <- as_datetime(1511870400)`
"2017-11-28 12:00:00 UTC"

2017-11-28 12:00:00

PARSE DATE-TIMES (Convert strings or numbers to date-times)

- Identify the order of the year (y), month (m), day (d), hour (h), minute (m) and second (s) elements in your data
- Use the function below whose name replicates the order. Each accepts a wide variety of input formats.

2017-11-28T14:02:00 `ymd_hms()` `ymd_hms("2017-11-28T14:02:00")`

2017-22-12 10:00:00 `ymd_hms()` `ymd_hms("2017-22-12 10:00:00")`

11/28/2017 1:02:03 `mdy_hms()` `mdy_hms("11/28/2017 1:02:03")`

1 Jan 2017 23:59:59 `dmym_hms()` `dmym_hms("1 Jan 2017 23:59:59")`

20170131 `ymd()` `ymd("20170131")`

July 4th, 2000 `mdy()` `mdy("July 4th, 2000")`

4th of July 99 `dmym()` `dmym("4th of July 99")`

2001:Q3 `qq()` `qq("2001:Q3")`

2:01 `hms_hms()` `hms_hms("2:01")`

2017.5 `date_decimal()` `date_decimal("2017.5")`

now(tzone = "UTC") `now(tzone = "UTC")`

today(tzone = "UTC") `today(tzone = "UTC")`

fast_strptime() `fast_strptime("2017-11-28T14:02:00")`

parse_date_time() `parse_date_time("2017-11-28T14:02:00")`

Round Date-times

floor_date(x, unit = "second") Round down to nearest unit. `floor_date(dt, unit = "month")`

round_date(x, unit = "second") Round to nearest unit. `round_date(dt, unit = "month")`

ceiling_date(x, unit = "second") Round up to nearest unit. `ceiling_date(dt, unit = "month")`

rollback(dates, roll, to, first = FALSE, preserve_hms = TRUE) Roll back to last day of previous month. `rollback(dt)`

Stamp Date-times

stamp() Derive a template from an example string and return a new function that will apply the template to date-times. Also `stamp_date()` and `stamp_time()`.

- Derive a template, create a function `st <- stamp("Created Sunday, Jan 17, 1993 3:34")`
- Apply the template to dates `st(ymd("2010-04-05"))`
"11 Created Monday, Apr 05, 2010 00:00"

Time Zones

R recognizes ~600 time zones. Each encodes the time zone, Daylight Savings Time, and historical calendar variations for an area. R assigns one time zone per vector.

Use the UTC time zone to avoid Daylight Savings.

olsonNames() Returns a list of valid time zone names. `olsonNames()`

with_tz(time, tzone = "UTC") Get the same instant in a new time zone (a new clock time). `with_tz(dt, "US/Pacific")`

force_tz(time, tzone = "UTC") Get the same clock time in a new time zone (a new instant). `force_tz(dt, "US/Pacific")`

R Studio

RStudio is a trademark of RStudio, Inc. • CC BY RStudio • info@rstudio.com • 844-448-1212 • rstudio.com • Learn more at lubridate.tidyverse.org • lubridate 1.6.0 • Updated: 2017-12

Figure 4.4: Dates and Times

- `readr` for reading text files
- `readxl` for reading Excel files
- `openxlsx` for writing to Excel files

4.1.3.2 Connecting to Databases

The new RStudio Connections Pane makes it possible to easily connect to a variety of data sources, and explore the objects and data inside the connection.

The recommended packages for connecting to databases (also used by RStudio) are:

- `DBI` provides a standard interface to any database
- `odbc` for connecting to databases using ODBC
- `dplyr` for transforming tables in a database

4.1.3.3 Web APIs and Pages

Obtaining data from the internet has two main approaches. If the website provides an application programming interface (API) then you can send and receive data through it. The data from web APIs is usually returned in JSON or XML format. Alternatively, you can scrape the website itself, extracting data from the pages. The recommended packages for these approaches are:

- `httr` for communicating with web APIs
- `jsonlite` for reading JSON formatted text
- `xml2` for reading XML formatted text
- `rvest` for scraping web pages

4.1.4 Tidying Data

Now you have the data you want, it probably requires some processing in order to get it into a structure that is useful for analysis. There are two main packages for this `tidyr` and `dplyr`:

4.1.4.1 `tidyr`

The focus of `tidyr` is to get the data into a tidy format. The main functions it provides to do this are:

- `gather()` takes multiple columns, and gathers them into key-value pairs: it makes “wide” data longer.
- `spread()` takes two columns (key & value) and spreads in to multiple columns, it makes “long” data wider.
- `separate()` and `extract()` pulls apart a column that represents multiple variables.
- `unite()` combines columns, useful if one is redundant.
- `separate_rows()` for separating a row that contains multiple observations.

In addition, `tidyr` provides a set of functions for dealing with implicit and explicit missing values:

- `drop_na()`: Drop rows containing missing values
- `replace_na()`: Replace missing values
- `fill()`: Fill in missing values.
- `complete()`: Complete a data frame with missing combinations of data.
- `expand()`, `crossing()`: Expand data frame to include all combinations of values

4.1.4.2 dplyr

`dplyr` is a grammar of data manipulation, providing a consistent set of verbs that help you solve the most common data manipulation challenges:

- `mutate()` adds new variables that are functions of existing variables.
- `select()` picks variables based on their names.
- `filter()` picks cases based on their values.
- `summarise()` reduces multiple values down to a single summary.
- `arrange()` changes the ordering of the rows.

These all combine naturally with `group_by()` which allows you to perform any operation “by group”.

`dplyr` also provides a set of functions for combining tables. Mutating joins combine tables based on matching a subset of common variables; Set operations expect the tables to have the same variables and combine observations like sets.

Mutating joins

- `inner_join(x, y)` only includes observations that match in both `x` and `y`.
- `left_join(x, y)` includes all observations in `x`, regardless of whether they match or not.
- `right_join(x, y)` includes all observations in `y`. It’s equivalent to `left_join(y, x)`, but the columns will be ordered differently.
- `full_join()` includes all observations from `x` and `y`.

Set operations

- `intersect(x, y)`: returns only observations in both `x` and `y`.
- `union(x, y)`: returns unique observations in `x` and `y`.
- `setdiff(x, y)`: return observations in `x`, but not in `y`.

4.1.4.3 Variable Types

There are also packages which simplify working with specific types of data:

- `stringr` for strings and regular expressions
- `forcats` for factors, used to handle categorical data
- `lubridate` for dates and date-times.
- `hms` for time-of-day values.

4.1.4.4 tibble

When using these packages you may notice they return a `tibble` rather than a `data.frame`. They are basically the same thing with three key differences:

1. Printing: Tibbles only show the first ten rows and all the columns that fit on one screen.
2. Subsetting: `[` always returns another tibble. `$` never uses partial matching.
3. Recycling: Only values of length 1 are recycled. Trying to combine columns of different length throws an error.

The package `tibble` has some useful functions for constructing tibbles. In particular, you can construct a tibble row-wise using the `tribble` function (transposed tibble).