

# TRES Tidyverse Tutorial

Raphael, Pratik and Theo

2020-06-08



# Contents

2	<b>Outline</b>	<b>5</b>
3	About . . . . .	5
4	Schedule . . . . .	5
5	Possible extras . . . . .	6
6	Join . . . . .	6
7	<b>1 Reading files and string manipulation</b>	<b>7</b>
8	1.1 Data import and export with <code>readr</code> . . . . .	7
9	1.2 String manipulation with <code>stringr</code> . . . . .	10
10	1.3 String interpolation with <code>glue</code> . . . . .	18
11	1.4 Strings in <code>ggplot</code> . . . . .	20
12	<b>2 Reshaping data tables in the tidyverse, and other things</b>	<b>23</b>
13	2.1 The new data frame: <code>tibble</code> . . . . .	24
14	2.2 The concept of tidy data . . . . .	26
15	2.3 Reshaping with <code>tidyr</code> . . . . .	28
16	2.4 Extra: factors and the <code>forcats</code> package . . . . .	35
17	2.5 External resources . . . . .	40
18	<b>3 Data manipulation with <code>dplyr</code></b>	<b>41</b>
19	3.1 Introduction . . . . .	41
20	3.2 Example data of the day . . . . .	41
21	3.3 Select variables with <code>select()</code> . . . . .	43
22	3.4 Select observations with <code>filter()</code> . . . . .	43
23	3.5 Create new variables with <code>mutate()</code> . . . . .	43
24	3.6 Grouped results with <code>group_by()</code> and <code>summarise()</code> . . . . .	44
25	3.7 Scoped variables . . . . .	44
26	3.8 More ! . . . . .	44
27	<b>4 Working with lists and iteration</b>	<b>45</b>
28	4.1 List columns with <code>tidyr</code> . . . . .	45
29	4.2 Iteration with <code>map</code> . . . . .	48
30	4.3 More <code>map</code> variants . . . . .	52
31	4.4 Other functions for working with lists . . . . .	55

32	4.5 To add: patchwork . . . . .	59
----	---------------------------------	----

# Outline

This is the readable version of the TRES tidyverse tutorial. A convenient PDF version can be downloaded by clicking the PDF document icon in the header bar.

## About

The TRES tidyverse tutorial is an online workshop on how to use the tidyverse, a set of packages in the R computing language designed at making data handling and plotting easier.

This tutorial will take the form of a one hour per week video stream via Google Meet, every Friday morning at 10.00 (Groningen time) starting from the 29th of May, 2020 and lasting for a couple of weeks (depending on the number of topics we want to cover, but there should be at least 5).

**PhD students from outside our department are welcome to attend.**

## Schedule

Topic	Package	Instructor	Date*
Reading data and string manipulation	readr, stringr, glue	Pratik	29/05/20
Data and reshaping	tibble, tidyr	Raphael	05/06/20
Manipulating data	dplyr	Theo	12/06/20
Working with lists and iteration	purrr	Pratik	19/06/20
Plotting	ggplot2	Raphael	26/06/20
Regular expressions	regex	Richel	03/07/20
Programming with the tidyverse	rlang	Pratik	10/07/20

## 47 Possible extras

- 48 • Reproducibility and package-making (with e.g. usethis)
- 49 • Embedding C++ code with Rcpp

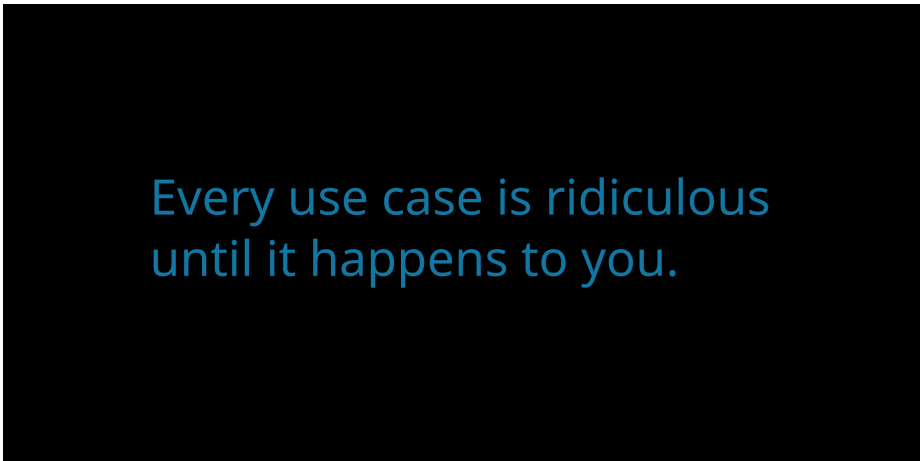
## 51 Join

52 Join the Slack by clicking this link (Slack account required).

53 \*Tentative dates.

## Chapter 1

# Reading files and string manipulation



Every use case is ridiculous  
until it happens to you.

Load the packages for the day.

```
library(readr)
library(stringr)
library(glue)
```

### 1.1 Data import and export with readr

Data in the wild with which ecologists and evolutionary biologists deal is most often in the form of a text file, usually with the extensions `.csv` or `.txt`. Often, such data has to be written to file from within R. `readr` contains a number of

functions to help with reading and writing text files.

### 1.1.1 Reading data

Reading in a csv file with `readr` is done with the `read_csv` function, a faster alternative to the base R `read.csv`. Here, `read_csv` is applied to the `mtcars` example.

```
# get the filepath of the example
some_example = readr_example("mtcars.csv")

# read the file in
some_example = read_csv(some_example)

head(some_example)
#> # A tibble: 6 x 11
#>   mpg   cyl  disp    hp  drat    wt   qsec    vs  am  gear  carb
#>   <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
#> 1  21     6   160   110   3.9   2.62  16.5     0    1     4     4
#> 2  21     6   160   110   3.9   2.88  17.0     0    1     4     4
#> 3 22.8     4   108    93   3.85   2.32  18.6     1    1     4     1
#> 4 21.4     6   258   110   3.08   3.22  19.4     1    0     3     1
#> 5 18.7     8   360   175   3.15   3.44  17.0     0    0     3     2
#> 6 18.1     6   225   105   2.76   3.46  20.2     1    0     3     1
```

The `read_csv2` function is useful when dealing with files where the separator between columns is a semicolon `;`, and where the decimal point is represented by a comma `,`.

Other variants include:

- `read_tsv` for tab-separated files, and
- `read_delim`, a general case which allows the separator to be specified manually.

`readr` import function will attempt to guess the column type from the first  $N$  lines in the data. This  $N$  can be set using the function argument `guess_max`. The `n_max` argument sets the number of rows to read, while the `skip` argument sets the number of rows to be skipped before reading data.

By default, the column names are taken from the first row of the data, but they can be manually specified by passing a character vector to `col_names`.

There are some other arguments to the data import functions, but the defaults usually *just work*.



### 83 1.1.2 Writing data

84 Writing data uses the `write_*` family of functions, with implementations for  
 85 `csv`, `csv2` etc. (represented by the asterisk), mirroring the import functions  
 86 discussed above. `write_*` functions offer the `append` argument, which allow a  
 87 data frame to be added to an existing file.

88 These functions are not covered here.

### 89 1.1.3 Reading and writing lines

90 Sometimes, there is text output generated in R which needs to be written to file,  
 91 but is not in the form of a dataframe. A good example is model outputs. It is  
 92 good practice to save model output as a text file, and add it to version control.  
 93 Similarly, it may be necessary to import such text, either for display to screen,  
 94 or to extract data.

95 This can be done using the `readr` functions `read_lines` and `write_lines`. Con-  
 96 sider the model summary from a simple linear model.

```
# get the model
model = lm(mpg ~ wt, data = mtcars)
```

97 The model summary can be written to file. When writing lines to file, BE  
 98 AWARE OF THE DIFFERENCES BETWEEN UNIX AND WINODWS line  
 99 separators. Usually, this causes no trouble.

```
# capture the model summary output
model_output = capture.output(summary(model))
```

```
# save it to file
write_lines(x = model_output,
  path = "model_output.txt")
```

100 This model output can be read back in for display, and each line of the model  
 101 output is an element in a character vector.

```
# read in the model output and display
model_output = read_lines("model_output.txt")

# use cat to show the model output as it would be on screen
cat(model_output, sep = "\n")
#>
#> Call:
#> lm(formula = mpg ~ wt, data = mtcars)
#>
#> Residuals:
#>      Min       1Q   Median       3Q      Max
#> -4.543  -2.365  -0.125   1.410   6.873
#>
```

```

#> Coefficients:
#>               Estimate Std. Error t value Pr(>|t|)
#> (Intercept)   37.285      1.878   19.86 < 2e-16 ***
#> wt            -5.344      0.559   -9.56 1.3e-10 ***
#> ---
#> Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
#>
#> Residual standard error: 3.05 on 30 degrees of freedom
#> Multiple R-squared:  0.753, Adjusted R-squared:  0.745
#> F-statistic: 91.4 on 1 and 30 DF,  p-value: 1.29e-10

```

These few functions demonstrate the most common uses of `readr`, but most other use cases for text data can be handled using different function arguments, including reading data off the web, unzipping compressed files before reading, and specifying the column types to control for type conversion errors.

## Excel files

Finally, data is often shared or stored by well meaning people in the form of Microsoft Excel sheets. Indeed, Excel (especially when synced regularly to remote storage) is a good way of noting down observational data in the field. The `readxl` package allows importing from Excel files, including reading in specific sheets.

## 1.2 String manipulation with `stringr`

`stringr` is the tidyverse package for string manipulation, and exists in an interesting symbiosis with the `stringi` package. For the most part, `stringr` is a wrapper around `stringi`, and is almost always more than sufficient for day-to-day needs.

`stringr` functions begin with `str_`.

### 1.2.1 Putting strings together

Concatenate two strings with `str_c`, and duplicate strings with `str_dup`. Flatten a list or vector of strings using `str_flatten`.

```

# str_c works like paste(), choose a separator
str_c("this string", "this other string", sep = "_")
#> [1] "this string_this other string"

# str_dup works like rep
str_dup("this string", times = 3)
#> [1] "this stringthis stringthis string"

# str_flatten works on lists and vectors

```

```

str_flatten(string = as.list(letters), collapse = "_")
#> [1] "a_b_c_d_e_f_g_h_i_j_k_l_m_n_o_p_q_r_s_t_u_v_w_x_y_z"
str_flatten(string = letters, collapse = "-")
#> [1] "a-b-c-d-e-f-g-h-i-j-k-l-m-n-o-p-q-r-s-t-u-v-w-x-y-z"

```

121 `str_flatten` is especially useful when displaying the type of an object that  
 122 returns a list when `class` is called on it.

```

# get the class of a tibble and display it as a single string
class_tibble = class(tibble::tibble(a = 1))
str_flatten(string = class_tibble, collapse = ", ")
#> [1] "tbl_df, tbl, data.frame"

```

### 123 1.2.2 Detecting strings

124 Count the frequency of a pattern in a string with `str_count`. Returns an integer.  
 125 Detect whether a pattern exists in a string with `str_detect`. Returns a logical  
 126 and can be used as a predicate.  
 127 Both are vectorised, i.e, automatically applied to a vector of arguments.

```

# there should be 5 a-s here
str_count(string = "ababababa", pattern = "a")
#> [1] 5

# vectorise over the input string
# should return a vector of length 2, with integers 5 and 3
str_count(string = c("ababbababa", "banana"), pattern = "a")
#> [1] 5 3

# vectorise over the pattern to count both a-s and b-s
str_count(string = "ababababa", pattern = c("a", "b"))
#> [1] 5 4

```

128 Vectorising over both string and pattern works as expected.

```

# vectorise over both string and pattern
# counts a-s in first input, and b-s in the second
str_count(string = c("ababababa", "banana"),
          pattern = c("a", "b"))
#> [1] 5 1

# provide a longer pattern vector to search for both a-s
# and b-s in both inputs
str_count(string = c("ababababa", "banana"),
          pattern = c("a", "b",
                      "b", "a"))
#> [1] 5 1 4 3

```

129 **str\_locate** locates the search pattern in a string, and returns the start and  
 130 end as a two column matrix.

```
# the behaviour of both str_locate and str_locate_all is
# to find the first match by default
str_locate(string = "banana", pattern = "ana")
#>      start end
#> [1,]      2  4

# str_detect detects a sequence in a string
str_detect(string = "Bananageddon is coming!",
           pattern = "na")
#> [1] TRUE

# str_detect is also vectorised and returns a two-element logical vector
str_detect(string = "Bananageddon is coming!",
           pattern = c("na", "don"))
#> [1] TRUE TRUE

# use any or all to convert a multi-element logical to a single logical
# here we ask if either of the patterns is detected
any(str_detect(string = "Bananageddon is coming!",
               pattern = c("na", "don")))
#> [1] TRUE
```

131 Detect whether a string starts or ends with a pattern. Also vectorised. Both  
 132 have a **negate** argument, which returns the negative, i.e., returns FALSE if the  
 133 search pattern is detected.

```
# taken straight from the examples, because they suffice
fruit <- c("apple", "banana", "pear", "pineapple")
# str_detect looks at the first character
str_starts(fruit, "p")
#> [1] FALSE FALSE TRUE TRUE

# str_ends looks at the last character
str_ends(fruit, "e")
#> [1] TRUE FALSE FALSE TRUE

# an example of negate = TRUE
str_ends(fruit, "e", negate = TRUE)
#> [1] FALSE TRUE TRUE FALSE
```

134 **str\_subset** [WHICH IS NOT RELATED TO **str\_sub**] helps with subsetting a  
 135 character vector based on a **str\_detect** predicate. In the example, all elements  
 136 containing “banana” are subset.

137 **str\_which** has the same logic except that it returns the vector position and not  
 138 the elements.

```

# should return a subset vector containing the first two elements
str_subset(c("banana",
             "bananageddon is coming",
             "applegeddon is not real"),
           pattern = "banana")
#> [1] "banana"                "bananageddon is coming"

# returns an integer vector
str_which(c("banana",
            "bananageddon is coming",
            "applegeddon is not real"),
          pattern = "banana")
#> [1] 1 2

```

### 1.2.3 Matching strings

139 `str_match` returns all positive matches of the pattern in the string. The return  
 140 type is a list, with one element per search pattern.

142 A simple case is shown below where the search pattern is the phrase “banana”.

```

str_match(string = c("banana",
                     "bananageddon",
                     "bananas are bad"),
          pattern = "banana")
#>      [,1]
#> [1,] "banana"
#> [2,] "banana"
#> [3,] "banana"

```

143 The search pattern can be extended to look for multiple subsets of the search  
 144 pattern. Consider searching for dates and times.

145 Here, the search pattern is a `regex` pattern that looks for a set of four digits  
 146 (`\\d{4}`) and a month name (`\\w+`) separated by a hyphen. There’s much more  
 147 to be explored in dealing with dates and times in `lubridate`, another `tidyverse`  
 148 package.

149 The return type is a list, each element is a character matrix where the first  
 150 column is the string subset matching the full search pattern, and then as many  
 151 columns as there are parts to the search pattern. The parts of interest in the  
 152 search pattern are indicated by wrapping them in parentheses. For example, in  
 153 the case below, wrapping `[-.]` in parentheses will turn it into a distinct part  
 154 of the search pattern.

```

# first with [-.] treated simply as a separator
str_match(string = c("1970-somemonth-01",
                     "1990-anothermonth-01",
                     "2010-thismonth-01"),

```

```

        pattern = "\\d{4}][-.](\\w+)"
#>      [,1]      [,2]      [,3]
#> [1,] "1970-somemonth" "1970" "somemonth"
#> [2,] "1990-anothermonth" "1990" "anothermonth"
#> [3,] "2010-thismonth" "2010" "thismonth"

# then with [-.] actively searched for
str_match(string = c("1970-somemonth-01",
                     "1990-anothermonth-01",
                     "2010-thismonth-01"),
          pattern = "\\d{4}([-.])(\\w+)"
#>      [,1]      [,2]      [,3] [,4]
#> [1,] "1970-somemonth" "1970" "-" "somemonth"
#> [2,] "1990-anothermonth" "1990" "-" "anothermonth"
#> [3,] "2010-thismonth" "2010" "-" "thismonth"

```

155 Multiple possible matches are dealt with using `str_match_all`. An example  
 156 case is uncertainty in date-time in raw data, where the date has been entered  
 157 as 1970-somemonth-01 or 1970/anothermonth/01.

158 The return type is a list, with one element per input string. Each element is a  
 159 character matrix, where each row is one possible match, and each column after  
 160 the first (the full match) corresponds to the parts of the search pattern.

```

# first with a single date entry
str_match_all(string = c("1970-somemonth-01 or maybe 1990/anothermonth/01"),
              pattern = "\\d{4}[\\-\\/](\\w+)"
#> [[1]]
#>      [,1]      [,2]      [,3]
#> [1,] "1970-somemonth" "1970" "somemonth"
#> [2,] "1990/anothermonth" "1990" "anothermonth"

# then with multiple date entries
str_match_all(string = c("1970-somemonth-01 or maybe 1990/anothermonth/01",
                        "1990-somemonth-01 or maybe 2001/anothermonth/01"),
              pattern = "\\d{4}[\\-\\/](\\w+)"
#> [[1]]
#>      [,1]      [,2]      [,3]
#> [1,] "1970-somemonth" "1970" "somemonth"
#> [2,] "1990/anothermonth" "1990" "anothermonth"
#>
#> [[2]]
#>      [,1]      [,2]      [,3]
#> [1,] "1990-somemonth" "1990" "somemonth"
#> [2,] "2001/anothermonth" "2001" "anothermonth"

```

### 1.2.4 Simpler pattern extraction

The full functionality of `str_match_*` can be boiled down to the most common use case, extracting one or more full matches of the search pattern using `str_extract` and `str_extract_all` respectively.

`str_extract` returns a character vector with the same length as the input string vector, while `str_extract_all` returns a list, with a character vector whose elements are the matches.

```
# extracting the first full match using str_extract
str_extract(string = c("1970-somemonth-01 or maybe 1990/anothermonth/01",
                        "1990-somemonth-01 or maybe 2001/anothermonth/01"),
             pattern = "(\\d{4})[\\-\\/]([a-z]+)")
#> [1] "1970-somemonth" "1990-somemonth"

# extracting all full matches using str_extract_all
str_extract_all(string = c("1970-somemonth-01 or maybe 1990/anothermonth/01",
                           "1990-somemonth-01 or maybe 2001/anothermonth/01"),
                 pattern = "(\\d{4})[\\-\\/]([a-z]+)")
#> [[1]]
#> [1] "1970-somemonth" "1990/anothermonth"
#>
#> [[2]]
#> [1] "1990-somemonth" "2001/anothermonth"
```

### 1.2.5 Breaking strings apart

`str_split`, `str_sub`, In the above date-time example, when reading filenames from a path, or when working sequences separated by a known pattern generally, `str_split` can help separate elements of interest.

The return type is a list similar to `str_match`.

```
# split on either a hyphen or a forward slash
str_split(string = c("1970-somemonth-01",
                     "1990/anothermonth/01"),
           pattern = "[\\-\\/]")
#> [[1]]
#> [1] "1970" "somemonth" "01"
#>
#> [[2]]
#> [1] "1990" "anothermonth" "01"
```

This can be useful in recovering simulation parameters from a filename, but may require some knowledge of `regex`.

```
# assume a simulation output file
filename = "sim_param1_0.01_param2_0.05_param3_0.01.ext"
```

```

# not quite there
str_split(filename, pattern = "_")
#> [[1]]
#> [1] "sim"      "param1"    "0.01"      "param2"    "0.05"      "param3"    "0.01.ext"

# not really
str_split(filename,
            pattern = "sim_")
#> [[1]]
#> [1] ""
#> [2] "param1_0.01_param2_0.05_param3_0.01.ext"

# getting there but still needs work
str_split(filename,
            pattern = "(sim_)|_*param\\d{1}|(.ext)")
#> [[1]]
#> [1] ""      ""      "0.01"  "0.05"  "0.01"  ""

```

175 **str\_split\_fixed** split the string into as many pieces as specified, and can be  
 176 especially useful dealing with filepaths.

```

# split on either a hyphen or a forward slash
str_split_fixed(string = "dir_level_1/dir_level_2/file.ext",
                 pattern = "/",
                 n = 2)
#>      [,1]      [,2]
#> [1,] "dir_level_1" "dir_level_2/file.ext"

```

### 177 1.2.6 Replacing string elements

178 **str\_replace** is intended to replace the search pattern, and can be co-opted  
 179 into the task of recovering simulation parameters or other data from regularly  
 180 named files. **str\_replace\_all** works the same way but replaces all matches of  
 181 the search pattern.

```

# replace all unwanted characters from this hypothetical filename with spaces
filename = "sim_param1_0.01_param2_0.05_param3_0.01.ext"
str_replace_all(filename,
                 pattern = "(sim_)|_*param\\d{1}|(.ext)",
                 replacement = " ")
#> [1] " 0.01 0.05 0.01 "

```

182 **str\_remove** is a wrapper around **str\_replace** where the replacement is set to  
 183 `""`. This is not covered here.

184 Having replaced unwanted characters in the filename with spaces, **str\_trim**  
 185 offers a way to remove leading and trailing whitespaces.



```

# trim whitespaces from this filename after replacing unwanted text
filename = "sim_param1_0.01_param2_0.05_param3_0.01.ext"
filename_with_spaces = str_replace_all(filename,
                                     pattern = "(sim_)|_*param\\d{1}|(.ext)",
                                     replacement = " ")
filename_without_spaces = str_trim(filename_with_spaces)
filename_without_spaces
#> [1] "0.01 0.05 0.01"

# the result can be split on whitespaces to return useful data
str_split(filename_without_spaces, " ")
#> [[1]]
#> [1] "0.01" "0.05" "0.01"

```

### 186 1.2.7 Subsetting within strings

187 When strings are highly regular, useful data can be extracted from a string using  
 188 `str_sub`. In the date-time example, the year is always represented by the first  
 189 four characters.

```

# get the year as characters 1 - 4
str_sub(string = c("1970-somemonth-01",
                  "1990-anothermonth-01",
                  "2010-thismonth-01"),
        start = 1, end = 4)
#> [1] "1970" "1990" "2010"

```

190 Similarly, it's possible to extract the last few characters using negative indices.

```

# get the day as characters -2 to -1
str_sub(string = c("1970-somemonth-01",
                  "1990-anothermonth-21",
                  "2010-thismonth-31"),
        start = -2, end = -1)
#> [1] "01" "21" "31"

```

191 Finally, it's also possible to replace characters within a string based on the  
 192 position. This requires using the assignment operator `<-`.

```

# replace all days in these dates to 01
date_times = c("1970-somemonth-25",
               "1990-anothermonth-21",
               "2010-thismonth-31")

# a strictly necessary use of the assignment operator
str_sub(date_times,
        start = -2, end = -1) <- "01"

```

```
date_times
#> [1] "1970-somemonth-01" "1990-anothermonth-01" "2010-thismonth-01"
```

### 193 1.2.8 Padding and truncating strings

194 Strings included in filenames or plots are often of unequal lengths, especially  
 195 when they represent numbers. `str_pad` can pad strings with suitable characters  
 196 to maintain equal length filenames, with which it is easier to work.

```
# pad so all values have three digits
str_pad(string = c("1", "10", "100"),
        width = 3,
        side = "left",
        pad = "0")
#> [1] "001" "010" "100"
```

197 Strings can also be truncated if they are too long.

```
str_trunc(string = c("bananas are great and wonderful
                      and more stuff about bananas and
                      it really goes on about bananas"),
          width = 27,
          side = "right", ellipsis = "etc. etc.")
#> [1] "bananas are great etc. etc."
```

### 198 1.2.9 Stringr aspects not covered here

199 Some `stringr` functions are not covered here. These include:

- 200 • `str_wrap` (of dubious use),
- 201 • `str_interp`, `str_glue*` (better to use `glue`; see below),
- 202 • `str_sort`, `str_order` (used in sorting a character vector),
- 203 • `str_to_case*` (case conversion), and
- 204 • `str_view*` (a graphical view of search pattern matches).
- 205 • `word`, `boundary` etc. The use of `word` is covered below.

206 `stringi`, of which `stringr` is a wrapper, offers a lot more flexibility and control.

## 207 1.3 String interpolation with glue

208 The idea behind string interpolation is to procedurally generate new complex  
 209 strings from pre-existing data.

210 `glue` is as simple as the example shown.

```
# print that each car name is a car model
cars = rownames(head(mtcars))
glue('The {cars} is a car model')
#> The Mazda RX4 is a car model
#> The Mazda RX4 Wag is a car model
#> The Datsun 710 is a car model
#> The Hornet 4 Drive is a car model
#> The Hornet Sportabout is a car model
#> The Valiant is a car model
```

211 This creates and prints a vector of car names stating each is a car model.

212 The related `glue_data` is even more useful in printing from a dataframe. In  
 213 this example, it can quickly generate command line arguments or filenames.

```
# use dataframes for now
parameter_combinations = data.frame(param1 = letters[1:5],
                                     param2 = 1:5)

# for command line arguments or to start multiple job scripts on the cluster
glue_data(parameter_combinations,
           'simulation-name {param1} {param2}')
#> simulation-name a 1
#> simulation-name b 2
#> simulation-name c 3
#> simulation-name d 4
#> simulation-name e 5

# for filenames
glue_data(parameter_combinations,
           'sim_data_param1_{param1}_param2_{param2}.ext')
#> sim_data_param1_a_param2_1.ext
#> sim_data_param1_b_param2_2.ext
#> sim_data_param1_c_param2_3.ext
#> sim_data_param1_d_param2_4.ext
#> sim_data_param1_e_param2_5.ext
```

214 Finally, the convenient `glue_sql` and `glue_data_sql` are used to safely write  
 215 SQL queries where variables from data are appropriately quoted. This is not  
 216 covered here, but it is good to know it exists.

217 `glue` has some more functions — `glue_safe`, `glue_collapse`, and `glue_col`,  
 218 but these are infrequently used. Their functionality can be found on the `glue`  
 219 github page.

## 1.4 Strings in ggplot

`ggplot` has two `geoms` (wait for the `ggplot` tutorial to understand more about `geoms`) that work with text: `geom_text` and `geom_label`. These `geoms` allow text to be pasted on to the main body of a plot.

Often, these may overlap when the data are closely spaced. The package `ggrepel` offers another `geom`, `geom_text_repel` (and the related `geom_label_repel`) that help arrange text on a plot so it doesn't overlap with other features. This is *not perfect*, but it works more often than not.

More examples can be found on the `ggrepel` website.

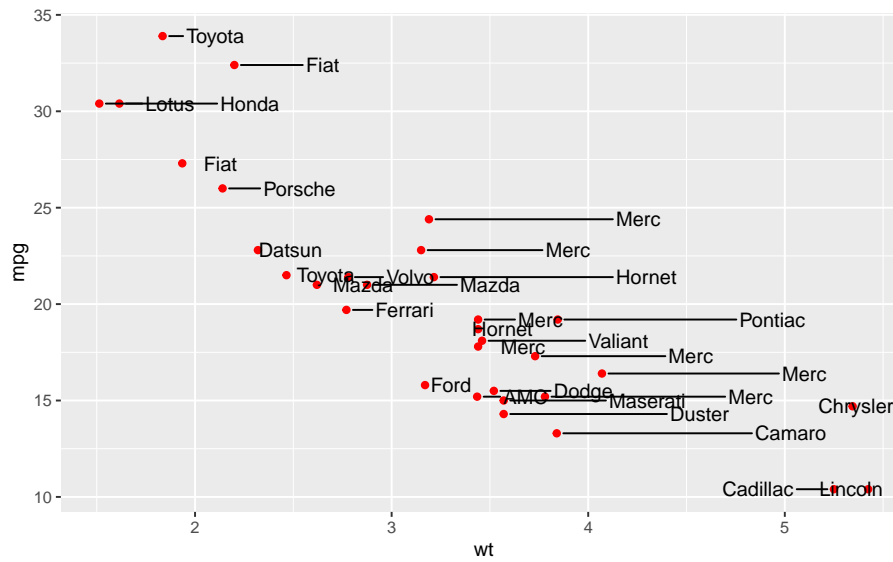
Here, the arguments to `geom_text_repel` are taken both from the `mtcars` data (position), as well as from the car brands extracted using the `stringr::word` (labels), which tries to separate strings based on a regular pattern.

The details of `ggplot` are covered in a later tutorial.

```
library(ggplot2)
library(ggrepel)

# prepare car labels using word function
car_labels = word(rownames(mtcars))

ggplot(mtcars,
       aes(x = wt, y = mpg,
           label = rownames(mtcars)))+
  geom_point(colour = "red")+
  geom_text_repel(aes(label = car_labels),
                 direction = "x",
                 nudge_x = 0.2,
                 box.padding = 0.5,
                 point.padding = 0.5)
```



235

236 This is not a good looking plot, because it breaks other rules of plot design,  
 237 such as whether this sort of plot should be made at all. Labels and text need  
 238 to be applied sparingly, for example drawing attention or adding information to  
 239 outliers.



## Chapter 2

# Reshaping data tables in the tidyverse, and other things

Raphael Scherrer

Every use case is ridiculous until it happens to you.

```
library(tibble)
library(tidyr)
```

In this chapter we will learn what *tidy* means in the context of the tidyverse, and how to reshape our data into a tidy format using the `tidyr` package. But first, let us take a detour and introduce the `tibble`.

## 2.1 The new data frame: tibble

The `tibble` is the recommended class to use to store tabular data in the tidyverse. Consider it as the operational unit of any data science pipeline. For most practical purposes, a `tibble` is basically a `data.frame`.

```
# Make a data frame
data.frame(who = c("Pratik", "Theo", "Raph"), chapt = c("1, 4", "3", "2, 5"))
#>      who chapt
#> 1 Pratik  1, 4
#> 2  Theo    3
#> 3  Raph   2, 5

# Or an equivalent tibble
tibble(who = c("Pratik", "Theo", "Raph"), chapt = c("1, 4", "3", "2, 5"))
#> # A tibble: 3 x 2
#>   who    chapt
#>   <chr> <chr>
#> 1 Pratik 1, 4
#> 2 Theo    3
#> 3 Raph   2, 5
```

The difference between `tibble` and `data.frame` is in its display and in the way it is subsetting, among others. Most functions working with `data.frame` will work with `tibble` and vice versa. Use the `as*` family of functions to switch back and forth between the two if needed, using e.g. `as.data.frame` or `as_tibble`.

In terms of display, the tibble has the advantage of showing the class of each column: `chr` for `character`, `fct` for `factor`, `int` for `integer`, `dbl` for `numeric` and `lgl` for `logical`, just to name the main atomic classes. This may be more important than you think, because many hard-to-find bugs in R are due to wrong variable types and/or cryptic type conversions. This especially happens with `factor` and `character`, which can cause quite some confusion. More about this in the extra section at the end of this chapter!

Note that you can build a tibble by rows rather than by columns with `tribble`:

```
tribble(
  ~who, ~chapt,
  "Pratik", "1, 4",
  "Theo", "3",
  "Raph", "2, 5"
)
#> # A tibble: 3 x 2
#>   who    chapt
#>   <chr> <chr>
#> 1 Pratik 1, 4
#> 2 Theo    3
```



```
#> 3 Raph 2, 5
```

As a rule of thumb, try to convert your tables to tibbles whenever you can, especially when the original table is *not* a data frame. For example, the principal component analysis function `prcomp` outputs a **matrix** of coordinates in principal component-space.

```
# Perform a PCA on mtcars
pca_scores <- prcomp(mtcars)$x
head(pca_scores) # looks like a data frame or a tibble...
#>           PC1  PC2  PC3  PC4  PC5  PC6  PC7  PC8
#> Mazda RX4      -79.60  2.13 -2.15 -2.707 -0.702 -0.3149 -0.09870 -0.0779
#> Mazda RX4 Wag  -79.60  2.15 -2.22 -2.178 -0.884 -0.4534 -0.00355 -0.0957
#> Datsun 710     -133.89 -5.06 -2.14  0.346  1.106  1.1730  0.00576  0.1362
#> Hornet 4 Drive   8.52 44.99  1.23  0.827  0.424 -0.0579 -0.02431  0.2212
#> Hornet Sportabout 128.69 30.82  3.34 -0.521  0.737 -0.3329  0.10630 -0.0530
#> Valiant        -23.22 35.11 -3.26  1.401  0.803 -0.0884  0.23895  0.4239
#>           PC9  PC10  PC11
#> Mazda RX4      -0.200 -0.2901  0.106
#> Mazda RX4 Wag  -0.353 -0.1928  0.107
#> Datsun 710     -0.198  0.0763  0.267
#> Hornet 4 Drive   0.356 -0.0906  0.209
#> Hornet Sportabout 0.153 -0.1886 -0.109
#> Valiant         0.101 -0.0377  0.276
class(pca_scores) # but is actually a matrix
#> [1] "matrix"

# Convert to tibble
as_tibble(pca_scores)
#> # A tibble: 32 x 11
#>       PC1  PC2  PC3  PC4  PC5  PC6  PC7  PC8  PC9  PC10
#>   <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
#> 1  -79.6  2.13 -2.15 -2.71 -0.702 -0.315 -0.0987 -0.0779 -0.200 -0.290
#> 2  -79.6  2.15 -2.22 -2.18 -0.884 -0.453 -0.00355 -0.0957 -0.353 -0.193
#> 3 -134.  -5.06 -2.14  0.346  1.11  1.17  0.00576  0.136 -0.198  0.0763
#> 4   8.52 45.0  1.23  0.827  0.424 -0.0579 -0.0243  0.221  0.356 -0.0906
#> 5  129.  30.8  3.34 -0.521  0.737 -0.333  0.106 -0.0530  0.153 -0.189
#> 6  -23.2 35.1 -3.26  1.40  0.803 -0.0884  0.239  0.424  0.101 -0.0377
#> # ... with 26 more rows, and 1 more variable: PC11 <dbl>
```

This is important because a **matrix** can contain only one type of values (e.g. only **numeric** or **character**), while **tibble** (and **data.frame**) allow you to have columns of different types.

So, in the tidyverse we are going to work with tibbles, got it. But what does “tidy” mean exactly?

## 2.2 The concept of tidy data

When it comes to putting data into tables, there are many ways one could organize a dataset. The *tidy* format is one such format. According to the formal definition, a table is tidy if each column is a variable and each row is an observation. In practice, however, I found that this is not a very operational definition, especially in ecology and evolution where we often record multiple variables per individual. So, let's dig in with an example.

Say we have a dataset of several morphometrics measured on Darwin's finches in the Galapagos islands. Let's first get this dataset.

```
# We first simulate random data
beak_lengths <- rnorm(100, mean = 5, sd = 0.1)
beak_widths <- rnorm(100, mean = 2, sd = 0.1)
body_weights <- rgamma(100, shape = 10, rate = 1)
islands <- rep(c("Isabela", "Santa Cruz"), each = 50)

# Assemble into a tibble
data <- tibble(
  id = 1:100,
  body_weight = body_weights,
  beak_length = beak_lengths,
  beak_width = beak_widths,
  island = islands
)

# Snapshot
data
#> # A tibble: 100 x 5
#>   id body_weight beak_length beak_width island
#>   <int>      <dbl>      <dbl>      <dbl> <chr>
#> 1     1      10.8        4.94      1.94 Isabela
#> 2     2      15.4        5.02      2.00 Isabela
#> 3     3      15.0        4.92      1.91 Isabela
#> 4     4       8.51        5.16      2.02 Isabela
#> 5     5      14.9        5.03      1.93 Isabela
#> 6     6       8.41        4.92      2.18 Isabela
#> # ... with 94 more rows
```

Here, we pretend to have measured `beak_length`, `beak_width` and `body_weight` on 100 birds, 50 of them from Isabela and 50 of them from Santa Cruz. In this tibble, each row is an individual bird. This is probably the way most scientists would record their data in the field. However, a single bird is not an “observation” in the sense used in the tidyverse. Our dataset is not tidy but *messy*.

The tidy equivalent of this dataset would be:

```

data <- pivot_longer(
  data,
  cols = c("body_weight", "beak_length", "beak_width"),
  names_to = "variable"
)
data
#> # A tibble: 300 x 4
#>   id island variable    value
#>   <int> <chr>   <chr>    <dbl>
#> 1     1  Isabela body_weight 10.8
#> 2     1  Isabela beak_length  4.94
#> 3     1  Isabela beak_width  1.94
#> 4     2  Isabela body_weight 15.4
#> 5     2  Isabela beak_length  5.02
#> 6     2  Isabela beak_width  2.00
#> # ... with 294 more rows

```

289 where each *measurement* (and not each *individual*) is now the unit of observation  
 290 (the rows). The `pivot_longer` function is the easiest way to get to this format.  
 291 It belongs to the `tidyr` package, which we'll cover in a minute.

292 As you can see our tibble now has three times as many rows and fewer columns.  
 293 This format is rather unintuitive and not optimal for display. However, it pro-  
 294 vides a very standardized and consistent way of organizing data that will be  
 295 understood (and expected) by pretty much all functions in the tidyverse. This  
 296 makes the tidyverse tools work well together and reduces the time you would  
 297 otherwise spend reformatting your data from one tool to the next.

298 That does not mean that the *messy* format is useless though. There may be  
 299 use-cases where you need to switch back and forth between formats. For this  
 300 reason I prefer referring to these formats using their other names: *long* (tidy)  
 301 versus *wide* (messy). For example, matrix operations work much faster on wide  
 302 data, and the wide format arguably looks nicer for display. Luckily the `tidyr`  
 303 package gives us the tools to reshape our data as needed, as we shall see shortly.

304 Another common example of wide-or-long dilemma is when dealing with *con-*  
 305 *tingency tables*. This would be our case, for example, if we asked how many  
 306 observations we have for each morphometric and each island. We use `table`  
 307 (from base R) to get the answer:

```

# Make a contingency table
ctg <- with(data, table(island, variable))
ctg
#>           variable
#> island  beak_length beak_width body_weight
#> Isabela           50          50          50
#> Santa Cruz        50          50          50

```

308 A variety of statistical tests can be used on contingency tables such as Fisher's

exact test, the chi-square test or the binomial test. Contingency tables are in the wide format by construction, but they too can be pivoted to the long format, and the tidyverse manipulation tools will expect you to do so. Actually, `tibble` knows that very well and does it by default if you convert your `table` into a `tibble`:

```
# Contingency table is pivoted to the long-format automatically
as_tibble(ctg)
#> # A tibble: 6 x 3
#>   island      variable      n
#>   <chr>      <chr>      <int>
#> 1 Isabela    beak_length    50
#> 2 Santa Cruz beak_length    50
#> 3 Isabela    beak_width     50
#> 4 Santa Cruz beak_width     50
#> 5 Isabela    body_weight    50
#> 6 Santa Cruz body_weight    50
```

## 2.3 Reshaping with `tidyr`

The `tidyr` package implements tools to easily switch between layouts and also perform a few other reshaping operations. Old school R users will be familiar with the `reshape` and `reshape2` packages, of which `tidyr` is the tidyverse equivalent. Beware that `tidyr` is about playing with the general *layout* of the dataset, while *operations* and *transformations* of the data are within the scope of the `dplyr` and `purrr` packages. All these packages work hand-in-hand really well, and analysis pipelines usually involve all of them. But today, we focus on the first member of this holy trinity, which is often the first one you'll need because you will want to reshape your data before doing other things. So, please hold your non-layout-related questions for the next chapters.

### 2.3.1 Pivoting

Pivoting a dataset between the long and wide layout is the main purpose of `tidyr` (check out the package's logo). We already saw the `pivot_longer` function above. This function converts a table from wide to long format. Similarly, there is a `pivot_wider` function that does exactly the opposite and takes you back to the wide format:

```
pivot_wider(
  data,
  names_from = "variable",
  values_from = "value",
  id_cols = c("id", "island")
)
#> # A tibble: 100 x 5
```

```
#>      id island  body_weight beak_length beak_width
#>   <int> <chr>      <dbl>      <dbl>      <dbl>
#> 1     1  Isabela      10.8        4.94        1.94
#> 2     2  Isabela      15.4        5.02        2.00
#> 3     3  Isabela      15.0        4.92        1.91
#> 4     4  Isabela       8.51        5.16        2.02
#> 5     5  Isabela      14.9        5.03        1.93
#> 6     6  Isabela       8.41        4.92        2.18
#> # ... with 94 more rows
```

331 The order of the columns is not exactly as it was, but this should not matter in  
 332 a data analysis pipeline where you should access columns by their names. It is  
 333 straightforward to change the order of the columns, but this is more within the  
 334 scope of the `dplyr` package.

335 If you are familiar with earlier versions of the tidyverse, `pivot_longer` and  
 336 `pivot_wider` are the respective equivalents of `gather` and `spread`, which are  
 337 now deprecated.

338 There are a few other reshaping operations from `tidyr` that are worth knowing.

### 339 2.3.2 Handling missing values

340 Say we have some missing measurements in the column “value” of our finch  
 341 dataset:

```
# We replace 100 random observations by NAs
ii <- sample(nrow(data), 100)
data$value[ii] <- NA
data
#> # A tibble: 300 x 4
#>      id island variable  value
#>   <int> <chr>   <chr>    <dbl>
#> 1     1  Isabela body_weight 10.8
#> 2     1  Isabela beak_length NA
#> 3     1  Isabela beak_width  NA
#> 4     2  Isabela body_weight NA
#> 5     2  Isabela beak_length 5.02
#> 6     2  Isabela beak_width  NA
#> # ... with 294 more rows
```

342 We could get rid of the rows that have missing values using `drop_na`:

```
drop_na(data, value)
#> # A tibble: 200 x 4
#>      id island variable  value
#>   <int> <chr>   <chr>    <dbl>
#> 1     1  Isabela body_weight 10.8
#> 2     2  Isabela beak_length 5.02
```

```
#> 3      3 Isabela body_weight 15.0
#> 4      3 Isabela beak_length 4.92
#> 5      4 Isabela body_weight 8.51
#> 6      4 Isabela beak_width 2.02
#> # ... with 194 more rows
```

343 Else, we could replace the NAs with some user-defined value:

```
replace_na(data, replace = list(value = -999))
#> # A tibble: 300 x 4
#>       id island variable      value
#>   <int> <chr>   <chr>    <dbl>
#> 1     1 Isabela body_weight 10.8
#> 2     1 Isabela beak_length -999
#> 3     1 Isabela beak_width -999
#> 4     2 Isabela body_weight -999
#> 5     2 Isabela beak_length 5.02
#> 6     2 Isabela beak_width -999
#> # ... with 294 more rows
```

344 where the **replace** argument takes a named list, and the names should refer to  
345 the columns to apply the replacement to.

346 We could also replace NAs with the most recent non-NA values:

```
fill(data, value)
#> # A tibble: 300 x 4
#>       id island variable      value
#>   <int> <chr>   <chr>    <dbl>
#> 1     1 Isabela body_weight 10.8
#> 2     1 Isabela beak_length 10.8
#> 3     1 Isabela beak_width 10.8
#> 4     2 Isabela body_weight 10.8
#> 5     2 Isabela beak_length 5.02
#> 6     2 Isabela beak_width 5.02
#> # ... with 294 more rows
```

347 Note that most functions in the tidyverse take a tibble as their first argument,  
348 and columns to which to apply the functions are usually passed as “objects”  
349 rather than character strings. In the above example, we passed the **value**  
350 column as **value**, not “value”. These column-objects are called by the tidyverse  
351 functions *in the context* of the data (the tibble) they belong to.

### 352 2.3.3 Splitting and combining cells

353 The **tidyr** package offers tools to split and combine columns. This is a nice  
354 extension to the string manipulations we saw last week in the **stringr** tutorial.

355 Say we want to add the specific dates when we took measurements on our birds

356 (we would normally do this using `dplyr` but for now we will stick to the old  
 357 way):

```
# Sample random dates for each observation
data$day <- sample(30, nrow(data), replace = TRUE)
data$month <- sample(12, nrow(data), replace = TRUE)
data$year <- sample(2019:2020, nrow(data), replace = TRUE)
data
#> # A tibble: 300 x 7
#>   id island variable    value  day month  year
#>   <int> <chr>   <chr>      <dbl> <int> <int> <int>
#> 1     1  Isabela body_weight 10.8     8     7  2020
#> 2     1  Isabela beak_length NA      19     7  2019
#> 3     1  Isabela beak_width  NA     17    12  2019
#> 4     2  Isabela body_weight  NA     20    12  2020
#> 5     2  Isabela beak_length 5.02    21    10  2020
#> 6     2  Isabela beak_width  NA     23     2  2020
#> # ... with 294 more rows
```

358 We could combine the `day`, `month` and `year` columns into a single `date` column,  
 359 with a dash as a separator, using `unite`:

```
data <- unite(data, day, month, year, col = "date", sep = "-")
data
#> # A tibble: 300 x 5
#>   id island variable    value date
#>   <int> <chr>   <chr>      <dbl> <chr>
#> 1     1  Isabela body_weight 10.8 8-7-2020
#> 2     1  Isabela beak_length NA   19-7-2019
#> 3     1  Isabela beak_width  NA   17-12-2019
#> 4     2  Isabela body_weight  NA   20-12-2020
#> 5     2  Isabela beak_length 5.02 21-10-2020
#> 6     2  Isabela beak_width  NA   23-2-2020
#> # ... with 294 more rows
```

360 Of course, we can revert back to the previous dataset by splitting the `date`  
 361 column with `separate`.

```
separate(data, date, into = c("day", "month", "year"))
#> # A tibble: 300 x 7
#>   id island variable    value day  month year
#>   <int> <chr>   <chr>      <dbl> <chr> <chr> <chr>
#> 1     1  Isabela body_weight 10.8   8     7    2020
#> 2     1  Isabela beak_length NA    19     7    2019
#> 3     1  Isabela beak_width  NA    17    12    2019
#> 4     2  Isabela body_weight  NA    20    12    2020
#> 5     2  Isabela beak_length 5.02   21    10    2020
#> 6     2  Isabela beak_width  NA    23     2    2020
```

```
#> # ... with 294 more rows
```

362 But note that the `day`, `month` and `year` columns are now of class `character` and  
 363 not `integer` anymore. This is because they result from the splitting of `date`,  
 364 which itself was a `character` column.

365 You can also separate a single column into multiple *rows* using `separate_rows`:

```
separate_rows(data, date)
#> # A tibble: 900 x 5
#>   id island variable value date
#>   <int> <chr>   <chr>     <dbl> <chr>
#> 1     1  Isabela body_weight  10.8 8
#> 2     1  Isabela body_weight  10.8 7
#> 3     1  Isabela body_weight  10.8 2020
#> 4     1  Isabela beak_length NA    19
#> 5     1  Isabela beak_length NA     7
#> 6     1  Isabela beak_length NA   2019
#> # ... with 894 more rows
```

### 366 2.3.4 Expanding tables using combinations

367 Instead of getting rid of rows with NAs, we may want to add rows with NAs,  
 368 for example, for combinations of parameters that we did not measure.

```
data <- separate(data, date, into = c("day", "month", "year"))
to_rm <- with(data, island == "Santa Cruz" & year == "2020")
data <- data[!to_rm,]
tail(data)
#> # A tibble: 6 x 7
#>   id island variable value day month year
#>   <int> <chr>   <chr>     <dbl> <chr> <chr> <chr>
#> 1    98 Santa Cruz beak_length  4.94 22   12   2019
#> 2    98 Santa Cruz beak_width   1.90 9     1   2019
#> 3    99 Santa Cruz body_weight  15.0 16    7   2019
#> 4    99 Santa Cruz beak_length NA    26   10   2019
#> 5    99 Santa Cruz beak_width   2.04 30    7   2019
#> 6   100 Santa Cruz beak_width   NA   23    3   2019
```

369 We could generate a tibble with all combinations of island, morphometric and  
 370 year using `expand_grid`:

```
expand_grid(
  island = c("Isabela", "Santa Cruz"),
  year = c("2019", "2020")
)
#> # A tibble: 4 x 2
#>   island year
#>   <chr>   <chr>
```



```
#> 1 Isabela 2019
#> 2 Isabela 2020
#> 3 Santa Cruz 2019
#> 4 Santa Cruz 2020
```

371 If we already have a tibble to work from that contains the variables to combine,  
372 we can use `expand` on that tibble:

```
expand(data, island, year)
#> # A tibble: 4 x 2
#>   island year
#>   <chr>   <chr>
#> 1 Isabela 2019
#> 2 Isabela 2020
#> 3 Santa Cruz 2019
#> 4 Santa Cruz 2020
```

373 As you can see, we get all the combinations of the variables of interest, even  
374 those that are missing. But sometimes you might be interested in variables  
375 that are *nested* within each other and not *crossed*. For example, say we have  
376 measured birds at different locations within each island:

```
nrow_Isabela <- with(data, length(which(island == "Isabela")))
nrow_SantaCruz <- with(data, length(which(island == "Santa Cruz")))
sites_Isabela <- sample(c("A", "B"), size = nrow_Isabela, replace = TRUE)
sites_SantaCruz <- sample(c("C", "D"), size = nrow_SantaCruz, replace = TRUE)
sites <- c(sites_Isabela, sites_SantaCruz)
data$site <- sites
data
#> # A tibble: 232 x 8
#>   id island variable value day month year site
#>   <int> <chr>   <chr>     <dbl> <chr> <chr> <chr> <chr>
#> 1     1 Isabela body_weight 10.8   8     7    2020 A
#> 2     1 Isabela beak_length NA     19    7    2019 B
#> 3     1 Isabela beak_width NA     17   12    2019 B
#> 4     2 Isabela body_weight NA     20   12    2020 A
#> 5     2 Isabela beak_length 5.02  21   10    2020 A
#> 6     2 Isabela beak_width NA     23    2    2020 A
#> # ... with 226 more rows
```

377 Of course, if sites A and B are on Isabela, they cannot be on Santa Cruz, where  
378 we have sites C and D instead. It would not make sense to `expand` assuming  
379 that `island` and `site` are crossed, instead, they are nested. We can therefore  
380 expand using the `nesting` function:

```
expand(data, nesting(island, site, year))
#> # A tibble: 6 x 3
#>   island site year
#>   <chr>   <chr> <chr>
```

```
#> 1 Isabela      A      2019
#> 2 Isabela      A      2020
#> 3 Isabela      B      2019
#> 4 Isabela      B      2020
#> 5 Santa Cruz C      2019
#> 6 Santa Cruz D      2019
```

381 But now the missing data for Santa Cruz in 2020 are not accounted for because  
 382 **expand** thinks the **year** is also nested within island. To get back the missing  
 383 combination, we use **crossing**, the complement of **nesting**:

```
expand(data, crossing(nesting(island, site), year)) # both can be used together
#> # A tibble: 8 x 3
#>   island      site year
#>   <chr>      <chr> <chr>
#> 1 Isabela      A      2019
#> 2 Isabela      A      2020
#> 3 Isabela      B      2019
#> 4 Isabela      B      2020
#> 5 Santa Cruz C      2019
#> 6 Santa Cruz C      2020
#> # ... with 2 more rows
```

384 Here, we specify that **site** is nested within **island** and these two are crossed  
 385 with **site**. Easy!

386 But wait a minute. These combinations are all very good, but our measurements  
 387 have disappeared! We can get them back by levelling up to the **complete**  
 388 function instead of using **expand**:

```
tail(complete(data, crossing(nesting(island, site), year)))
#> # A tibble: 6 x 8
#>   island      site year   id variable    value day  month
#>   <chr>      <chr> <chr> <int> <chr>      <dbl> <chr> <chr>
#> 1 Santa Cruz D      2019    95 beak_width  NA    13    10
#> 2 Santa Cruz D      2019    98 beak_length  4.94  22    12
#> 3 Santa Cruz D      2019    99 body_weight 15.0   16     7
#> 4 Santa Cruz D      2019    99 beak_length NA    26    10
#> 5 Santa Cruz D      2019    99 beak_width  2.04  30     7
#> 6 Santa Cruz D      2020    NA <NA>      NA    <NA> <NA>
# the last row has been added, full of NAs
```

389 which nicely keeps the rest of the columns in the tibble and just adds the missing  
 390 combinations.

### 391 2.3.5 Nesting

392 The **tidyr** package has yet another feature that makes the tidyverse very pow-  
 393 erful: the **nest** function. However, it makes little sense without combining it

394 with the functions in the `purrr` package, so we will not cover it in this chapter  
 395 but rather in the `purrr` chapter.

## 396 2.4 Extra: factors and the forcats package

```
library(forcats)
```

397 Categorical variables can be stored in R as character strings in `character` or  
 398 `factor` objects. A `factor` looks like a `character`, but it actually is an `integer`  
 399 vector, where each `integer` is mapped to a `character` label. With this respect  
 400 it is sort of an enhanced version of `character`. For example,

```
my_char_vec <- c("Pratik", "Theo", "Raph")
my_char_vec
#> [1] "Pratik" "Theo"   "Raph"
```

401 is a `character` vector, recognizable to its double quotes, while

```
my_fact_vec <- factor(my_char_vec) # as.factor would work too
my_fact_vec
#> [1] Pratik Theo   Raph
#> Levels: Pratik Raph Theo
```

402 is a `factor`, of which the *labels* are displayed. The *levels* of the factor are the  
 403 unique values that appear in the vector. If I added an extra occurrence of my  
 404 name:

```
factor(c(my_char_vec, "Raph"))
#> [1] Pratik Theo   Raph   Raph
#> Levels: Pratik Raph Theo
```

405 we would still have the the same levels. Note that the levels are returned as a  
 406 `character` vector in alphabetical order by the `levels` function:

```
levels(my_fact_vec)
#> [1] "Pratik" "Raph"   "Theo"
```

407 Why does it matter? Well, most operations on categorical variables can be  
 408 performed on `character` or `factor` objects, so it does not matter so much  
 409 which one you use for your own data. However, some functions in R require  
 410 you to provide categorical variables in one specific format, and others may even  
 411 implicitly convert your variables. In `ggplot2` for example, `character` vectors  
 412 are converted into factors by default. So, it is always good to remember the  
 413 differences and what type your variables are.

414 But this is a tidyverse tutorial, so I would like to introduce here the package  
 415 `forcats`, which offers tools to manipulate factors. First of all, most tools from  
 416 `stringr` *will work* on factors. The `forcats` functions expand the string manip-  
 417 ulation toolbox with factor-specific utilities. Similar in philosophy to `stringr`  
 418 where functions started with `str_`, in `forcats` most functions start with `fct_`.

419 I see two main ways `forcats` can come handy in the kind of data most people  
 420 deal with: playing with the order of the levels of a factor and playing with the  
 421 levels themselves. We will show here a few examples, but the full breadth of  
 422 factor manipulations can be found online or in the excellent `forcats` cheatsheet.

### 423 2.4.1 Change the order of the levels

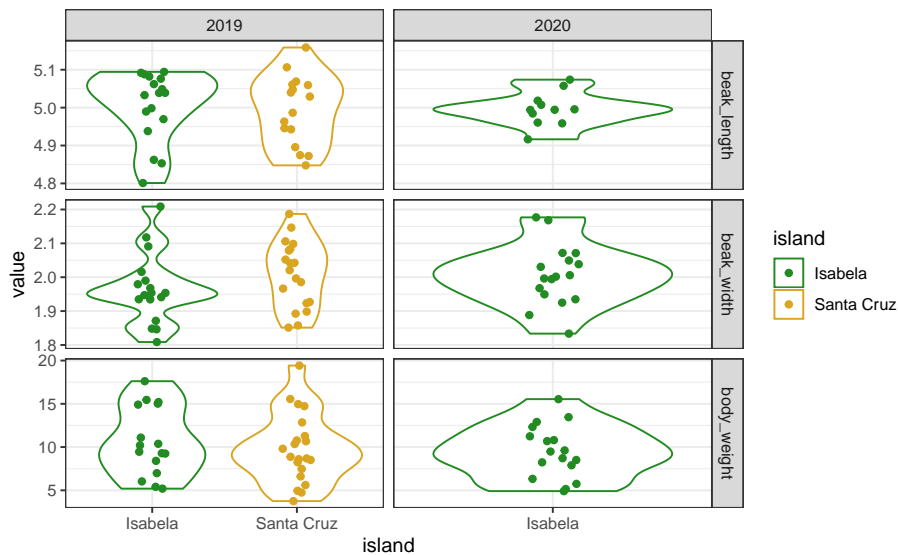
424 One example use-case where you would want to change the order of the levels  
 425 of a factor is when plotting. Your categorical variable, for example, may not be  
 426 plotted in the order you want. If we plot the distribution of each variable across  
 427 islands, we get

```
# Make the plotting code a function so we can re-use it without copying and pasting
my_plot <- function(data) {

  # We do not cover the ggplot functions in this chapter, this is just to
# illustrate our use-case, wait until chapter 5!
  library(ggplot2)
  ggplot(data, aes(x = island, y = value, color = island)) +
    geom_violin() +
    geom_jitter(width = 0.1) +
    facet_grid(variable ~ year, scales = "free") +
    theme_bw() +
    scale_color_manual(values = c("forestgreen", "goldenrod"))

}

my_plot(data)
# Remember that data are missing from Santa Cruz in 2020
```



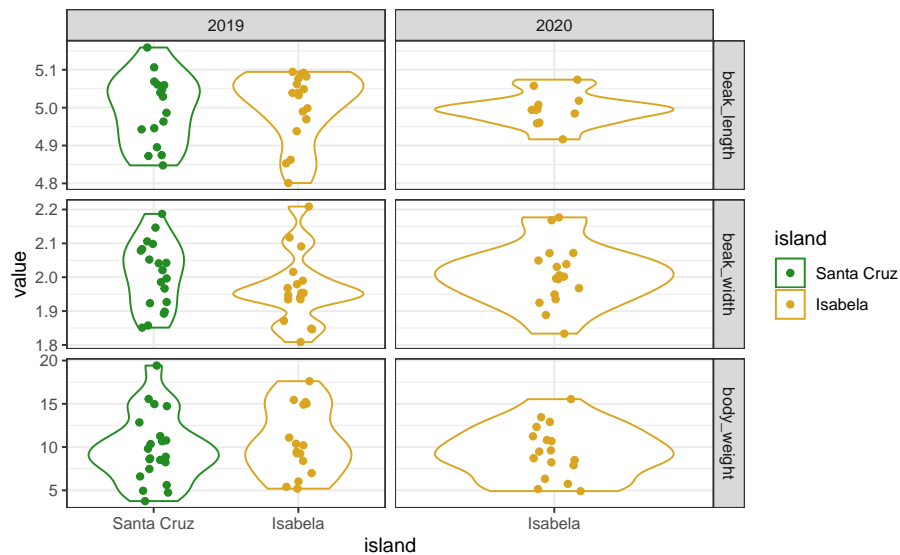
429

430 Here, the islands (horizontal axis) and the variables (the facets) are displayed  
 431 in alphabetical order. When making a figure you may want to customize these  
 432 orders in such a way that your message is optimally conveyed by your figure,  
 433 and this may involve playing with the order of levels.

434 Use `fct_relevel` to manually change the order of the levels:

```
data$island <- as.factor(data$island) # turn this column into a factor
data$island <- fct_relevel(data$island, c("Santa Cruz", "Isabela"))
my_plot(data) # order of islands has changed!
```

435



436

437 Beware that reordering a factor *does not change* the order of the items within  
 438 the vector, only the order of the *levels*. So, it does not introduce any mismatch  
 439 between the `island` column and the other columns! It only matters when the  
 440 levels are called, for example, in a `ggplot`. As you can see:

```
data$island[1:10]
#> [1] Isabela Isabela Isabela Isabela Isabela Isabela Isabela Isabela Isabela Isabela
#> [10] Isabela
#> Levels: Santa Cruz Isabela
fct_relevel(data$island, c("Isabela", "Santa Cruz"))[1:10] # same thing, different lev
#> [1] Isabela Isabela Isabela Isabela Isabela Isabela Isabela Isabela Isabela Isabela
#> [10] Isabela
#> Levels: Isabela Santa Cruz
```

441 Alternatively, use `fct_inorder` to set the order of the levels to the order in  
 442 which they appear:

```
data$variable <- as.factor(data$variable)
levels(data$variable)
#> [1] "beak_length" "beak_width" "body_weight"
levels(fct_inorder(data$variable))
#> [1] "body_weight" "beak_length" "beak_width"
```

443 or `fct_rev` to reverse the order of the levels:

```
levels(fct_rev(data$island)) # back in the alphabetical order
#> [1] "Isabela" "Santa Cruz"
```

444 Other variants exist to do more complex reordering, all present in the `forcats`  
 445 cheatsheet, for example: \* `fct_infreq` to re-order according to the frequency

of each level (how many observation on each island?) \* `fct_shift` to shift the order of all levels by a certain rank (in a circular way so that the last one becomes the first one or vice versa) \* `fct_shuffle` if you want your levels in random order \* `fct_reorder`, which reorders based on an associated variable (see `fct_reorder2` for even more complex relationship between the factor and the associated variable)

## 2.4.2 Change the levels themselves

Changing the levels of a factor will change the labels in the actual vector. It is similar to performing a string substitution in `stringr`. One can change the levels of a factor using `fct_recode`:

```
fct_recode(
  my_fact_vec,
  "Pratik Gupte" = "Pratik",
  "Theo Pannetier" = "Theo",
  "Raphael Scherrer" = "Raph"
)
#> [1] Pratik Gupte      Theo Pannetier  Raphael Scherrer
#> Levels: Pratik Gupte Raphael Scherrer Theo Pannetier
```

or collapse factor levels together using `fct_collapse`:

```
fct_collapse(my_fact_vec, EU = c("Theo", "Raph"), NonEU = "Pratik")
#> [1] NonEU EU      EU
#> Levels: NonEU EU
```

Again, we do not provide an exhaustive list of `forcats` functions here but the most usual ones, to give a glimpse of many things that one can do with factors. So, if you are dealing with factors, remember that `forcats` may have handy tools for you. Among others: \* `fct_anon` to “anonymize”, i.e. replace the levels by random integers \* `fct_lump` to collapse levels together based on their frequency (e.g. the two most frequent levels together)

## 2.4.3 Dropping levels

If you use factors in your tibble and get rid of one level, for any reason, the factor will usually remember the old levels, which may cause some problems when applying functions to your data.

```
data <- data[data$island == "Santa Cruz",] # keep only one island
unique(data$island) # Isabela is gone from the labels
#> [1] Santa Cruz
#> Levels: Santa Cruz Isabela
levels(data$island) # but not from the levels
#> [1] "Santa Cruz" "Isabela"
```

467 Use `droplevels` (from base R) to make sure you get rid of levels that are not  
 468 in your data anymore:

```
data <- droplevels(data)
levels(data$island)
#> [1] "Santa Cruz"
```

469 Fortunately, most functions within the tidyverse will not complain about missing  
 470 levels, and will automatically get rid of those inexistant levels for you. But  
 471 because factors are such common causes of bugs, keep this in mind!

472 Note that this is equivalent to doing:

```
data$island <- fct_drop(data$island)
```

#### 473 2.4.4 Other things

474 Among other things you can use in `forcats`: \* `fct_count` to get the frequency  
 475 of each level \* `fct_c` to combine factors together

#### 476 2.4.5 Take home message for forcats

477 Use this package to manipulate your factors. Do you need factors? Or are  
 478 character vectors enough? That is your call, and may depend on the kind of  
 479 analyses you want to do and what they require. We saw here that for plotting,  
 480 having factors can allow you to do quite some tweaking of the display. If you  
 481 encounter a situation where the order of encoding of your character vector starts  
 482 to matter, then maybe converting into a factor would make your life easier. And  
 483 if you do so, remember that lots of tools to perform all kinds of manipulation  
 484 are available to you with both `stringr` and `forcats`.

### 485 2.5 External resources

486 Find lots of additional info by looking up the following links:

- 487 • The `readr/tibble/tidyr` and `forcats` cheatsheets.
- 488 • This link on the concept of tidy data
- 489 • The `tibble`, `tidyr` and `forcats` websites



## Chapter 3

# Data manipulation with dplyr

```
# load the tidyverse
library(tidyverse)
```

### 3.1 Introduction

Reminders from last weeks: pipe operator, tidy tables, ggplot  
Why dplyr ? dplyr vs base R

### 3.2 Example data of the day

Through this tutorial, we will be using mammal trait data from the Phylacine database. The dataset contains information on mass, diet, life habit, etc, for more than all living species of mammals. Let's have a look.

```
phylacine <- readr::read_csv("data/phylacine_traits.csv")
phylacine
#> # A tibble: 5,831 x 24
#>   Binomial.1.2 Order.1.2 Family.1.2 Genus.1.2 Species.1.2 Terrestrial Marine
#>   <chr>           <chr>      <chr>      <chr>      <chr>           <dbl> <dbl>
#> 1 Abditomys_l~ Rodentia Muridae   Abditomys latidens         1      0
#> 2 Abeomelomys~ Rodentia Muridae   Abeomelo~ sevia           1      0
#> 3 Abrawayaomy~ Rodentia Cricetidae Abrawaya~ ruschii         1      0
#> 4 Abrocoma_be~ Rodentia Abrocomid~ Abrocoma bennettii        1      0
#> 5 Abrocoma_bo~ Rodentia Abrocomid~ Abrocoma boliviensis        1      0
#> 6 Abrocoma_bu~ Rodentia Abrocomid~ Abrocoma budini          1      0
#> # ... with 5,825 more rows, and 17 more variables: Freshwater <dbl>,
```

```
#> #   Aerial <dbl>, Life.Habit.Method <chr>, Life.Habit.Source <chr>,
#> #   Mass.g <dbl>, Mass.Method <chr>, Mass.Source <chr>, Mass.Comparison <chr>,
#> #   Mass.Comparison.Source <chr>, Island.Endemicity <chr>,
#> #   IUCN.Status.1.2 <chr>, Added.IUCN.Status.1.2 <chr>, Diet.Plant <dbl>,
#> #   Diet.Vertebrate <dbl>, Diet.Invertebrate <dbl>, Diet.Method <chr>,
#> #   Diet.Source <chr>
```

500 Note the friendly output given by the `tibble` (as opposed to a `data.frame`).  
 501 `readr` automatically stores the content it reads in a `tibble`, tidyverse oblige.  
 502 You should know however that `dplyr` doesn't require your data to be in a tibble,  
 503 a regular `data.frame` will work just as fine.

504 Most of the `dplyr` verbs covered in the next sections assume your data is *tidy*:  
 505 wide format, variables as column, 1 observation per row. Not that they won't  
 506 work if your data isn't tidy, but the results could be very different from what  
 507 I'm going to show here. Fortunately, the phylacine trait dataset appears to be  
 508 tidy: there is one unique entry for each species.

509 The first operation I'm going to run on this table is changing the names with  
 510 `rename()`. Some people prefer their tea without sugar, and I prefer my variable  
 511 names without uppercase characters, dots or (if possible) numbers. This will  
 512 give me the opportunity to introduce the trivial syntax of `dplyr` verbs.

```
phylacine <- phylacine %>%
  dplyr::rename(
    "binomial" = Binomial.1.2,
    "order" = Order.1.2,
    "family" = Family.1.2,
    "genus" = Genus.1.2,
    "species" = Species.1.2,
    "terrestrial" = Terrestrial,
    "marine" = Marine,
    "freshwater" = Freshwater,
    "aerial" = Aerial,
    "life_habit_method" = Life.Habit.Method,
    "life_habit_source" = Life.Habit.Source,
    "mass_g" = Mass.g,
    "mass_method" = Mass.Method,
    "mass_source" = Mass.Source,
    "mass_comparison" = Mass.Comparison,
    "mass_comparison_source" = Mass.Comparison.Source,
    "island_endemicity" = Island.Endemicity,
    "iucn_status" = IUCN.Status.1.2, # not even for acronyms
    "added_iucn_status" = Added.IUCN.Status.1.2,
    "diet_plant" = Diet.Plant,
    "diet_vertibrate" = Diet.Vertebrate,
    "diet_invertebrate" = Diet.Invertebrate,
    "diet_method" = Diet.Method,
```

```

    "diet_source" = Diet.Source
  )

```

513 For convenience, I'm going to use the pipe operator (`%>%`) that we've seen before,  
 514 through this chapter. All `dplyr` functions are built to work with the pipe (i.e.,  
 515 their first argument is always `data`), but again, this is not compulsory. I could  
 516 do

```

phylacine <- dplyr::rename(
  data = phylacine,
  "binomial" = Binomial.1.2,
  # ...
)

```

517 Note how columns are referred to. Once the data has been passed as an argument,  
 518 no need to refer to it anymore, `dplyr` understands that you're dealing with  
 519 variables inside that data frame. So drop that `data$var`, `data[, "var"]`, and,  
 520 if you've read *The R book*, forget the very existence of `attach()`.

521 Finally, I should mention that you can refer to variables names either with  
 522 strings or directly as objects, whether you're reading or creating them:

```

phylacine2 <- readr::read_csv("data/phylacine_traits.csv")

phylacine2 %>%
  dplyr::rename(
    # this works
    binomial = Binomial.1.2
  )
phylacine2 %>%
  dplyr::rename(
    # this works too!
    binomial = "Binomial.1.2"
  )
phylacine2 %>%
  dplyr::rename(
    # guess what
    "binomial" = "Binomial.1.2"
  )

```

### 523 3.3 Select variables with `select()`

### 524 3.4 Select observations with `filter()`

### 525 3.5 Create new variables with `mutate()`

526 can also edit existing ones

527 drop existing variables with `transmute()`

528 **3.6 Grouped results with `group_by()` and**  
529 **`summarise()`**

530 **3.7 Scoped variables**

```
data(mtcars)
mtcars %>% select_all(toupper)

is_whole <- function(x) all(floor(x) == x)
mtcars %>% select_if() # select integers only

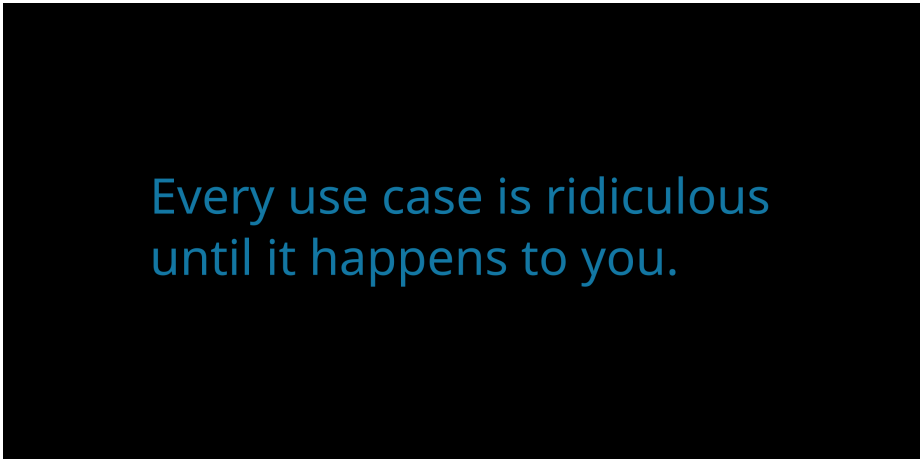
mtcars %>% select_at(vars(-contains("ar")))
mtcars %>% select_at(vars(-contains("ar"), starts_with("c")))
```

531 **3.8 More !**

532 `dollar` sign `x` point operator variables values -> `dplyr::distinct()` eq. to  
533 `base::unique()` `sample()` `slice()`

## 534 Chapter 4

# 535 Working with lists and 536 iteration



Every use case is ridiculous  
until it happens to you.

537

```
# load the tidyverse  
library(tidyverse)
```

## 538 4.1 List columns with `tidyr`

### 539 4.1.1 Nesting data

540 It may become necessary to indicate the groups of a tibble in a somewhat more  
541 explicit way than simply using `dplyr::group_by`. `tidyr` offers the option to  
542 create nested tibbles, that is, to store complex objects in the columns of a tibble.  
543 This includes other tibbles, as well as model objects and plots.

544 *NB:* Nesting data is done using `tidyr::nest`, which is different from the simi-  
 545 larly named `tidyr::nesting`.

546 The example below shows how `mtcars` can be converted into a nested tibble.

```
# nest mtcars into a list of dataframes based on number of cylinders
nested_cars = as_tibble(mtcars,
                        rownames = "car_name") %>%
  group_by(cyl) %>%
  nest()

nested_cars
#> # A tibble: 3 x 2
#> # Groups:   cyl [3]
#>   cyl data
#>   <dbl> <list>
#> 1     6 <tibble [7 x 11]>
#> 2     4 <tibble [11 x 11]>
#> 3     8 <tibble [14 x 11]>

# get column class
sapply(nested_cars, class)
#>      cyl      data
#> "numeric"  "list"
```

547 `mtcars` is now a nested data frame. The class of each of its columns is respec-  
 548 tively, a numeric (number of cylinders) and a list (the data of all cars with as  
 549 many cylinders as in the corresponding row).

550 While `nest` can be used without first grouping the tibble, it's just much easier  
 551 to group first.

### 552 4.1.2 Unnesting data

553 A nested tibble can be converted back into the original, or into a processed form,  
 554 using `tidyr::unnest`. The original groups are retained.

```
# use unnest to recover the original data frame
unnest(nested_cars, cols = "data")
#> # A tibble: 32 x 12
#> # Groups:   cyl [3]
#>   cyl car_name      mpg  disp    hp  drat    wt  qsec    vs    am  gear  carb
#>   <dbl> <chr>      <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
#> 1     6 Mazda RX4      21   160   110  3.9   2.62  16.5    0     1     4     4
#> 2     6 Mazda RX4 W~  21   160   110  3.9   2.88  17.0    0     1     4     4
#> 3     6 Hornet 4 Dr~  21.4 258   110  3.08  3.22  19.4    1     0     3     1
#> 4     6 Valiant      18.1 225   105  2.76  3.46  20.2    1     0     3     1
#> 5     6 Merc 280     19.2 168.   123  3.92  3.44  18.3    1     0     4     4
```

```
#> 6      6 Merc 280C      17.8 168.   123 3.92 3.44 18.9      1      0      4      4
#> # ... with 26 more rows
```

```
# unnesting preserves groups
groups(unnest(nested_cars, cols = "data"))
#> [[1]]
#> cyl
```

555 The `unnest_longer` and `unnest_wider` variants of `unnest` are maturing func-  
 556 tions, that is, not in their final form. They allow interesting variations on  
 557 unnesting — these are shown here but advised against.

558 Unnest the data first, and then convert it to the form needed.

```
unnest_longer(nested_cars, col = "data") %>%
  head()
#> # A tibble: 6 x 2
#> # Groups:   cyl [1]
#>   cyl data$car_name    $mpg $disp   $hp $drat   $wt $qsec   $vs   $am $gear
#>   <dbl> <chr>          <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
#> 1     6 Mazda RX4      21    160   110 3.9    2.62 16.5     0     1     4
#> 2     6 Mazda RX4 Wag  21    160   110 3.9    2.88 17.0     0     1     4
#> 3     6 Hornet 4 Dri~ 21.4  258   110 3.08   3.22 19.4     1     0     3
#> 4     6 Valiant       18.1  225   105 2.76   3.46 20.2     1     0     3
#> 5     6 Merc 280      19.2  168.   123 3.92   3.44 18.3     1     0     4
#> 6     6 Merc 280C     17.8  168.   123 3.92   3.44 18.9     1     0     4
#> # ... with 1 more variable: $carb <dbl>

unnest_wider(nested_cars, col = "data")
#> # A tibble: 3 x 12
#> # Groups:   cyl [3]
#>   cyl car_name mpg    disp  hp    drat  wt    qsec  vs    am    gear  carb
#>   <dbl> <list>    <list> <list> <list> <lis> <lis> <lis> <lis> <lis> <lis> <lis>
#> 1     6 <chr [7]> <dbl ~ <dbl ~ <dbl ~ <dbl~ <dbl~ <dbl~ <dbl~ <dbl~ <dbl~ <dbl~ <dbl~
#> 2     4 <chr [11~ <dbl ~ <dbl ~ <dbl ~ <dbl~ <dbl~ <dbl~ <dbl~ <dbl~ <dbl~ <dbl~
#> 3     8 <chr [14~ <dbl ~ <dbl ~ <dbl ~ <dbl~ <dbl~ <dbl~ <dbl~ <dbl~ <dbl~ <dbl~
```

### 559 4.1.3 Working with list columns

560 The class of a list column is `list`, and working with list columns (and lists, and  
 561 list-like objects such as vectors) makes iteration necessary, since this is one of  
 562 the only ways to operate on lists.

563 Two examples are shown below when getting the class and number of rows of  
 564 the nested tibbles in the list column.

```
# how many rows in each nested tibble?
for (i in seq_along(nested_cars$data)) {
```

```

    print(nrow(nested_cars$data[[i]]))
  }
#> [1] 7
#> [1] 11
#> [1] 14

# what is the class of each element?
lapply(X = nested_cars$data, FUN = class)
#> [[1]]
#> [1] "tbl_df"      "tbl"        "data.frame"
#>
#> [[2]]
#> [1] "tbl_df"      "tbl"        "data.frame"
#>
#> [[3]]
#> [1] "tbl_df"      "tbl"        "data.frame"

```

## 565 Functionals

566 The second example uses `lapply`, and this is a *functional*. *Functionals* are func-  
 567 tions that take another function as one of their arguments. Base R functionals  
 568 include the `*apply` family of functions: `apply`, `lapply`, `vapply` and so on.

## 569 4.2 Iteration with map

570 The `tidyverse` replaces traditional loop-based iteration with *functionals* from  
 571 the `purrr` package. A good reason to use `purrr` functionals instead of base  
 572 R functionals is their consistent and clear naming, which always indicates how  
 573 they should be used. This is explained in the examples below.

574 How `map` is different from `for` and `lapply` are best explained in the **Advanced**  
 575 **R Book**.

### 576 4.2.1 Basic use of map

577 `map` works very similarly to `lapply`, where `.x` is object on whose elements to  
 578 apply the function `.f`.

```

# get the number of rows in data
map(.x = nested_cars$data, .f = nrow)
#> [[1]]
#> [1] 7
#>
#> [[2]]
#> [1] 11
#>

```



```
#> [[3]]
#> [1] 14
```

579 `map` works on any list-like object, which includes vectors, and always returns a  
 580 list. `map` takes two arguments, the object on which to operate, and the function  
 581 to apply to each element.

```
# get the square root of each integer 1 - 10
some_numbers = 1:3
map(some_numbers, sqrt)
#> [[1]]
#> [1] 1
#>
#> [[2]]
#> [1] 1.41
#>
#> [[3]]
#> [1] 1.73
```

## 582 4.2.2 map variants returning vectors

583 Though `map` always returns a list, it has variants named `map_*` where the suffix  
 584 indicates the return type. `map_chr`, `map_dbl`, `map_int`, and `map_lgl` return  
 585 character, double (numeric), integer, and logical vectors.

```
# use map_dbl to get a vector of square roots
some_numbers = 1:10
map_dbl(some_numbers, sqrt)
#> [1] 1.00 1.41 1.73 2.00 2.24 2.45 2.65 2.83 3.00 3.16

# map_chr will convert the output to a character
map_chr(some_numbers, sqrt)
#> [1] "1.000000" "1.414214" "1.732051" "2.000000" "2.236068" "2.449490"
#> [7] "2.645751" "2.828427" "3.000000" "3.162278"

# map_lgl returns TRUE/FALSE values
some_numbers = c(NA, 1:3, NA, NaN, Inf, -Inf)
map_lgl(some_numbers, is.na)
#> [1] TRUE FALSE FALSE FALSE TRUE TRUE FALSE FALSE
```

## 586 4.2.3 map variants returning data frames

587 `map_df` returns data frames, and by default binds dataframes by rows, while  
 588 `map_dfr` does this explicitly, and `map_dfc` does returns a dataframe bound by  
 589 column.

```
# split mtcars into 3 dataframes, one per cylinder number
some_list = split(mtcars, mtcars$cyl)
```

```
# get the first two rows of each dataframe
map_df(some_list, head, n = 2)
#>      mpg cyl disp  hp drat   wt  qsec vs am gear carb
#> Datsun 710  22.8  4  108  93 3.85 2.32 18.6 1  1   4   1
#> Merc 240D  24.4  4  147  62 3.69 3.19 20.0 1  0   4   2
#> Mazda RX4  21.0  6  160 110 3.90 2.62 16.5 0  1   4   4
#> Mazda RX4 Wag 21.0  6  160 110 3.90 2.88 17.0 0  1   4   4
#> Hornet Sportabout 18.7  8  360 175 3.15 3.44 17.0 0  0   3   2
#> Duster 360  14.3  8  360 245 3.21 3.57 15.8 0  0   3   4
```

590 map accepts arguments to the function being mapped, such as in the example  
 591 above, where head() accepts the argument n = 2.

592 map\_dfr behaves the same as map\_df.

```
# the same as above but with a pipe
some_list %>%
  map_dfr(head, n = 2)
#>      mpg cyl disp  hp drat   wt  qsec vs am gear carb
#> Datsun 710  22.8  4  108  93 3.85 2.32 18.6 1  1   4   1
#> Merc 240D  24.4  4  147  62 3.69 3.19 20.0 1  0   4   2
#> Mazda RX4  21.0  6  160 110 3.90 2.62 16.5 0  1   4   4
#> Mazda RX4 Wag 21.0  6  160 110 3.90 2.88 17.0 0  1   4   4
#> Hornet Sportabout 18.7  8  360 175 3.15 3.44 17.0 0  0   3   2
#> Duster 360  14.3  8  360 245 3.21 3.57 15.8 0  0   3   4
```

593 map\_dfc binds the resulting 3 data frames of two rows each by column, and  
 594 automatically repairs the column names, adding a suffix to each duplicate.

```
some_list %>%
  map_dfc(head, n = 2)
#>      mpg...1 cyl...2 disp...3 hp...4 drat...5 wt...6 qsec...7 vs...8
#> Datsun 710  22.8      4      108      93      3.85      2.32      18.6      1
#> Merc 240D  24.4      4      147      62      3.69      3.19      20.0      1
#>      am...9 gear...10 carb...11 mpg...12 cyl...13 disp...14 hp...15
#> Datsun 710      1      4      1      21      6      160      110
#> Merc 240D      0      4      2      21      6      160      110
#>      drat...16 wt...17 qsec...18 vs...19 am...20 gear...21 carb...22
#> Datsun 710      3.9      2.62      16.5      0      1      4      4
#> Merc 240D      3.9      2.88      17.0      0      1      4      4
#>      mpg...23 cyl...24 disp...25 hp...26 drat...27 wt...28 qsec...29
#> Datsun 710      18.7      8      360      175      3.15      3.44      17.0
#> Merc 240D      14.3      8      360      245      3.21      3.57      15.8
#>      vs...30 am...31 gear...32 carb...33
#> Datsun 710      0      0      3      2
#> Merc 240D      0      0      3      4
```

### 595 4.2.4 Working with list columns using map

596 The various `map` versions integrate well with list columns to make syn-  
 597 thetic/summary data. In the example, the `dplyr::mutate` function is used to  
 598 add three columns to the nested tibble: the number of rows, the mean mileage,  
 599 and the name of the first car.

600 In each of these cases, the vectors added are generated using `purrr` functions.

```
# get the number of rows per dataframe, the mean mileage, and the first car
nested_cars = nested_cars %>%
  mutate(
    # use the int return to get the number of rows
    n_rows = map_int(data, nrow),

    # double return for mean mileage
    mean_mpg = map_dbl(data, function(df) {mean(df$mpg)}),

    # character return to get first car
    first_car = map_chr(data, function(df) {first(df$car_name)})
  )

# examine the output
nested_cars
#> # A tibble: 3 x 5
#> # Groups:   cyl [3]
#>   cyl data                n_rows mean_mpg first_car
#>   <dbl> <list>                <int>    <dbl> <chr>
#> 1     6 <tibble [7 x 11]>         7      19.7 Mazda RX4
#> 2     4 <tibble [11 x 11]>      11      26.7 Datsun 710
#> 3     8 <tibble [14 x 11]>     14      15.1 Hornet Sportabout
```

### 601 4.2.5 Selective mapping using map variants

602 `map_at` and `map_if` work like other `*_at` and `*_if` functions. Here, `map_if` is  
 603 used to run a linear model only on those tibbles which have sufficient data. The  
 604 predicate is specified by `.p`.

605 In this example, the nested tibble is given a new column using `dplyr::mutate`,  
 606 where the data to be added is a mixed list.

```
# split mtcars by cylinder number and run an lm only if there are more than 10 rows
data = nest(mtcars, data = -cyl)

data = mutate(data,
  model = map_if(.x = data,
    .p = function(x){
```

```

nrow(x) > 10
},
.f = function(x){
  lm(mpg ~ wt, data = x)
}))

# check the data structure
data
#> # A tibble: 3 x 3
#>   cyl data      model
#>   <dbl> <list>    <list>
#> 1     6 <tibble [7 x 10]> <tibble [7 x 10]>
#> 2     4 <tibble [11 x 10]> <lm>
#> 3     8 <tibble [14 x 10]> <lm>

```

607 The first element is a tibble of the corresponding element in `mtcars$cars`, which  
 608 has not been operated on because it has fewer than 10 rows. The remaining  
 609 elements are `lm` objects.

### 610 4.3 More map variants

611 `map` also has variants along the axis of how many elements are operated upon.  
 612 `map2` operates on two vectors or list-like elements, and returns a single list as  
 613 output, while `pmap` operates on a list of list-like elements. The output has as  
 614 many elements as the input lists, which must be of the same length.

#### 615 4.3.1 Mapping over two inputs with `map2`

616 `map2` has the same variants as `map`, allowing for different return types. Here  
 617 `map2_int` returns an integer vector.

```

# consider 2 vectors and replicate the simple vector addition using map2
map2_int(.x = 1:5,
         .y = 6:10,
         .f = sum)
#> [1]  7  9 11 13 15

```

618 `map2` doesn't have `_at` and `_if` variants.

619 One use case for `map2` is to deal with both a list element and its index, as shown  
 620 in the example. This may be necessary when the list index is removed in a  
 621 `split` or `nest`. This can also be done with `imap`, where the index is referred to  
 622 as `.y`.

```

# make a named list for this example
this_list = list(a = "first letter",
                 b = "second letter")

# a not particularly useful example

```

```

map2(this_list, names(this_list),
      function(x, y) {
        glue::glue('{x} : {y}')
      })
#> $a
#> first letter : a
#>
#> $b
#> second letter : b

# imap can also do this
imap(this_list,
      function(x, .y){
        glue::glue('{x} : {.y}')}
      })
#> $a
#> first letter : a
#>
#> $b
#> second letter : b

```

### 623 4.3.2 Mapping over multiple inputs with pmap

624 pmap instead operates on a list of multiple list-like objects, and also comes with  
 625 the same return type variants as map. The example shows both aspects of pmap  
 626 using pmap\_chr.

```

# operate on three different lists
list_01 = as.list(1:3)
list_02 = as.list(letters[1:3])
list_03 = as.list(rainbow(3))

# print a few statements
pmap_chr(list(list_01, list_02, list_03),
          function(l1, l2, l3){
            glue::glue('number {l1}, letter {l2}, colour {l3}')}
          })
#> [1] "number 1, letter a, colour #FF0000FF"
#> [2] "number 2, letter b, colour #00FF00FF"
#> [3] "number 3, letter c, colour #0000FFFF"

```

### 627 4.3.3 Mapping at depth

628 Lists are often nested, that is, a list element may itself be a list. It is possible  
 629 to map a function over elements as a specific depth.

630 In the example, mtcars is split by cylinders, and then by gears, creating a

631 two-level list, with the second layer operated on.

```
# use map to make a 2 level list
this_list = split(mtcars, mtcars$cyl) %>%
  map(function(df){ split(df, df$gear) })

# map over the second level to count the number of
# cars with N gears in the set of cars with M cylinders
# display only for cyl = 4
map_depth(this_list[1], 2, nrow)
#> $`4`
#> $`4`$`3`
#> [1] 1
#>
#> $`4`$`4`
#> [1] 8
#>
#> $`4`$`5`
#> [1] 2
```

#### 632 4.3.4 Iteration without a return

633 `map` and its variants have a return type, which is either a list or a vector. How-  
 634 ever, it is often necessary to iterate a function over a list-like object for that  
 635 function's side effects, such as printing a message to screen, plotting a series of  
 636 figures, or saving to file.

637 `walk` is the function for this task. It has only the variants `walk2`, `iwalk`, and  
 638 `pwalk`, whose logic is similar to `map2`, `imap`, and `pmap`. In the example, the  
 639 function applied to each list element is intended to print a message.

```
this_list = split(mtcars, mtcars$cyl)

iwalk(this_list,
  function(df, .y){
    message(glue::glue('{nrow(df)} cars with {.y} cylinders'))
  })
```

#### 640 4.3.5 Modify rather than map

641 When the return type is expected to be the same as the input type, that is, a  
 642 list returning a list, or a character vector returning the same, `modify` can help  
 643 with keeping strictly to those expectations.

644 In the example, simply adding 2 to each vector element produces an error,  
 645 because the output is a `numeric`, or `double`. `modify` helps ensure some type  
 646 safety in this way.

```

vec = as.integer(1:10)

tryCatch(
  expr = {

    # this is what we want you to look at

    modify(vec, function(x) { (x + 2) })

  },

  # do not pay attention to this
  error = function(e){
    print(toString(e))
  }
)
#> [1] "Error: Can't coerce element 1 from a double to a integer\n"

```

647 Converting the output to an integer, which was the original input type, serves  
 648 as a solution.

```

modify(vec, function(x) { as.integer(x + 2) })
#> [1] 3 4 5 6 7 8 9 10 11 12

```

#### 649 A note on invoke

650 `invoke` used to be a wrapper around `do.call`, and can still be found with its  
 651 family of functions in `purrr`. It is however retired in favour of functionality  
 652 already present in `map` and `rlang::exec`, the latter of which will be covered in  
 653 another session.

## 654 4.4 Other functions for working with lists

655 `purrr` has a number of functions to work with lists, especially lists that are not  
 656 nested list-columns in a tibble.

### 657 4.4.1 Filtering lists

658 Lists can be filtered on any predicate using `keep`, while the special case `compact`  
 659 is applied when the empty elements of a list are to be filtered out. `discard` is  
 660 the opposite of `keep`, and keeps only elements not satisfying a condition. Again,  
 661 the predicate is specified by `.p`.

```

# a list containing numbers
this_list = list(a = 1, b = -1, c = 2, d = NULL, e = NA)

# remove the empty element

```

```

# this must be done before using keep on the list
this_list = compact(this_list)

# use discard to remove the NA
this_list = discard(this_list, .p =is.na)

# keep list elements which are positive
keep(this_list, .p = function(x){ x > 0 })
#> $a
#> [1] 1
#>
#> $c
#> [1] 2

```

662 `head_while` is bit of an odd case, which returns all elements of a list-like object  
 663 in sequence until the first one fails to satisfy a predicate, specified by `.p`.

```

1:10 %>%
  head_while(.p = function(x) x < 5)
#> [1] 1 2 3 4

```

## 664 4.4.2 Summarising lists

665 The `purrr` functions `every`, `some`, `has_element`, `detect`, `detect_index`, and  
 666 `vec_depth` help determine whether a list passes a certain logical test or not.  
 667 These are seldom used and are not discussed here.

## 668 4.4.3 Reduction and accumulation

669 `reduce` helps combine elements along a list using a specific function. Consider  
 670 the example below where list elements are concatenated into a single vector.

```

this_list = list(a = 1:3, b = 3:4, c = 5:10)

reduce(this_list, c)
#> [1] 1 2 3 3 4 5 6 7 8 9 10

```

671 This can also be applied to data frames. Consider some random samples of  
 672 `mtcars`, each with only 5 cars removed. The objective is to find the cars present  
 673 in all 10 samples.

674 The way `reduce` works in the example below is to take the first element and find  
 675 its intersection with the second, and to take the result and find its intersection  
 676 with the third and so on.

```

# sample mtcars
mtcars = as_tibble(mtcars, rownames = "car")
sampled_data = map(1:10, function(x){sample_n(mtcars, nrow(mtcars)-5)})

```



```
# get cars which appear in all samples
sampled_data = reduce(sampled_data, dplyr::inner_join)
```

677 **accumulate** works very similarly, except it retains the intermediate products.  
 678 The first element is retained as is. **accumulate2** and **reduce2** work on two lists,  
 679 following the same logic as **map2** etc. Both functions can be used in much more  
 680 complex ways than demonstrated here.

```
# make a list
this_list = list(a = 1:3, b = 3:6, c = 5:10, d = c(1,2,5,10,12))

# a multiple accumulate can help
accumulate(this_list, union, .dir = "forward")
#> $a
#> [1] 1 2 3
#>
#> $b
#> [1] 1 2 3 4 5 6
#>
#> $c
#> [1] 1 2 3 4 5 6 7 8 9 10
#>
#> $d
#> [1] 1 2 3 4 5 6 7 8 9 10 12
```

#### 681 4.4.4 Miscellaneous operation

682 **purrr** offers a few more functions to work with lists (or list like objects).  
 683 **prepend** works very similarly to **append**, except it adds to the head of a list.  
 684 **splice** adds multiple objects together in a list. **splice** will break the existing  
 685 list structure of input lists.

```
# use prepend to add values to the head of a list
prepend(x = list("a", "b"), values = list("1", "2"))
#> [[1]]
#> [1] "1"
#>
#> [[2]]
#> [1] "2"
#>
#> [[3]]
#> [1] "a"
#>
#> [[4]]
#> [1] "b"

# use splice to add multiple elements together
```

```
splice(list("a", "b"), list("1", "2"), "something else")
#> [[1]]
#> [1] "a"
#>
#> [[2]]
#> [1] "b"
#>
#> [[3]]
#> [1] "1"
#>
#> [[4]]
#> [1] "2"
#>
#> [[5]]
#> [1] "something else"
```

686 **flatten** has a similar behaviour, and converts a list of vectors or list of lists to a  
 687 single list-like object. **flatten\_\*** options allow the output type to be specified.

```
this_list = list(a = rep("a", 3),
                  b = rep("b", 4))
```

```
this_list
#> $a
#> [1] "a" "a" "a"
#>
#> $b
#> [1] "b" "b" "b" "b"

# use flatten chr to get a character vector
flatten_chr(this_list)
#> [1] "a" "a" "a" "b" "b" "b" "b"
```

688 **transpose** shifts the index order in multi-level lists. This is seen in the example,  
 689 where the **gear** goes from being the index of the second level to the index of the  
 690 first.

```
this_list = split(mtcars, mtcars$cyl) %>%
  map(function(df) split(df, df$gear))

# from a list of lists where cars are divided by cylinders and then
# gears, this is now a list of lists where cars are divided by
# gears and then cylinders
transpose(this_list[1])
#> $`3`
#> $`3`$`4`
#> # A tibble: 1 x 12
#>   car      mpg   cyl  disp    hp  drat    wt   qsec    vs    am gear carb
```

```

#>   <chr>           <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
#> 1 Toyota Coro~  21.5    4  120.    97   3.7  2.46  20.0    1    0    3    1
#>
#> $`4`
#> $`4`$`4`
#> # A tibble: 8 x 12
#>   car      mpg   cyl  disp    hp  drat    wt  qsec    vs  am  gear  carb
#>   <chr>    <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
#> 1 Datsun 710  22.8    4  108    93  3.85  2.32  18.6    1    1    4    1
#> 2 Merc 240D  24.4    4  147    62  3.69  3.19  20     1    0    4    2
#> 3 Merc 230   22.8    4  141    95  3.92  3.15  22.9    1    0    4    2
#> 4 Fiat 128   32.4    4   78.7   66  4.08  2.2   19.5    1    1    4    1
#> 5 Honda Civic 30.4    4   75.7   52  4.93  1.62  18.5    1    1    4    2
#> 6 Toyota Coro~ 33.9    4   71.1   65  4.22  1.84  19.9    1    1    4    1
#> # ... with 2 more rows
#>
#>
#> $`5`
#> $`5`$`4`
#> # A tibble: 2 x 12
#>   car      mpg   cyl  disp    hp  drat    wt  qsec    vs  am  gear  carb
#>   <chr>    <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
#> 1 Porsche 914~  26     4  120.    91  4.43  2.14  16.7    0    1    5    2
#> 2 Lotus Europa 30.4    4   95.1  113  3.77  1.51  16.9    1    1    5    2

```

## 691 4.5 To add: patchwork

### 692 4.5.0.1 Final words

693 In general, an iteration based problem can usually be solved with `purrr`.