

# TRES Tidyverse Tutorial

Raphael, Pratik and Theo

2020-06-01



# 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> . . . . .	19
11	1.4 Strings in <code>ggplot</code> . . . . .	20
12	<b>2 Reshaping data tables in the tidyverse</b>	<b>23</b>
13	2.1 1. The new data frame: <code>tibble</code> . . . . .	24
14	2.2 2. The concept of tidy data . . . . .	26
15	2.3 3. Reshaping with <code>tidyr</code> . . . . .	29
16	2.4 4. Extra: factors and the <code>forcats</code> package . . . . .	35
17	2.5 5. External resources . . . . .	39
18	<b>3 Working with lists and iteration</b>	<b>41</b>
19	3.1 Basic iteration with <code>map</code> . . . . .	41
20	3.2 More <code>map</code> variants . . . . .	45
21	3.3 Modification in place . . . . .	45
22	3.4 Working with lists . . . . .	46
23	<b>4 Data manipulation with <code>dplyr</code></b>	<b>47</b>
24	4.1 Introduction . . . . .	47
25	4.2 Example data of the day . . . . .	47
26	4.3 Select variables with <code>select()</code> . . . . .	50
27	4.4 Select observations with <code>filter()</code> . . . . .	53
28	4.5 Create new variables with <code>mutate()</code> . . . . .	57
29	4.6 Grouped results with <code>group_by()</code> and <code>summarise()</code> . . . . .	58
30	4.7 Scoped variables . . . . .	58
31	4.8 More ! . . . . .	58



# Outline

This is the readable version of the TRES tidyverse tutorial.

## 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

## 44 Possible extras

- 45 • Reproducibility and package-making (with e.g. `usethis`)
- 46 • Embedding C++ code with `Rcpp`

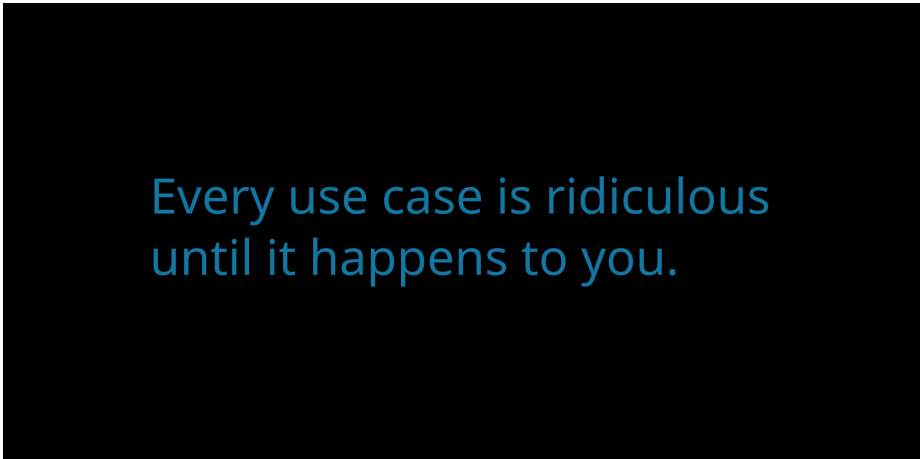
## 48 Join

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

50 \*Tentative dates.

## 51 Chapter 1

# 52 Reading files and string 53 manipulation



Every use case is ridiculous  
until it happens to you.

54

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

## 55 1.1 Data import and export with readr

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

## 60 1.1.1 Reading data

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

```

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

# read the file in
some_example = read_csv(some_example)

## Parsed with column specification:
## cols(
##   mpg = col_double(),
##   cyl = col_double(),
##   disp = col_double(),
##   hp = col_double(),
##   drat = col_double(),
##   wt = col_double(),
##   qsec = col_double(),
##   vs = col_double(),
##   am = col_double(),
##   gear = col_double(),
##   carb = col_double()
## )

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

```

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

90 Other variants include:

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



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

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

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

## 102 1.1.2 Writing data

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

107 These functions are not covered here.

## 108 1.1.3 Reading and writing lines

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

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

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

116 The model summary can be written to file. When writing lines to file, BE  
 117 AWARE OF THE DIFFERENCES BETWEEN UNIX AND WINDOWS line  
 118 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")
```

119 This model output can be read back in for display, and each line of the model  
 120 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
## -4.5432 -2.3647 -0.1252  1.4096  6.8727
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  37.2851     1.8776   19.858 < 2e-16 ***
## wt          -5.3445     0.5591   -9.559 1.29e-10 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.046 on 30 degrees of freedom
## Multiple R-squared:  0.7528, Adjusted R-squared:  0.7446
## F-statistic: 91.38 on 1 and 30 DF,  p-value: 1.294e-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.

154 `stringr` functions begin with `str_`.

### 155 1.2.1 Putting strings together

156 Concatenate two strings with `str_c`, and duplicate strings with `str_dup`. Flat-  
157 ten a list or vector of strings using `str_flatten`.

```

158 # str_c works like paste(), choose a separator
159 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 = "_")

160 ## [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 = "-")

161 ## [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"

162 str_flatten is especially useful when displaying the type of an object that
163 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 = ", ")

164 ## [1] "tbl_df, tbl, data.frame"
```

### 165 1.2.2 Detecting strings

166 Count the frequency of a pattern in a string with `str_count`. Returns an integer.  
167 Detect whether a pattern exists in a string with `str_detect`. Returns a logical  
168 and can be used as a predicate.

169 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")

170 ## [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")
```

```

171 ## [1] 5 3

      # vectorise over the pattern to count both a-s and b-s
      str_count(string = "ababababa", pattern = c("a", "b"))

172 ## [1] 5 4

173 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"))

174 ## [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"))

175 ## [1] 5 1 4 3

176 str_locate locates the search pattern in a string, and returns the start and
177 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")

178 ##          start end
179 ## [1,]         2   4

      # str_detect detects a sequence in a string
      str_detect(string = "Bananageddon is coming!",
                pattern = "na")

180 ## [1] TRUE

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

181 ## [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")))

182 ## [1] TRUE

```

183 Detect whether a string starts or ends with a pattern. Also vectorised. Both  
 184 have a `negate` argument, which returns the negative, i.e., returns `FALSE` if the  
 185 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")

186 ## [1] FALSE FALSE  TRUE  TRUE

# str_ends looks at the last character
str_ends(fruit, "e")

187 ## [1]  TRUE FALSE FALSE  TRUE

# an example of negate = TRUE
str_ends(fruit, "e", negate = TRUE)

188 ## [1] FALSE  TRUE  TRUE FALSE

189 str_subset [WHICH IS NOT RELATED TO str_sub] helps with subsetting a
190 character vector based on a str_detect predicate. In the example, all elements
191 containing “banana” are subset.

192 str_which has the same logic except that it returns the vector position and not
193 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")

194 ## [1] "banana"                "bananageddon is coming"

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

195 ## [1] 1 2
```

### 196 1.2.3 Matching strings

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

199 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"
```

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

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

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

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

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

```

227 # first with a single date entry
228 str_match_all(string = c("1970-somemonth-01 or maybe 1990/anothermonth/01"),
229               pattern = "(\\d{4})[\\-\\\\/](\\[a-z\\]+)")

230 ## [[1]]
231 ##           [,1]           [,2]   [,3]
232 ## [1,] "1970-somemonth"   "1970" "somemonth"
233 ## [2,] "1990/anothermonth" "1990" "anothermonth"

234 # then with multiple date entries
235 str_match_all(string = c("1970-somemonth-01 or maybe 1990/anothermonth/01",
236                           "1990-somemonth-01 or maybe 2001/anothermonth/01"),
237               pattern = "(\\d{4})[\\-\\\\/](\\[a-z\\]+)")

238 ## [[1]]
239 ##           [,1]           [,2]   [,3]
240 ## [1,] "1970-somemonth"   "1970" "somemonth"
241 ## [2,] "1990/anothermonth" "1990" "anothermonth"
242 ##
243 ## [[2]]
244 ##           [,1]           [,2]   [,3]
245 ## [1,] "1990-somemonth"   "1990" "somemonth"
246 ## [2,] "2001/anothermonth" "2001" "anothermonth"

```

### 243 1.2.4 Simpler pattern extraction

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

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

```

247 # extracting the first full match using str_extract
248 str_extract(string = c("1970-somemonth-01 or maybe 1990/anothermonth/01",
249                           "1990-somemonth-01 or maybe 2001/anothermonth/01"),
250               pattern = "(\\d{4})[\\-\\\\/](\\[a-z\\]+)")

251 ## [1] "1970-somemonth" "1990-somemonth"

252 # extracting all full matches using str_extract_all
253 str_extract_all(string = c("1970-somemonth-01 or maybe 1990/anothermonth/01",
254                             "1990-somemonth-01 or maybe 2001/anothermonth/01"),
255                  pattern = "(\\d{4})[\\-\\\\/](\\[a-z\\]+)")

```

```

251 ## [[1]]
252 ## [1] "1970-somemonth"      "1990/anothermonth"
253 ##
254 ## [[2]]
255 ## [1] "1990-somemonth"      "2001/anothermonth"

```

### 256 1.2.5 Breaking strings apart

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

260 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 = "[\\-\\/]" )

```

```

261 ## [[1]]
262 ## [1] "1970"      "somemonth" "01"
263 ##
264 ## [[2]]
265 ## [1] "1990"      "anothermonth" "01"

```

266 This can be useful in recovering simulation parameters from a filename, but may  
 267 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 = "_")

```

```

268 ## [[1]]
269 ## [1] "sim"      "param1"   "0.01"     "param2"   "0.05"     "param3"   "0.01.ext"

```

```

# not really
str_split(filename,
          pattern = "sim_")

```

```

270 ## [[1]]
271 ## [1] ""
272 ## [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)")

```



```

273 ## [[1]]
274 ## [1] ""      ""      "0.01" "0.05" "0.01" ""

275 str_split_fixed splits the string into as many pieces as specified, and can be
276 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)

```

```

277 ##      [,1]      [,2]
278 ## [1,] "dir_level_1" "dir_level_2/file.ext"

```

### 1.2.6 Replacing string elements

```

279
280 str_replace is intended to replace the search pattern, and can be co-opted
281 into the task of recovering simulation parameters or other data from regularly
282 named files. str_replace_all works the same way but replaces all matches of
283 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 = " ")

```

```

284 ## [1] " 0.01 0.05 0.01 "

```

```

285 str_remove is a wrapper around str_replace where the replacement is set to
286 "". This is not covered here.

```

```

287 Having replaced unwanted characters in the filename with spaces, str_trim
288 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

```

```

289 ## [1] "0.01 0.05 0.01"

```

```

# the result can be split on whitespaces to return useful data
str_split(filename_without_spaces, " ")

```

```

290 ## [[1]]
291 ## [1] "0.01" "0.05" "0.01"

```

### 292 1.2.7 Subsetting within strings

293 When strings are highly regular, useful data can be extracted from a string using  
 294 `str_sub`. In the date-time example, the year is always represented by the first  
 295 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)
```

296 ## [1] "1970" "1990" "2010"

297 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)
```

298 ## [1] "01" "21" "31"

299 Finally, it's also possible to replace characters within a string based on the  
 300 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
```

301 ## [1] "1970-somemonth-01" "1990-anothermonth-01" "2010-thismonth-01"

### 302 1.2.8 Padding and truncating strings

303 Strings included in filenames or plots are often of unequal lengths, especially  
 304 when they represent numbers. `str_pad` can pad strings with suitable characters  
 305 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")
306 ## [1] "001" "010" "100"
307 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.")
308 ## [1] "bananas are great etc. etc."

```

### 309 1.2.9 Stringr aspects not covered here

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

- 311 • `str_wrap` (of dubious use),
  - 312 • `str_interp`, `str_glue*` (better to use `glue`; see below),
  - 313 • `str_sort`, `str_order` (used in sorting a character vector),
  - 314 • `str_to_case*` (case conversion), and
  - 315 • `str_view*` (a graphical view of search pattern matches).
  - 316 • `word`, `boundary` etc. The use of `word` is covered below.
- 317 `stringi`, of which `stringr` is a wrapper, offers a lot more flexibility and control.

## 318 1.3 String interpolation with glue

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

321 `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')

322 ## The Mazda RX4 is a car model
323 ## The Mazda RX4 Wag is a car model
324 ## The Datsun 710 is a car model
325 ## The Hornet 4 Drive is a car model
326 ## The Hornet Sportabout is a car model
327 ## The Valiant is a car model

```

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

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

```

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

329 # for command line arguments or to start multiple job scripts on the cluster
glue_data(parameter_combinations,
330          'simulation-name {param1} {param2}')

331 ## simulation-name a 1
332 ## simulation-name b 2
333 ## simulation-name c 3
334 ## simulation-name d 4
335 ## simulation-name e 5

336 # for filenames
glue_data(parameter_combinations,
337          'sim_data_param1_{param1}_param2_{param2}.ext')

338 ## sim_data_param1_a_param2_1.ext
339 ## sim_data_param1_b_param2_2.ext
340 ## sim_data_param1_c_param2_3.ext
341 ## sim_data_param1_d_param2_4.ext
342 ## sim_data_param1_e_param2_5.ext

```

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

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

## 347 1.4 Strings in ggplot

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

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

356 More examples can be found on the ggplot website.

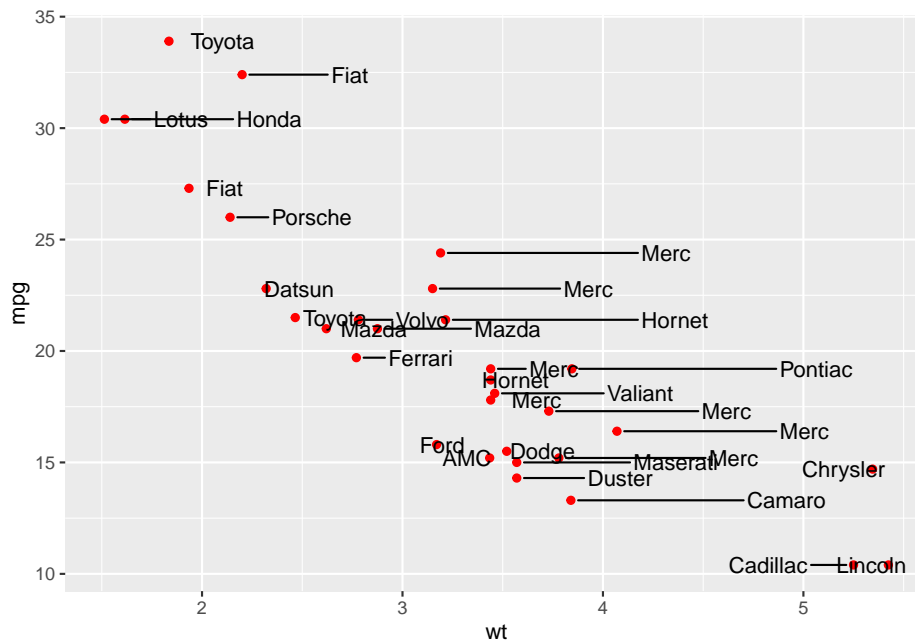
357 Here, the arguments to `geom_text_repel` are taken both from the `mtcars` data  
 358 (position), as well as from the car brands extracted using the `stringr::word`  
 359 (labels), which tries to separate strings based on a regular pattern.

360 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)
```



361

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

365 outliers.

## Chapter 2

# Reshaping data tables in the tidyverse

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 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
```



```
404 ## 2 Theo    3
405 ## 3 Raph    2, 5
```

406 As a rule of thumb, try to convert your tables to tibbles whenever you can,  
 407 especially when the original table is *not* a data frame. For example, the prin-  
 408 cipal component analysis function `prcomp` outputs a `matrix` of coordinates in  
 409 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
## Mazda RX4	-79.596425	2.132241	-2.153336	-2.7073437	-0.7023522
## Mazda RX4 Wag	-79.598570	2.147487	-2.215124	-2.1782888	-0.8843859
## Datsun 710	-133.894096	-5.057570	-2.137950	0.3460330	1.1061111
## Hornet 4 Drive	8.516559	44.985630	1.233763	0.8273631	0.4240145
## Hornet Sportabout	128.686342	30.817402	3.343421	-0.5211000	0.7365801
## Valiant	-23.220146	35.106518	-3.259562	1.4005360	0.8029768

	PC6	PC7	PC8	PC9	PC10
## Mazda RX4	-0.31486106	-0.098695018	-0.07789812	-0.2000092	-0.29008191
## Mazda RX4 Wag	-0.45343873	-0.003554594	-0.09566630	-0.3533243	-0.19283553
## Datsun 710	1.17298584	0.005755581	0.13624782	-0.1976423	0.07634353
## Hornet 4 Drive	-0.05789705	-0.024307168	0.22120800	0.3559844	-0.09057039
## Hornet Sportabout	-0.33290957	0.106304777	-0.05301719	0.1532714	-0.18862217
## Valiant	-0.08837864	0.238946304	0.42390551	0.1012944	-0.03769010

	PC11
## Mazda RX4	0.1057706
## Mazda RX4 Wag	0.1069047
## Datsun 710	0.2668713
## Hornet 4 Drive	0.2088354
## Hornet Sportabout	-0.1092563
## Valiant	0.2757693

```
class(pca_scores) # but is actually a matrix

## [1] "matrix" "array"

# 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
```

```

440 ## 6 -23.2 35.1 -3.26 1.40 0.803 -0.0884 0.239 0.424 0.101 -0.0377
441 ## 7 159. -32.3 0.649 0.199 0.786 0.0687 -0.530 -0.0593 0.221 -0.313
442 ## 8 -113. 39.7 -0.465 0.338 -1.24 0.280 -0.146 0.320 0.279 0.190
443 ## 9 -104. 7.51 -1.59 4.02 -1.14 0.0279 0.595 -0.233 -0.126 -0.349
444 ## 10 -67.0 -6.21 -3.61 -0.320 -0.960 -0.529 -0.0174 -0.182 0.543 0.412
445 ## # ... with 22 more rows, and 1 more variable: PC11 <dbl>

```

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

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

## 451 2.2 2. The concept of tidy data

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

458 Say we have a dataset of several morphometrics measured on Darwin’s finches  
 459 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,
  beak_length = beak_lengths,
  beak_width = beak_widths,
  body_weight = body_weights,
  island = islands
)

# Snapshot
data

460 ## # A tibble: 100 x 5
461 ##       id beak_length beak_width body_weight island

```

```

462 ##      <int>      <dbl>      <dbl>      <dbl> <chr>
463 ##    1      1      5.19      1.97      7.66 Isabela
464 ##    2      2      5.03      2.05     10.6 Isabela
465 ##    3      3      5.09      1.95     12.4 Isabela
466 ##    4      4      5.03      1.93     14.8 Isabela
467 ##    5      5      5.10      1.98     10.6 Isabela
468 ##    6      6      5.06      2.08     19.9 Isabela
469 ##    7      7      5.08      1.82     15.6 Isabela
470 ##    8      8      5.02      1.96      4.30 Isabela
471 ##    9      9      4.92      2.04      9.48 Isabela
472 ##   10     10      5.08      1.96      6.51 Isabela
473 ## # ... with 90 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("beak_length", "beak_width", "body_weight"),
  names_to = "variable"
)
data

480 ## # A tibble: 300 x 4
481 ##       id island variable    value
482 ##   <int> <chr>   <chr>      <dbl>
483 ## 1     1 Isabela beak_length  5.19
484 ## 2     1 Isabela beak_width   1.97
485 ## 3     1 Isabela body_weight  7.66
486 ## 4     2 Isabela beak_length  5.03
487 ## 5     2 Isabela beak_width   2.05
488 ## 6     2 Isabela body_weight 10.6
489 ## 7     3 Isabela beak_length  5.09
490 ## 8     3 Isabela beak_width   1.95
491 ## 9     3 Isabela body_weight 12.4
492 ## 10    4 Isabela beak_length  5.03
493 ## # ... with 290 more rows

```

where each *measurement* (and not each *individual*) is now the unit of observation (the rows). We will come back to the `pivot_longer` function later.

As you can see our tibble now has three times as many rows and fewer columns. This format is rather unintuitive and not optimal for display. However, it provides a very standardized and consistent way of organizing data that will be

understood (and expected) by pretty much all functions in the tidyverse. This makes the tidyverse tools work well together and reduces the time you would otherwise spend reformatting your data from one tool to the next.

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

Another common example of wide-or-long dilemma is when dealing with *contingency tables*. This would be our case, for example, if we asked how many observations we have for each morphometric and each island. We use `table` (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

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)
```

#	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

## 531 2.3 3. Reshaping with tidyr

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

### 542 2.3.1 3.1. Pivoting

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

```

548 pivot_wider(
    data,
    names_from = "variable",
    values_from = "value",
    id_cols = c("id", "island")
)

```

```

548 ## # A tibble: 100 x 5
549 ##       id island beak_length beak_width body_weight
550 ##   <int> <chr>      <dbl>      <dbl>      <dbl>
551 ## 1     1  1 Isabela      5.19      1.97      7.66
552 ## 2     2  2 Isabela      5.03      2.05     10.6
553 ## 3     3  3 Isabela      5.09      1.95     12.4
554 ## 4     4  4 Isabela      5.03      1.93     14.8
555 ## 5     5  5 Isabela      5.10      1.98     10.6
556 ## 6     6  6 Isabela      5.06      2.08     19.9
557 ## 7     7  7 Isabela      5.08      1.82     15.6
558 ## 8     8  8 Isabela      5.02      1.96      4.30
559 ## 9     9  9 Isabela      4.92      2.04      9.48
560 ## 10    10 10 Isabela      5.08      1.96     6.51
561 ## # ... with 90 more rows

```

562 The order of the columns is not exactly as it was, but this should not matter in  
 563 a data analysis pipeline where you should access columns by their names. It is

straightforward to change the order of the columns, but this is more within the scope of the `dplyr` package.

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

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

### 2.3.2 3.2. Handling missing values

Say we have some missing measurements in the column “value” of our finch 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     1 Isabela beak_length  5.19
## 2     2     1 Isabela beak_width  1.97
## 3     3     1 Isabela body_weight NA
## 4     4     2 Isabela beak_length NA
## 5     5     2 Isabela beak_width  2.05
## 6     6     2 Isabela body_weight 10.6
## 7     7     3 Isabela beak_length  5.09
## 8     8     3 Isabela beak_width  1.95
## 9     9     3 Isabela body_weight 12.4
## 10    10     4 Isabela beak_length  5.03
## # ... with 290 more rows
```

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     1 Isabela beak_length  5.19
## 2     2     1 Isabela beak_width  1.97
## 3     3     2 Isabela beak_width  2.05
## 4     4     2 Isabela body_weight 10.6
## 5     5     3 Isabela beak_length  5.09
## 6     6     3 Isabela beak_width  1.95
## 7     7     3 Isabela body_weight 12.4
```

```

598 ## 8      4 Isabela beak_length  5.03
599 ## 9      4 Isabela beak_width   1.93
600 ## 10     5 Isabela beak_length  5.10
601 ## # ... with 190 more rows

```

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

```

        replace_na(data, replace = list(value = -999))

603 ## # A tibble: 300 x 4
604 ##       id island variable      value
605 ##   <int> <chr>   <chr>      <dbl>
606 ## 1     1  1 Isabela beak_length  5.19
607 ## 2     1  1 Isabela beak_width  1.97
608 ## 3     1  1 Isabela body_weight -999
609 ## 4     2  2 Isabela beak_length -999
610 ## 5     2  2 Isabela beak_width  2.05
611 ## 6     2  2 Isabela body_weight 10.6
612 ## 7     3  3 Isabela beak_length  5.09
613 ## 8     3  3 Isabela beak_width  1.95
614 ## 9     3  3 Isabela body_weight 12.4
615 ## 10    4  4 Isabela beak_length  5.03
616 ## # ... with 290 more rows

```

617 where the `replace` argument takes a named list, and the names should refer to  
 618 the columns to apply the replacement to.

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

```

        fill(data, value)

620 ## # A tibble: 300 x 4
621 ##       id island variable      value
622 ##   <int> <chr>   <chr>      <dbl>
623 ## 1     1  1 Isabela beak_length  5.19
624 ## 2     1  1 Isabela beak_width  1.97
625 ## 3     1  1 Isabela body_weight 1.97
626 ## 4     2  2 Isabela beak_length 1.97
627 ## 5     2  2 Isabela beak_width  2.05
628 ## 6     2  2 Isabela body_weight 10.6
629 ## 7     3  3 Isabela beak_length  5.09
630 ## 8     3  3 Isabela beak_width  1.95
631 ## 9     3  3 Isabela body_weight 12.4
632 ## 10    4  4 Isabela beak_length  5.03
633 ## # ... with 290 more rows

```

634 Note that most functions in the tidyverse take a tibble as their first argument,  
 635 and columns to which to apply the functions are usually passed as “objects”  
 636 rather than character strings. In the above example, we passed the `value`

column as `value`, not `"value"`. These column-objects are called by the tidyverse functions *in the context* of the data (the tibble) they belong to.

### 2.3.3 3.3. Splitting and combining cells

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

Say we want to add the specific dates when we took measurements on our birds (we would normally do this using `dplyr` but for now we will stick to the old 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     1 Isabela beak_length  5.19     5     1  2019
## 2     2     1 Isabela beak_width   1.97     5     9  2020
## 3     3     1 Isabela body_weight NA        13     7  2020
## 4     4     2 Isabela beak_length NA         6     5  2019
## 5     5     2 Isabela beak_width   2.05    19     6  2019
## 6     6     2 Isabela body_weight 10.6      2     2  2020
## 7     7     3 Isabela beak_length  5.09    20     9  2020
## 8     8     3 Isabela beak_width   1.95    24     7  2019
## 9     9     3 Isabela body_weight 12.4      3     9  2020
## 10    10     4 Isabela beak_length  5.03    10     1  2020
## # ... with 290 more rows
```

We could combine the `day`, `month` and `year` columns into a single `date` column, 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     1 Isabela beak_length  5.19 5-1-2019
## 2     2     1 Isabela beak_width   1.97 5-9-2020
## 3     3     1 Isabela body_weight NA    13-7-2020
## 4     4     2 Isabela beak_length NA     6-5-2019
## 5     5     2 Isabela beak_width   2.05 19-6-2019
## 6     6     2 Isabela body_weight 10.6  2-2-2020
```



```

670 ## 7      3 Isabela beak_length  5.09 20-9-2020
671 ## 8      3 Isabela beak_width   1.95 24-7-2019
672 ## 9      3 Isabela body_weight 12.4  3-9-2020
673 ## 10     4 Isabela beak_length  5.03 10-1-2020
674 ## # ... with 290 more rows

```

Of course, we can revert back to the previous dataset by splitting the date 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     1 Isabela beak_length  5.19 5      1     2019
## 2     2     1 Isabela beak_width   1.97 5      9     2020
## 3     3     1 Isabela body_weight NA    13     7     2020
## 4     4     2 Isabela beak_length NA     6     5     2019
## 5     5     2 Isabela beak_width   2.05 19     6     2019
## 6     6     2 Isabela body_weight 10.6  2      2     2020
## 7     7     3 Isabela beak_length  5.09 20     9     2020
## 8     8     3 Isabela beak_width   1.95 24     7     2019
## 9     9     3 Isabela body_weight 12.4  3      9     2020
## 10    10    4 Isabela beak_length  5.03 10     1     2020
## # ... with 290 more rows

```

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

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     1 Isabela beak_length  5.19 5
## 2     1     1 Isabela beak_length  5.19 1
## 3     1     1 Isabela beak_length  5.19 2019
## 4     1     1 Isabela beak_width   1.97 5
## 5     1     1 Isabela beak_width   1.97 9
## 6     1     1 Isabela beak_width   1.97 2020
## 7     1     1 Isabela body_weight NA    13
## 8     1     1 Isabela body_weight NA     7
## 9     1     1 Isabela body_weight NA    2020
## 10    2     2 Isabela beak_length NA     6
## # ... with 890 more rows

```

### 2.3.4 3.4. Expanding tables using combinations

Sometimes one may need to quickly create a table with all combinations of a set of variables. We could generate a tibble with all combinations of island-by-morphometric using `expand_grid`:

```
expand_grid(
  island = c("Isabela", "Santa Cruz"),
  variable = c("beak_length", "beak_width", "body_weight")
)

## # A tibble: 6 x 2
##   island      variable
##   <chr>      <chr>
## 1 Isabela    beak_length
## 2 Isabela    beak_width
## 3 Isabela    body_weight
## 4 Santa Cruz beak_length
## 5 Santa Cruz beak_width
## 6 Santa Cruz body_weight
```

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

```
expand(data, island, variable)

## # A tibble: 6 x 2
##   island      variable
##   <chr>      <chr>
## 1 Isabela    beak_length
## 2 Isabela    beak_width
## 3 Isabela    body_weight
## 4 Santa Cruz beak_length
## 5 Santa Cruz beak_width
## 6 Santa Cruz body_weight
```

As an extension of this, the function `complete` can come particularly handy if we need to add missing combinations to our tibble:

```
complete(data, island, variable)

## # A tibble: 300 x 5
##   island variable      id value date
##   <chr>  <chr>      <int> <dbl> <chr>
## 1 Isabela beak_length     1  5.19 5-1-2019
## 2 Isabela beak_length     2  NA    6-5-2019
## 3 Isabela beak_length     3  5.09 20-9-2020
## 4 Isabela beak_length     4  5.03 10-1-2020
## 5 Isabela beak_length     5  5.10 15-11-2020
```

```

743 ## 6 Isabela beak_length      6 NA      19-11-2020
744 ## 7 Isabela beak_length      7 5.08 29-6-2020
745 ## 8 Isabela beak_length      8 5.02 2-10-2020
746 ## 9 Isabela beak_length      9 4.92 17-6-2019
747 ## 10 Isabela beak_length     10 NA      16-10-2020
748 ## # ... with 290 more rows

```

749 which does nothing here because we already have all combinations of `island`  
 750 and `variable`.

### 751 2.3.5 3.5. Nesting

752 The `tidyr` package has yet another feature that makes the tidyverse very pow-  
 753 erful: the `nest` function. However, it makes little sense without combining it  
 754 with the functions in the `purrr` package, so we will not cover it in this chapter  
 755 but rather in the `purrr` chapter.

## 756 2.4 4. Extra: factors and the forcats package

```

library(forcats)

```

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

```

my_char_vec <- c("Pratik", "Theo", "Raph")
my_char_vec

```

```

761 ## [1] "Pratik" "Theo"   "Raph"

```

762 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

```

```

763 ## [1] Pratik Theo   Raph
764 ## Levels: Pratik Raph Theo

```

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

```

factor(c(my_char_vec, "Raph"))
768 ## [1] Pratik Theo   Raph   Raph
769 ## Levels: Pratik Raph Theo

```

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

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

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

But this is a tidyverse tutorial, so I would like to introduce here the package **forcats**, which offers tools to manipulate factors. First of all, most tools from **stringr** *will work* on factors. The **forcats** functions expand the string manipulation toolbox with factor-specific utilities. Similar in philosophy to **stringr** where functions started with **str\_**, in **forcats** most functions start with **fct\_**.

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

#### 2.4.1 4.1. Reordering a factor

Use **fct\_relevel** to manually change the order of the levels:

```
fct_relevel(my_fact_vec, c("Pratik", "Theo", "Raph"))
## [1] Pratik Theo   Raph
## Levels: Pratik Theo Raph
```

Alternatively, use **fct\_inorder** to set the order of the levels to the order in which they appear:

```
fct_inorder(my_fact_vec)
## [1] Pratik Theo   Raph
## Levels: Pratik Theo Raph
```

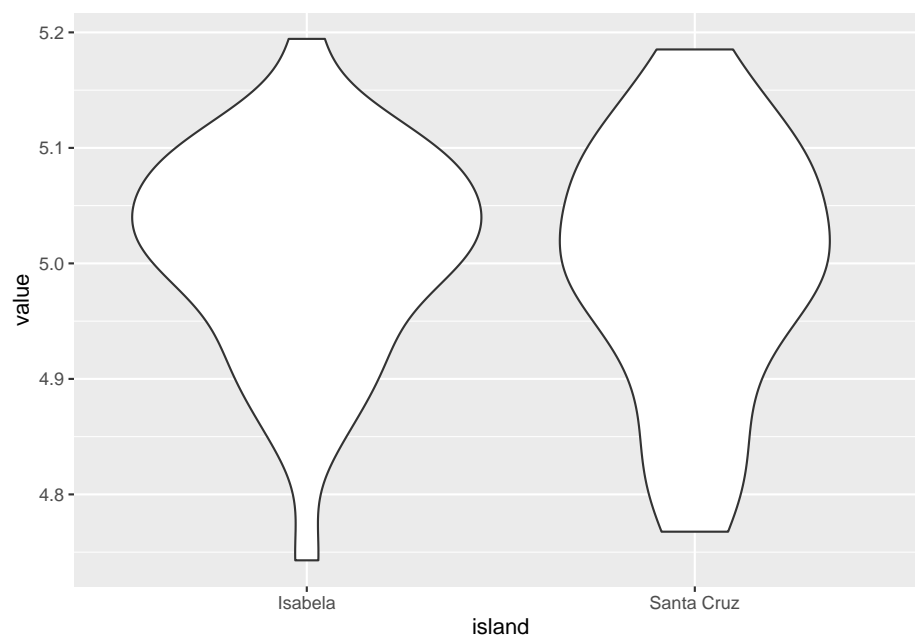
or **fct\_rev** to reverse the order of the levels:

```
fct_rev(my_fact_vec)
## [1] Pratik Theo   Raph
## Levels: Theo Raph Pratik
```

Factor reordering may come useful when plotting categorical variables, for example. Say we want to plot **beak\_length** against **island** in our finch dataset:

```
library(ggplot2)
ggplot(data[data$variable == "beak_length",], aes(x = island, y = value)) +
  geom_violin()
```

802 ## Warning: Removed 27 rows containing non-finite values (stat\_ydensity).



803

804 We could use factor reordering to change the order of the violins:

```
data$island <- fct_relevel(data$island, c("Santa Cruz", "Isabela"))
ggplot(data[data$variable == "beak_length",], aes(x = island, y = value)) +
  geom_violin()
```

805 ## Warning: Removed 27 rows containing non-finite values (stat\_ydensity).



806

807 Lots of other variants exist for reordering (e.g. reordering by association with  
 808 a variable), which we do not cover here. Please refer to the cheatsheet or the  
 809 online documentation for more examples.

## 810 2.4.2 4.2. Factor levels

811 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"
)
```

```
812 ## [1] Pratik Gupte      Theo Pannetier    Raphael Scherrer
813 ## Levels: Pratik Gupte Raphael Scherrer Theo Pannetier
```

814 or collapse factor levels together using `fct_collapse`:

```
fct_collapse(my_fact_vec, EU = c("Theo", "Raph"), NonEU = "Pratik")
```

```
815 ## [1] NonEU EU      EU
816 ## Levels: NonEU EU
```

817 Again, we do not provide an exhaustive list of `forcats` functions here but the  
 818 most usual ones, to give a glimpse of many things that one can do with factors.

819 So, if you are dealing with factors, remember that `forcats` may have handy  
820 tools for you.

### 821 2.4.3 4.3. Bonus: dropping levels

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

```
data <- data[data$island == "Santa Cruz",]  
unique(data$island) # Isabela is gone from the labels  
825 ## [1] Santa Cruz  
826 ## Levels: Santa Cruz Isabela  
  
levels(data$island) # but not from the levels  
827 ## [1] "Santa Cruz" "Isabela"
```

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

```
data <- droplevels(data)  
levels(data$island)  
830 ## [1] "Santa Cruz"
```

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

## 834 2.5 5. External resources

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

- 836 • The `readr/tibble/tidyr` and `forcats` cheatsheets.
- 837 • This link on the concept of tidy data
- 838 • The tibble, tidyr and forcats websites





## Chapter 3

# Working with lists and iteration

Every use case is ridiculous until it happens to you.

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

### 3.1 Basic iteration with map

Iteration in base R is commonly done with `for` and `while` loops. There is no readymade alternative to `while` loops in the tidyverse. However, the functionality of `for` loops is spread over the `map` family of functions.

`purrr` functions are *functionals*, i.e., functions that take another function as an argument. The closest equivalent in R is the `*apply` family of functions: `apply`,

849 `lapply`, `vapply` and so on.

850 A good reason to use `purrr` functions instead of base R functions is their consis-  
851 tent and clear naming, which always indicates how they should be used. This  
852 is explained in the examples below.

853 These reasons, as well as how `map` is different from `for` and `lapply` are best  
854 explained in the Advanced R book.

### 855 3.1.1 `map` basic use

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

```
# get the square root of each integer 1 - 10
some_numbers = 1:10
map(some_numbers, sqrt)

859 ## [[1]]
860 ## [1] 1
861 ##
862 ## [[2]]
863 ## [1] 1.414214
864 ##
865 ## [[3]]
866 ## [1] 1.732051
867 ##
868 ## [[4]]
869 ## [1] 2
870 ##
871 ## [[5]]
872 ## [1] 2.236068
873 ##
874 ## [[6]]
875 ## [1] 2.44949
876 ##
877 ## [[7]]
878 ## [1] 2.645751
879 ##
880 ## [[8]]
881 ## [1] 2.828427
882 ##
883 ## [[9]]
884 ## [1] 3
885 ##
886 ## [[10]]
```

```
887 ## [1] 3.162278
```

### 888 3.1.2 map variants returning vectors

889 Though `map` always returns a list, it has variants named `map_*` where the suffix  
 890 indicates the return type. `map_chr`, `map_dbl`, `map_int`, and `map_lgl` return  
 891 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)

892 ## [1] 1.000000 1.414214 1.732051 2.000000 2.236068 2.449490 2.645751 2.828427
893 ## [9] 3.000000 3.162278

# map_chr will convert the output to a character
map_chr(some_numbers, sqrt)

894 ## [1] "1.000000" "1.414214" "1.732051" "2.000000" "2.236068" "2.449490"
895 ## [7] "2.645751" "2.828427" "3.000000" "3.162278"

# map_int will NOT round the output to an integer

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

896 ## [1] TRUE FALSE FALSE FALSE TRUE TRUE FALSE FALSE
```

### 897 Integrating map and tidyr::nest

898 The example show how each map variant can be used. This integrates  
 899 `tidyr::nest` with `map`, and the two are especially complementary.

```

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

# get the number of rows per dataframe
# the mean mileage
# and the first car
some_data = some_data %>%
  mutate(n_rows = map_int(data, nrow),
         mean_mpg = map_dbl(data, ~mean(.$mpg)),
         first_car = map_chr(data, ~first(.$car_name)))
```

```

some_data

900 ## # A tibble: 3 x 5
901 ## # Groups:   cyl [3]
902 ##   cyl data          n_rows mean_mpg first_car
903 ##   <dbl> <list>      <int>   <dbl> <chr>
904 ## 1     6 <tibble [7 x 11]>     7    19.7 Mazda RX4
905 ## 2     4 <tibble [11 x 11]>    11    26.7 Datsun 710
906 ## 3     8 <tibble [14 x 11]>   14    15.1 Hornet Sportabout

907 map accepts multiple functions that are applied in sequence to the input list-like
908 object, but this is confusing to the reader and ill advised.

```

### 909 3.1.3 map variants returning dataframes

910 `map_df` returns data frames, and by default binds dataframes by rows, while  
 911 `map_dfr` does this explicitly, and `map_dfc` does returns a dataframe bound by  
 912 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)

```

```

913 ##   mpg cyl  disp  hp drat   wt  qsec vs am gear carb
914 ## 1 22.8  4 108.0  93 3.85 2.320 18.61  1  1   4     1
915 ## 2 24.4  4 146.7  62 3.69 3.190 20.00  1  0   4     2
916 ## 3 21.0  6 160.0 110 3.90 2.620 16.46  0  1   4     4
917 ## 4 21.0  6 160.0 110 3.90 2.875 17.02  0  1   4     4
918 ## 5 18.7  8 360.0 175 3.15 3.440 17.02  0  0   3     2
919 ## 6 14.3  8 360.0 245 3.21 3.570 15.84  0  0   3     4

```

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

922 `map_dfr` behaves the same as `map_df`.

```

# the same as above but with a pipe
some_list %>%
  map_dfr(head, n = 2)

```

```

923 ##   mpg cyl  disp  hp drat   wt  qsec vs am gear carb
924 ## 1 22.8  4 108.0  93 3.85 2.320 18.61  1  1   4     1
925 ## 2 24.4  4 146.7  62 3.69 3.190 20.00  1  0   4     2
926 ## 3 21.0  6 160.0 110 3.90 2.620 16.46  0  1   4     4
927 ## 4 21.0  6 160.0 110 3.90 2.875 17.02  0  1   4     4
928 ## 5 18.7  8 360.0 175 3.15 3.440 17.02  0  0   3     2

```

```
929 ## 6 14.3    8 360.0 245 3.21 3.570 15.84  0  0    3    4
```

930 `map_dfc` binds the resulting 3 data frames of two rows each by column, and  
 931 automatically repairs the column names, adding a suffix to each duplicate.

```
some_list %>%
  map_dfc(head, n = 2)

932 ##      mpg cyl  disp hp drat   wt  qsec vs am gear carb mpg1 cyl1 disp1 hp1 drat1
933 ## 1 22.8   4 108.0 93 3.85 2.32 18.61 1  1   4   1  21   6  160 110   3.9
934 ## 2 24.4   4 146.7 62 3.69 3.19 20.00 1  0   4   2  21   6  160 110   3.9
935 ##      wt1 qsec1 vs1 am1 gear1 carb1 mpg2 cyl2 disp2 hp2 drat2  wt2 qsec2 vs2 am2
936 ## 1 2.620 16.46   0   1     4     4 18.7   8  360 175  3.15 3.44 17.02   0   0
937 ## 2 2.875 17.02   0   1     4     4 14.3   8  360 245  3.21 3.57 15.84   0   0
938 ##      gear2 carb2
939 ## 1         3     2
940 ## 2         3     4
```

### 941 3.1.4 Selective mapping

- 942 • `map_at` and `map_if`

## 943 3.2 More map variants

### 944 3.2.1 `map2`

945 `imap` here

### 946 3.2.2 `pmap`

### 947 3.2.3 `walk`

948 `walk2` and `pwalk`

## 949 3.3 Modification in place

950 `modify`

## 951 **3.4 Working with lists**

### 952 **3.4.1 Filtering lists**

### 953 **3.4.2 Summarising lists**

### 954 **3.4.3 Reduction and accumulation**

### 955 **3.4.4 Miscellaneous operation**

## Chapter 4

# Data manipulation with dplyr

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

### 4.1 Introduction

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

### 4.2 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

```

973 ## 5 Abrocoma_bo~ Rodentia Abrocomid~ Abrocoma boliviensis      1      0
974 ## 6 Abrocoma_bu~ Rodentia Abrocomid~ Abrocoma budini          1      0
975 ## 7 Abrocoma_ci~ Rodentia Abrocomid~ Abrocoma cinerea         1      0
976 ## 8 Abrocoma_fa~ Rodentia Abrocomid~ Abrocoma famatina        1      0
977 ## 9 Abrocoma_sh~ Rodentia Abrocomid~ Abrocoma shistacea        1      0
978 ## 10 Abrocoma_us~ Rodentia Abrocomid~ Abrocoma uspollata       1      0
979 ## # ... with 5,821 more rows, and 17 more variables: Freshwater <dbl>,
980 ## # Aerial <dbl>, Life.Habit.Method <chr>, Life.Habit.Source <chr>,
981 ## # Mass.g <dbl>, Mass.Method <chr>, Mass.Source <chr>, Mass.Comparison <chr>,
982 ## # Mass.Comparison.Source <chr>, Island.Endemicity <chr>,
983 ## # IUCN.Status.1.2 <chr>, Added.IUCN.Status.1.2 <chr>, Diet.Plant <dbl>,
984 ## # Diet.Vertebrate <dbl>, Diet.Invertebrate <dbl>, Diet.Method <chr>,
985 ## # Diet.Source <chr>

```

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

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

The first operation I'm going to run on this table is changing the names with `rename()`. Some people prefer their tea without sugar, and I prefer my variable names without uppercase characters, dots or (if possible) numbers. This will 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
  )

```

999 For convenience, I'm going to use the pipe operator (`%>%`) that we've seen before,  
 1000 through this chapter. All `dplyr` functions are built to work with the pipe (i.e.,  
 1001 their first argument is always `data`), but again, this is not compulsory. I could  
 1002 do

```

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

```

1003 Note how columns are referred to. Once the data has been passed as an argument,  
 1004 no need to refer to it anymore, `dplyr` understands that you're dealing with  
 1005 variables inside that data frame. So drop that `data$var`, `data[, "var"]`, and,  
 1006 if you've read *The R book*, forget the very existence of `attach()`.

1007 Finally, I should mention that you can refer to variables names either with  
 1008 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"
  )

```

### 1009 4.3 Select variables with `select()`

1010 To extract a set of variables (i.e. columns), use the conveniently-named  
 1011 `select()`.

```

1012 phylacine_subset <- phylacine %>%
1013   dplyr::select(
1014     binomial,
1015     order,
1016     terrestrial,
1017     marine,
1018     freshwater,
1019     aerial
1020   )
1021 phylacine_subset
1022
1023 ## # A tibble: 5,831 x 6
1024 ##   binomial      order terrestrial marine freshwater aerial
1025 ##   <chr>      <chr>      <dbl>  <dbl>      <dbl>  <dbl>
1026 ## 1 Abditomys_latidens Rodentia      1      0          0      0
1027 ## 2 Abeomelomys_sevia Rodentia      1      0          0      0
1028 ## 3 Abrawayaomys_ruschii Rodentia      1      0          0      0
1029 ## 4 Abrocoma_bennettii Rodentia      1      0          0      0
1030 ## 5 Abrocoma_boliviensis Rodentia      1      0          0      0
1031 ## 6 Abrocoma_budini Rodentia      1      0          0      0
1032 ## 7 Abrocoma_cinerea Rodentia      1      0          0      0
1033 ## 8 Abrocoma_famatina Rodentia      1      0          0      0
1034 ## 9 Abrocoma_shistacea Rodentia      1      0          0      0
1035 ## 10 Abrocoma_ustallata Rodentia      1      0          0      0
1036 ## # ... with 5,821 more rows
1037
1038 # Single variable
1039 phylacine %>% dplyr::select(family)
1040
1041 ## # A tibble: 5,831 x 1
1042 ##   family
1043 ##   <chr>
1044 ## 1 Muridae
1045 ## 2 Muridae
1046 ## 3 Cricetidae
1047 ## 4 Abrocomidae
1048 ## 5 Abrocomidae
1049 ## 6 Abrocomidae
1050 ## 7 Abrocomidae
1051 ## 8 Abrocomidae
1052 ## 9 Abrocomidae
1053 ## 10 Abrocomidae

```

```

1039 ## # ... with 5,821 more rows

      # A set of variables
phylacine %>% dplyr::select(genus, species, mass_g)

1040 ## # A tibble: 5,831 x 3
1041 ##   genus      species      mass_g
1042 ##   <chr>      <chr>      <dbl>
1043 ## 1 Abditomys latidens      269
1044 ## 2 Abeomelomys sevia        52
1045 ## 3 Abrawayaomys ruschii      63
1046 ## 4 Abrocoma bennettii      250
1047 ## 5 Abrocoma boliviensis     158
1048 ## 6 Abrocoma budini        361.
1049 ## 7 Abrocoma cinerea        250
1050 ## 8 Abrocoma famatina       233.
1051 ## 9 Abrocoma shistacea       276.
1052 ## 10 Abrocoma uspallata      246.
1053 ## # ... with 5,821 more rows

      # A range of variables
phylacine %>% dplyr::select(family:terrestrial)

1054 ## # A tibble: 5,831 x 4
1055 ##   family      genus      species      terrestrial
1056 ##   <chr>      <chr>      <chr>      <dbl>
1057 ## 1 Muridae    Abditomys latidens      1
1058 ## 2 Muridae    Abeomelomys sevia      1
1059 ## 3 Cricetidae Abrawayaomys ruschii      1
1060 ## 4 Abrocomidae Abrocoma bennettii      1
1061 ## 5 Abrocomidae Abrocoma boliviensis      1
1062 ## 6 Abrocomidae Abrocoma budini      1
1063 ## 7 Abrocomidae Abrocoma cinerea      1
1064 ## 8 Abrocomidae Abrocoma famatina      1
1065 ## 9 Abrocomidae Abrocoma shistacea      1
1066 ## 10 Abrocomidae Abrocoma uspallata      1
1067 ## # ... with 5,821 more rows

1068 select() can also exclude variables:

phylacine %>% dplyr::select(-binomial)

1069 ## # A tibble: 5,831 x 23
1070 ##   order family genus species terrestrial marine freshwater aerial
1071 ##   <chr> <chr> <chr> <chr>      <dbl> <dbl> <dbl> <dbl>
1072 ## 1 Rode~ Murid~ Abdi~ latide~      1     0     0     0
1073 ## 2 Rode~ Murid~ Abeo~ sevia      1     0     0     0
1074 ## 3 Rode~ Crice~ Abra~ ruschii      1     0     0     0
1075 ## 4 Rode~ Abroc~ Abro~ bennet~      1     0     0     0

```

```

1076 ## 5 Rode~ Abroc~ Abro~ bolivi~          1      0      0      0
1077 ## 6 Rode~ Abroc~ Abro~ budini          1      0      0      0
1078 ## 7 Rode~ Abroc~ Abro~ cinerea         1      0      0      0
1079 ## 8 Rode~ Abroc~ Abro~ famati~         1      0      0      0
1080 ## 9 Rode~ Abroc~ Abro~ shista~         1      0      0      0
1081 ## 10 Rode~ Abroc~ Abro~ uspoll~         1      0      0      0
1082 ## # ... with 5,821 more rows, and 15 more variables: life_habit_method <chr>,
1083 ## #   life_habit_source <chr>, mass_g <dbl>, mass_method <chr>,
1084 ## #   mass_source <chr>, mass_comparison <chr>, mass_comparison_source <chr>,
1085 ## #   island_endemicity <chr>, iucn_status <chr>, added_iucn_status <chr>,
1086 ## #   diet_plant <dbl>, diet_vertebrate <dbl>, diet_invertebrate <dbl>,
1087 ## #   diet_method <chr>, diet_source <chr>

phylacine %>% dplyr::select(-(binomial:species))

1088 ## # A tibble: 5,831 x 19
1089 ##   terrestrial marine freshwater aerial life_habit_meth~ life_habit_sour~ mass_g
1090 ##   <dbl> <dbl> <dbl> <dbl> <chr> <chr> <dbl>
1091 ## 1      1      0      0      0 Reported IUCN. 2016. IUC~ 269
1092 ## 2      1      0      0      0 Reported IUCN. 2016. IUC~ 52
1093 ## 3      1      0      0      0 Reported IUCN. 2016. IUC~ 63
1094 ## 4      1      0      0      0 Reported IUCN. 2016. IUC~ 250
1095 ## 5      1      0      0      0 Reported IUCN. 2016. IUC~ 158
1096 ## 6      1      0      0      0 Reported IUCN. 2016. IUC~ 361.
1097 ## 7      1      0      0      0 Reported IUCN. 2016. IUC~ 250
1098 ## 8      1      0      0      0 Reported IUCN. 2016. IUC~ 233.
1099 ## 9      1      0      0      0 Reported IUCN. 2016. IUC~ 276.
1100 ## 10     1      0      0      0 Reported IUCN. 2016. IUC~ 246.
1101 ## # ... with 5,821 more rows, and 12 more variables: mass_method <chr>,
1102 ## #   mass_source <chr>, mass_comparison <chr>, mass_comparison_source <chr>,
1103 ## #   island_endemicity <chr>, iucn_status <chr>, added_iucn_status <chr>,
1104 ## #   diet_plant <dbl>, diet_vertebrate <dbl>, diet_invertebrate <dbl>,
1105 ## #   diet_method <chr>, diet_source <chr>

select() and rename() are pretty similar, and in fact, select() can also
1106 rename variables along the way:

1107 phylacine %>% dplyr::select("fam" = family)

1108 ## # A tibble: 5,831 x 1
1109 ##   fam
1110 ##   <chr>
1111 ## 1 Muridae
1112 ## 2 Muridae
1113 ## 3 Cricetidae
1114 ## 4 Abrocomidae
1115 ## 5 Abrocomidae
1116 ## 6 Abrocomidae

```

```

1117 ## 7 Abrocomidae
1118 ## 8 Abrocomidae
1119 ## 9 Abrocomidae
1120 ## 10 Abrocomidae
1121 ## # ... with 5,821 more rows

```

1122 And you can mix all of that at once:

```

phylacine %>% dplyr::select(
  "fam" = family,
  genus:freshwater,
  -terrestrial
)

1123 ## # A tibble: 5,831 x 5
1124 ##   fam      genus      species    marine freshwater
1125 ##   <chr>    <chr>    <chr>      <dbl>      <dbl>
1126 ## 1 Muridae  Abditomys latidens      0          0
1127 ## 2 Muridae  Abeomelomys sevia      0          0
1128 ## 3 Cricetidae Abrawayaomys ruschii      0          0
1129 ## 4 Abrocomidae Abrocoma bennettii      0          0
1130 ## 5 Abrocomidae Abrocoma boliviensis      0          0
1131 ## 6 Abrocomidae Abrocoma budini      0          0
1132 ## 7 Abrocomidae Abrocoma cinerea      0          0
1133 ## 8 Abrocomidae Abrocoma famatina      0          0
1134 ## 9 Abrocomidae Abrocoma shistacea      0          0
1135 ## 10 Abrocomidae Abrocoma uspallata      0          0
1136 ## # ... with 5,821 more rows

```

## 1137 4.4 Select observations with filter()

1138 Conditional selection of observations is performed through `filter()`. This is  
 1139 arguably the most useful function in the entire package. The syntax uses  
 1140 conditions involving the variables, just as you would use for `if` statements or  
 1141 `while` loops.

1142 For example, I might want to extract mammals that are above a certain mass:

```

# megafauna
phylacine %>%
  dplyr::filter(mass_g > 1e5) # 100 kg

1143 ## # A tibble: 302 x 24
1144 ##   binomial order family genus species terrestrial marine freshwater aerial
1145 ##   <chr>    <chr> <chr> <chr> <chr>      <dbl>  <dbl>      <dbl>  <dbl>
1146 ## 1 Ailurop~ Carn~ Ursid~ Ailu~ melano~      1    0          0    0
1147 ## 2 Alcelap~ Ceta~ Bovid~ Alce~ busela~      1    0          0    0

```

```

1148 ## 3 Alces_a~ Ceta~ Cervi~ Alces alces          1      0      0      0
1149 ## 4 Archaeo~ Prim~ Palae~ Arch~ fontoy~          1      0      0      0
1150 ## 5 Arctoce~ Carn~ Otari~ Arct~ forste~          1      1      0      0
1151 ## 6 Arctoce~ Carn~ Otari~ Arct~ pusill~          1      1      0      0
1152 ## 7 Arctoce~ Carn~ Otari~ Arct~ townse~          1      1      0      0
1153 ## 8 Arctodu~ Carn~ Ursid~ Arct~ simus           1      0      0      0
1154 ## 9 Arctoth~ Carn~ Ursid~ Arct~ tarije~          1      0      0      0
1155 ## 10 Babyrou~ Ceta~ Suidae Baby~ togean~          1      0      1      0
1156 ## # ... with 292 more rows, and 15 more variables: life_habit_method <chr>,
1157 ## #   life_habit_source <chr>, mass_g <dbl>, mass_method <chr>,
1158 ## #   mass_source <chr>, mass_comparison <chr>, mass_comparison_source <chr>,
1159 ## #   island_endemicity <chr>, iucn_status <chr>, added_iucn_status <chr>,
1160 ## #   diet_plant <dbl>, diet_vertibrate <dbl>, diet_invertebrate <dbl>,
1161 ## #   diet_method <chr>, diet_source <chr>

# living megafauna
phylacine %>%
  dplyr::filter(mass_g > 1e5, iucn_status != "EP")

1162 ## # A tibble: 178 x 24
1163 ##   binomial order family genus species terrestrial marine freshwater aerial
1164 ##   <chr>      <chr> <chr> <chr> <chr>          <dbl>  <dbl>          <dbl>  <dbl>
1165 ## 1 Ailurop~ Carn~ Ursid~ Ailu~ melano~          1      0              0      0
1166 ## 2 Alcelap~ Ceta~ Bovid~ Alce~ busela~          1      0              0      0
1167 ## 3 Alces_a~ Ceta~ Cervi~ Alces alces          1      0              0      0
1168 ## 4 Arctoce~ Carn~ Otari~ Arct~ forste~          1      1              0      0
1169 ## 5 Arctoce~ Carn~ Otari~ Arct~ pusill~          1      1              0      0
1170 ## 6 Arctoce~ Carn~ Otari~ Arct~ townse~          1      1              0      0
1171 ## 7 Babyrou~ Ceta~ Suidae Baby~ togean~          1      0              1      0
1172 ## 8 Balaena~ Ceta~ Balae~ Bala~ mystic~          0      1              0      0
1173 ## 9 Balaeno~ Ceta~ Balae~ Bala~ acutor~          0      1              0      0
1174 ## 10 Balaeno~ Ceta~ Balae~ Bala~ bonaer~          0      1              0      0
1175 ## # ... with 168 more rows, and 15 more variables: life_habit_method <chr>,
1176 ## #   life_habit_source <chr>, mass_g <dbl>, mass_method <chr>,
1177 ## #   mass_source <chr>, mass_comparison <chr>, mass_comparison_source <chr>,
1178 ## #   island_endemicity <chr>, iucn_status <chr>, added_iucn_status <chr>,
1179 ## #   diet_plant <dbl>, diet_vertibrate <dbl>, diet_invertebrate <dbl>,
1180 ## #   diet_method <chr>, diet_source <chr>

# Are there any flying mammals besides bats?
phylacine %>%
  dplyr::filter(aerial == 1, order != "Chiroptera")

1181 ## # A tibble: 0 x 24
1182 ## # ... with 24 variables: binomial <chr>, order <chr>, family <chr>,
1183 ## #   genus <chr>, species <chr>, terrestrial <dbl>, marine <dbl>,
1184 ## #   freshwater <dbl>, aerial <dbl>, life_habit_method <chr>,
1185 ## #   life_habit_source <chr>, mass_g <dbl>, mass_method <chr>,

```

```

1186 ## # mass_source <chr>, mass_comparison <chr>, mass_comparison_source <chr>,
1187 ## # island_endemicity <chr>, iucn_status <chr>, added_iucn_status <chr>,
1188 ## # diet_plant <dbl>, diet_vertibrate <dbl>, diet_invertebrate <dbl>,
1189 ## # diet_method <chr>, diet_source <chr>

# no :(

# That one species
phylacine %>% dplyr::filter(binomial == "Homo_sapiens")

1190 ## # A tibble: 1 x 24
1191 ##   binomial order family genus species terrestrial marine freshwater aerial
1192 ##   <chr>      <chr> <chr> <chr> <chr>          <dbl> <dbl>          <dbl> <dbl>
1193 ## 1 Homo_sa~ Prim~ Homin~ Homo sapiens          1      0              0      0
1194 ## # ... with 15 more variables: life_habit_method <chr>, life_habit_source <chr>,
1195 ## # mass_g <dbl>, mass_method <chr>, mass_source <chr>, mass_comparison <chr>,
1196 ## # mass_comparison_source <chr>, island_endemicity <chr>, iucn_status <chr>,
1197 ## # added_iucn_status <chr>, diet_plant <dbl>, diet_vertibrate <dbl>,
1198 ## # diet_invertebrate <dbl>, diet_method <chr>, diet_source <chr>

1199 Tip: dplyr introduces the useful function between() that does exactly what
1200 the name implies

between(1:5, 2, 4)

1201 ## [1] FALSE TRUE TRUE TRUE FALSE

# Mesofauna
phylacine %>% dplyr::filter(mass_g > 1e3, mass_g < 1e5) # base

1202 ## # A tibble: 1,126 x 24
1203 ##   binomial order family genus species terrestrial marine freshwater aerial
1204 ##   <chr>      <chr> <chr> <chr> <chr>          <dbl> <dbl>          <dbl> <dbl>
1205 ## 1 Acerodo~ Chir~ Ptero~ Acer~ jubatus          0      0              0      1
1206 ## 2 Acinony~ Carn~ Felid~ Acin~ jubatus          1      0              0      0
1207 ## 3 Acratoc~ Pilo~ Megal~ Acra~ odontr~          1      0              0      0
1208 ## 4 Acratoc~ Pilo~ Megal~ Acra~ ye              1      0              0      0
1209 ## 5 Addax_n~ Ceta~ Bovid~ Addax nasoma~          1      0              0      0
1210 ## 6 Aepycer~ Ceta~ Bovid~ Aepy~ melamp~          1      0              0      0
1211 ## 7 Aepypry~ Dipr~ Potor~ Aepy~ rufesc~          1      0              0      0
1212 ## 8 Aeromys~ Rode~ Sciur~ Aero~ tephro~          1      0              0      0
1213 ## 9 Aeromys~ Rode~ Sciur~ Aero~ thomasi          1      0              0      0
1214 ## 10 Agalmac~ Ceta~ Cervi~ Agal~ blicki          1      0              0      0
1215 ## # ... with 1,116 more rows, and 15 more variables: life_habit_method <chr>,
1216 ## # life_habit_source <chr>, mass_g <dbl>, mass_method <chr>,
1217 ## # mass_source <chr>, mass_comparison <chr>, mass_comparison_source <chr>,
1218 ## # island_endemicity <chr>, iucn_status <chr>, added_iucn_status <chr>,
1219 ## # diet_plant <dbl>, diet_vertibrate <dbl>, diet_invertebrate <dbl>,
1220 ## # diet_method <chr>, diet_source <chr>

```

```

1221 phylacine %>% dplyr::filter(mass_g %>% between(1e3, 1e5)) # dplyr
1222 ## # A tibble: 1,148 x 24
1223 ##   binomial order family genus species terrestrial marine freshwater aerial
1224 ##   <chr>      <chr> <chr> <chr> <chr>          <dbl>  <dbl>        <dbl>  <dbl>
1225 ## 1 Acerodo~ Chir~ Ptero~ Acer~ jubatus          0      0            0      1
1226 ## 2 Acinony~ Carn~ Felid~ Acin~ jubatus          1      0            0      0
1227 ## 3 Acratoc~ Pilo~ Megal~ Acra~ odontr~          1      0            0      0
1228 ## 4 Acratoc~ Pilo~ Megal~ Acra~ ye            1      0            0      0
1229 ## 5 Addax_n~ Ceta~ Bovid~ Addax nasoma~          1      0            0      0
1230 ## 6 Aepycer~ Ceta~ Bovid~ Aepy~ melamp~          1      0            0      0
1231 ## 7 Aepypry~ Dipr~ Potor~ Aepy~ rufesc~          1      0            0      0
1232 ## 8 Aeromys~ Rode~ Sciur~ Aero~ tephro~          1      0            0      0
1233 ## 9 Aeromys~ Rode~ Sciur~ Aero~ thomasi          1      0            0      0
1234 ## 10 Agalmac~ Ceta~ Cervi~ Agal~ blicki          1      0            0      0
1235 ## # ... with 1,138 more rows, and 15 more variables: life_habit_method <chr>,
1236 ## #   life_habit_source <chr>, mass_g <dbl>, mass_method <chr>,
1237 ## #   mass_source <chr>, mass_comparison <chr>, mass_comparison_source <chr>,
1238 ## #   island_endemicity <chr>, iucn_status <chr>, added_iucn_status <chr>,
1239 ## #   diet_plant <dbl>, diet_vertibrate <dbl>, diet_invertebrate <dbl>,
1240 ## #   diet_method <chr>, diet_source <chr>
1241
1240 If instead of selecting observations based on conditions, you want to get the nth
1241 row, use slice()

1242 phylacine %>% dplyr::slice(1:3)
1243 ## # A tibble: 3 x 24
1244 ##   binomial order family genus species terrestrial marine freshwater aerial
1245 ##   <chr>      <chr> <chr> <chr> <chr>          <dbl>  <dbl>        <dbl>  <dbl>
1246 ## 1 Abditom~ Rode~ Murid~ Abdi~ latide~          1      0            0      0
1247 ## 2 Abeomel~ Rode~ Murid~ Abeo~ sevia          1      0            0      0
1248 ## 3 Abaway~ Rode~ Crice~ Abra~ ruschii          1      0            0      0
1249 ## # ... with 15 more variables: life_habit_method <chr>, life_habit_source <chr>,
1250 ## #   mass_g <dbl>, mass_method <chr>, mass_source <chr>, mass_comparison <chr>,
1251 ## #   mass_comparison_source <chr>, island_endemicity <chr>, iucn_status <chr>,
1252 ## #   added_iucn_status <chr>, diet_plant <dbl>, diet_vertibrate <dbl>,
1253 ## #   diet_invertebrate <dbl>, diet_method <chr>, diet_source <chr>
1254
1255 # can also be used to exclude rows
1256 phylacine %>%
1257   dplyr::slice(1:3) %>%
1258   dplyr::slice(-2)
1259
1260 ## # A tibble: 2 x 24
1261 ##   binomial order family genus species terrestrial marine freshwater aerial
1262 ##   <chr>      <chr> <chr> <chr> <chr>          <dbl>  <dbl>        <dbl>  <dbl>
1263 ## 1 Abditom~ Rode~ Murid~ Abdi~ latide~          1      0            0      0
1264 ## 2 Abaway~ Rode~ Crice~ Abra~ ruschii          1      0            0      0

```



```

1258 ## # ... with 15 more variables: life_habit_method <chr>, life_habit_source <chr>,
1259 ## #   mass_g <dbl>, mass_method <chr>, mass_source <chr>, mass_comparison <chr>,
1260 ## #   mass_comparison_source <chr>, island_endemicity <chr>, iucn_status <chr>,
1261 ## #   added_iucn_status <chr>, diet_plant <dbl>, diet_vertibrate <dbl>,
1262 ## #   diet_invertebrate <dbl>, diet_method <chr>, diet_source <chr>

```

1263 You may also want the row number to be included as a variable, to be used as  
 1264 part of some condition or calculation. `tibble` is there for you:

```
phylacine %>% tibble::rownames_to_column()
```

```

1265 ## # A tibble: 5,831 x 25
1266 ##   rowname binomial order family genus species terrestrial marine freshwater
1267 ##   <chr>   <chr>   <chr> <chr> <chr> <chr>         <dbl>  <dbl>      <dbl>
1268 ## 1 1      Abditom~ Rode~ Murid~ Abdi~ latide~         1      0          0
1269 ## 2 2      Abeomel~ Rode~ Murid~ Abeo~ sevia         1      0          0
1270 ## 3 3      Abraway~ Rode~ Crice~ Abra~ ruschii         1      0          0
1271 ## 4 4      Abrocom~ Rode~ Abroc~ Abro~ bennet~         1      0          0
1272 ## 5 5      Abrocom~ Rode~ Abroc~ Abro~ bolivi~         1      0          0
1273 ## 6 6      Abrocom~ Rode~ Abroc~ Abro~ budini         1      0          0
1274 ## 7 7      Abrocom~ Rode~ Abroc~ Abro~ cinerea         1      0          0
1275 ## 8 8      Abrocom~ Rode~ Abroc~ Abro~ famati~         1      0          0
1276 ## 9 9      Abrocom~ Rode~ Abroc~ Abro~ shista~         1      0          0
1277 ## 10 10     Abrocom~ Rode~ Abroc~ Abro~ uspall~         1      0          0
1278 ## # ... with 5,821 more rows, and 16 more variables: aerial <dbl>,
1279 ## #   life_habit_method <chr>, life_habit_source <chr>, mass_g <dbl>,
1280 ## #   mass_method <chr>, mass_source <chr>, mass_comparison <chr>,
1281 ## #   mass_comparison_source <chr>, island_endemicity <chr>, iucn_status <chr>,
1282 ## #   added_iucn_status <chr>, diet_plant <dbl>, diet_vertibrate <dbl>,
1283 ## #   diet_invertebrate <dbl>, diet_method <chr>, diet_source <chr>

```

## 1284 4.5 Create new variables with mutate()

1285 can also edit existing ones

1286 drop existing variables with `transmute()`

## 1287 4.6 Grouped results with `group_by()` and 1288 `summarise()`

## 1289 4.7 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")))
```

## 1290 4.8 More !

1291 dolla sign x point operator variables values -> `dplyr::distinct()` eq. to  
1292 `base::unique()` `sample_n()` `sample_frac()` first, last, nth