

Data manipulation

Contents

Introduction	1
Data types	1
Lists and data frames	2
Loading, viewing, and summarising data	4
Data transformation with dplyr verbs	6
Data processing pipelines	9
Grouping and summarisation	11
Final exercise	13

Introduction

This is the first programming practical. If you haven't yet done so, open the project file `02_Data_manipulation.Rproj` in RStudio. You can choose to write the answers to your exercises in either an `.R` file or in an `.Rmd` file. Example answer files are provided in the project directory (`example_answers.Rmd` and `example_answers.R`). You can open these from the files pane and use them as a starting point. While working through the exercises, write down your code in one of these files. Use proper style and provide comments so you can read it back later and still understand what is happening.

The practicals always start with the packages we are going to use. Be sure to run these lines in your session to load their functions before you continue. If there are packages that you have not yet installed, first install them with `install.packages()`.

```
library(ISLR)
library(tidyverse)
library(haven)
library(readxl)
```

Data types

There are several data types in R. Here is a table with the most common ones:

Type	Short	Example
Integer	int	0, 1, 2, 3, -4, -5
Numeric / Double	dbl	0.1, -2.5, 123.456
Character	chr	"dav is a cool course"
Logical	lgl	TRUE / FALSE
Factor	fctr	low, medium, high

The `class()` function can give you an idea about what type of data each variable contains.

1. Run the following code in R and inspect their data types using the `class()` function. Try to guess beforehand what their types will be!

```
object_1 <- 1:5
object_2 <- 1L:5L
object_3 <- "-123.456"
object_4 <- as.numeric(object_2)
object_5 <- letters[object_1]
object_6 <- as.factor(rep(object_5, 2))
object_7 <- c(1, 2, 3, "4", "5", "6")
```

the factor data type is special to R and uncommon in other programming languages. It is used to represent categorical variables with fixed possible values. For example, when there is a multiple choice question with 5 possible choices (a to e) and 10 students answer the question, we may get a result as in `object_6`.

Vectors can have only a single data type. Note that the first three elements in `object_7` have been converted. We can convert to different data types using the `as.<class>()` functions.

2. Convert `object_7` back to a vector of numbers using the `as.numeric()` function

```
object_7 <- as.numeric(object_7)
```

Lists and data frames

A list is a collection of objects. The elements may have names, but it is not necessary. Each element of a list can have a different data type, unlike vectors.

3. Make a list called `objects` containing object 1 to 7 using the `list()` function.

```
objects <- list(object_1, object_2, object_3, object_4, object_5, object_6,
               object_7)
```

You can select elements of a list using its name (`objects$elementname`) or using its index (`objects[[1]]` for the first element).

A special type of list is the `data.frame`. It is the same as a list, but each element is forced to have the same length and a name. The elements of a `data.frame` are the columns of a dataset. In the tidyverse, `data.frames` are called `tibbles` and they are printed in a nice way, as we will see later.

-
4. **Make a data frame out of `object_1`, `object_2`, and `object_5` using the `data.frame()` function**

```
dat <- data.frame(Var1 = object_1, Var2 = object_2, Var3 = object_5)
dat
```

```
##   Var1 Var2 Var3
## 1     1     1    a
## 2     2     2    b
## 3     3     3    c
## 4     4     4    d
## 5     5     5    e
```

Just like a list, the columns in a data frame (the variables in a dataset) can be accessed using their name `df$columnname` or their index `df[[1]]`. Additionally, the tenth row can be selected using `df[10,]`, the second column using `df[, 2]` and cell number 10, 2 can be accessed using `df[10, 2]`. This is because data frames also behave like the `matrix` data type in addition to the `list` type.

-
5. **Useful functions for determining the size of a data frame are `ncol()` and `nrow()`. Try them out!**

```
ncol(dat)
```

```
## [1] 3
```

```
nrow(dat)
```

```
## [1] 5
```

Loading, viewing, and summarising data

We are going to use a dataset from Kaggle - the Google play store apps data by user lava18. We have downloaded it into the data folder already from <https://www.kaggle.com/lava18/google-play-store-apps> (downloaded on 2018-09-28).

Tidyverse contains many data loading functions – each for their own file type – in the packages readr (default file types), readxl (excel files), and haven (external file types such as from SPSS or Stata). The most common file type is csv, which is what we use here.

-
6. Use the function `read_csv()` to import the file “data/googleplaystore.csv” and store it in a variable called `apps`.
-

```
apps <- read_csv("data/googleplaystore.csv")
```

```
## Parsed with column specification:
## cols(
##   App = col_character(),
##   Category = col_character(),
##   Rating = col_double(),
##   Reviews = col_integer(),
##   Size = col_character(),
##   Installs = col_character(),
##   Type = col_character(),
##   Price = col_character(),
##   `Content Rating` = col_character(),
##   Genres = col_character(),
##   `Last Updated` = col_character(),
##   `Current Ver` = col_character(),
##   `Android Ver` = col_character()
## )
```

If necessary, use the help files. These import functions from the tidyverse are fast and safe: they display informative errors if anything goes wrong. `read_csv()` also displays a message with information on how each column is imported: which variable type each column gets.

-
7. Did any column get a variable type you did not expect?
-

```
# Several columns such as price and number of installs were imported as
# character data types, but they are numbers.
```

8. Use the function `head()` to look at the first few rows of the apps dataset

```
head(apps)
```

```
## # A tibble: 6 x 13
##   App    Category Rating Reviews Size  Installs Type  Price `Content Rating`
##   <chr> <chr>    <dbl>   <int> <chr> <chr>    <chr> <chr> <chr>
## 1 Phot~ ART_AND~    4.1     159 19M    10,000+ Free  0    Everyone
## 2 Colo~ ART_AND~    3.9     967 14M    500,000+ Free  0    Everyone
## 3 "U L~ ART_AND~    4.7   87510 8.7M    5,000,0~ Free  0    Everyone
## 4 Sket~ ART_AND~    4.5  215644 25M    50,000,~ Free  0    Teen
## 5 Pixe~ ART_AND~    4.3     967 2.8M   100,000+ Free  0    Everyone
## 6 Pape~ ART_AND~    4.4     167 5.6M    50,000+ Free  0    Everyone
## # ... with 4 more variables: Genres <chr>, `Last Updated` <chr>, `Current
## #   Ver` <chr>, `Android Ver` <chr>
```

9. Repeat steps 5, 6, and 7 but now for “data/students.xlsx” (NB: You’ll need a function from the package `readxl`). Also try out the function `tail()` and `View()` (with a capital V).

```
students <- read_xlsx("data/students.xlsx")
head(students)
```

```
## # A tibble: 6 x 3
##   student_number grade programme
##           <dbl> <dbl> <chr>
## 1      5117250  6.54 A
## 2      6562582  7.57 A
## 3      6000241  6.08 B
## 4      4862862  7.71 A
## 5      6561723  6.57 B
## 6      5625916  7.90 B
```

```
tail(students)
```

```
## # A tibble: 6 x 3
##   student_number grade programme
##           <dbl> <dbl> <chr>
## 1      5062746  7.43 B
## 2      6560954  7.04 B
## 3      6120285  6.71 A
## 4      6553913  8.24 A
## 5      4181101  5.62 B
## 6      4639846  4.84 A
```

-
10. **Create a summary of the three columns in the students dataset using the `summary()` function. What is the range of the grades achieved by the students?**
-

```
summary(students)
```

```
## student_number      grade      programme
## Min.   :4011659   Min.    :4.844   Length:37
## 1st Qu.:4862862   1st Qu.:6.390   Class :character
## Median :6000241   Median :7.151   Mode  :character
## Mean   :5686729   Mean    :6.991
## 3rd Qu.:6553913   3rd Qu.:7.573
## Max.   :6997130   Max.    :9.291
```

Data transformation with dplyr verbs

The tidyverse package dplyr contains functions to transform, rearrange, and filter data frames.

Filter

The first verb is `filter()`, which selects rows from a data frame. [Chapter 5 of R4DS](#) states that to use filtering effectively, you have to know how to select the observations that you want using the comparison operators. R provides the standard suite: `>`, `>=`, `<`, `<=`, `!=` (not equal), and `==` (equal).

When you're starting out with R, the easiest mistake to make is to use `=` instead of `==` when testing for equality.

-
11. **Look at the help pages for `filter()` (especially the examples) and show the students with a grade lower than 5.5**
-

```
filter(students, grade < 5.5)
```

```
## # A tibble: 3 x 3
##   student_number grade programme
##           <dbl> <dbl> <chr>
## 1       6114656  5.16 A
## 2       5265402  5.49 B
## 3       4639846  4.84 A
```

12. Show only the students with a grade higher than 8 from programme A

If you are unsure how to proceed, read [Section 5.2.2 from R4DS](#).

```
filter(students, grade > 8, programme == "A")
```

```
## # A tibble: 5 x 3
##   student_number grade programme
##         <dbl> <dbl> <chr>
## 1      6352581  8.09 A
## 2      6165611  8.02 A
## 3      4133949  8.40 A
## 4      4011659  8.94 A
## 5      6553913  8.24 A
```

Arrange

The second verb is `arrange()`, which sorts a data frame by one or more columns.

13. Sort the students dataset such that the students from programme A are on top of the data frame and within the programmes the highest grades come first.

```
arrange(students, programme, -grade)
```

```
## # A tibble: 37 x 3
##   student_number grade programme
##         <dbl> <dbl> <chr>
## 1      4011659  8.94 A
## 2      4133949  8.40 A
## 3      6553913  8.24 A
## 4      6352581  8.09 A
## 5      6165611  8.02 A
## 6      6997130  7.75 A
## 7      4862862  7.71 A
## 8      6562582  7.57 A
## 9      4483974  7.46 A
## 10     5128923  7.26 A
## # ... with 27 more rows
```

Select

The third verb is `select()`, which selects columns of interest.

14. Show only the `student_number` and `programme` columns from the `students` dataset

```
select(students, student_number, programme)
```

```
## # A tibble: 37 x 2
##   student_number programme
##           <dbl> <chr>
## 1         5117250 A
## 2         6562582 A
## 3         6000241 B
## 4         4862862 A
## 5         6561723 B
## 6         5625916 B
## 7         4096023 A
## 8         6114656 A
## 9         5265402 B
## 10        5977188 B
## # ... with 27 more rows
```

```
# or, equivalently: select(students, -grade)
```

Mutate

With `mutate()` you can compute new columns and transform existing columns as functions of the columns in your dataset. For example, we may create a new logical column in the `students` dataset to indicate whether a student has passed or failed:

```
students <- mutate(students, pass = grade > 5.5)
students
```

```
## # A tibble: 37 x 4
##   student_number grade programme pass
##           <dbl> <dbl> <chr>    <lgl>
## 1         5117250  6.54 A      TRUE
## 2         6562582  7.57 A      TRUE
## 3         6000241  6.08 B      TRUE
## 4         4862862  7.71 A      TRUE
## 5         6561723  6.57 B      TRUE
## 6         5625916  7.90 B      TRUE
```



```
## 7      4096023  5.92 A      TRUE
## 8      6114656  5.16 A     FALSE
## 9      5265402  5.49 B     FALSE
## 10     5977188  7.26 B      TRUE
## # ... with 27 more rows
```

Now, the students dataset has an extra column named “pass”.

You can also transform existing columns with the `mutate()` function. For example, we may want to transform the programme column to an actual programme name according to this table:

Code	Name
A	Science
B	Social Science

-
15. Use `mutate()` and `recode()` to change the codes in the programme column of the students dataset to their names. Store the result in a variable called `students_recoded`
-

```
students_recoded <- mutate(students,
  programme = recode(programme, "A" = "Science", "B" = "Social Science")
)
```

[Chapter 5 of R4DS](#) neatly summarises the five key dplyr functions that allow you to solve the vast majority of your data manipulation challenges:

- Pick observations by their values (`filter()`).
- Reorder the rows (`arrange()`).
- Pick variables by their names (`select()`).
- Create new variables with functions of existing variables (`mutate()`).

Cleaning data files and extracting the most useful information is essential to any downstream steps such as plotting or analysis. Make sure you know exactly which variable types are in your tibbles / data frames!

Data processing pipelines

A very useful feature in tidyverse is the pipe `%>%`. The pipe inputs the result from the left-hand side as the first argument of the right-hand side function: `filter(students, grade > 5.5)` becomes `students %>% filter(grade > 5.5)`. With the pipe, a set of processing steps becomes a neatly legible data processing pipeline!

Different tasks we have performed on the students dataset can be done in one pipeline like so:

```
students_dataset <-
  read_xlsx("data/students.xlsx") %>%
  mutate(prog = recode(programme, "A" = "Science", "B" = "Social Science")) %>%
  filter(grade > 5.5) %>%
  arrange(programme, -grade) %>%
  select(student_number, prog, grade)

students_dataset
```

```
## # A tibble: 34 x 3
##   student_number prog    grade
##           <dbl> <chr>   <dbl>
## 1         4011659 Science  8.94
## 2         4133949 Science  8.40
## 3         6553913 Science  8.24
## 4         6352581 Science  8.09
## 5         6165611 Science  8.02
## 6         6997130 Science  7.75
## 7         4862862 Science  7.71
## 8         6562582 Science  7.57
## 9         4483974 Science  7.46
## 10        5128923 Science  7.26
## # ... with 24 more rows
```

In one statement, we have loaded the dataset from disk, recoded the programme variable, filtered only students that pass, reordered the rows and selected the relevant columns only. We did not need to save intermediate results or nest functions deeply.

-
16. Create a data processing pipeline that (a) loads the apps dataset, (b) parses the number of downloads using `mutate` and `parse_number()`, (c) shows only apps with more than 500 000 000 downloads, (d) orders them by rating (best on top), and (e) shows only the relevant columns (you can choose which are relevant, but select at least the Rating and Category variables). Save the result under the name `popular_apps`.
-

If you find duplicates, you may need to use `distinct(App, .keep_all = TRUE)` as the last step in your pipeline to remove duplicate app names. Tip: `ctrl/cmd + shift + M` inserts a pipe operator in RStudio.

```
popular_apps <-
  read_csv("data/googleplaystore.csv") %>%
  mutate(Downloads = parse_number(Installs)) %>%
  filter(Downloads > 5e8) %>% # 5e8 is the same as 5 x 10^8
  arrange(-Rating) %>%
```

```
select(App, Rating, Reviews, Downloads, Category) %>%
distinct(App, .keep_all = TRUE)
```

popular_apps

```
## # A tibble: 20 x 5
##   App                                Rating  Reviews Downloads Category
##   <chr>                            <dbl>    <int>    <dbl> <chr>
## 1 Subway Surfers                    4.5  27722264  1.00e9 GAME
## 2 Instagram                        4.5  66577313  1.00e9 SOCIAL
## 3 Google Photos                    4.5  10858556  1.00e9 PHOTOGRAPHY
## 4 WhatsApp Messenger               4.4  69119316  1.00e9 COMMUNICATION
## 5 Google                           4.4   8033493  1.00e9 TOOLS
## 6 Google Drive                     4.4   2731171  1.00e9 PRODUCTIVITY
## 7 Google Chrome: Fast & Secure     4.3   9642995  1.00e9 COMMUNICATION
## 8 Gmail                            4.3  4604324  1.00e9 COMMUNICATION
## 9 Google Play Games                4.3  7165362  1.00e9 ENTERTAINMENT
## 10 Maps - Navigate & Explore        4.3  9235155  1.00e9 TRAVEL_AND_LO~
## 11 YouTube                         4.3  25655305  1.00e9 VIDEO_PLAYERS
## 12 Google+                         4.2   4831125  1.00e9 SOCIAL
## 13 Google Street View               4.2   2129689  1.00e9 TRAVEL_AND_LO~
## 14 Skype - free IM & video calls    4.1  10484169  1.00e9 COMMUNICATION
## 15 Facebook                        4.1  78158306  1.00e9 SOCIAL
## 16 "Messenger \x96 Text and Vide~  4    56642847  1.00e9 COMMUNICATION
## 17 Hangouts                        4     3419249  1.00e9 COMMUNICATION
## 18 Google Play Books                3.9   1433233  1.00e9 BOOKS_AND_REF~
## 19 Google News                     3.9    877635  1.00e9 NEWS_AND_MAGA~
## 20 Google Play Movies & TV         3.7    906384  1.00e9 VIDEO_PLAYERS
```

Grouping and summarisation

We have now seen how we can transform and clean our datasets. The next step is to start exploring the dataset by computing relevant summary statistics, such as means, ranges, variances, differences, etc. We have already used the function `summary()` which comes with R, but `dplyr` has extra summary functionality in the form of the `summarise()` (or `summarize()`) verb.

An example to get the mean grade of the `students_dataset` we made earlier is below:

```
students_dataset %>%
  summarise(
    mean = mean(grade),
    variance = var(grade),
    min = min(grade),
```

```

    max = max(grade)
  )

## # A tibble: 1 x 4
##   mean variance   min   max
##   <dbl>     <dbl> <dbl> <dbl>
## 1  7.15     0.817  5.53  9.29

```

17. Show the median, minimum, and maximum for the popular apps dataset you made in the previous assignment._

```

popular_apps %>%
  summarise(
    med = median(Rating),
    min = min(Rating),
    max = max(Rating),
  )

```

```

## # A tibble: 1 x 3
##   med   min   max
##   <dbl> <dbl> <dbl>
## 1  4.3   3.7   4.5

```

The summarise() function works with any function that takes a vector of numbers and outputs a single number. For example, we can create our own [Median Absolute Deviation \(MAD\)](#) function:

```

mad <- function(x) {
  median(abs(x - median(x)))
}

students_dataset %>% summarise(mad = mad(grade))

```

```

## # A tibble: 1 x 1
##   mad
##   <dbl>
## 1 0.591

```

18. Add the median absolute deviation to the summaries you made before

```

popular_apps %>%
  summarise(
    med = median(Rating),

```

```

    min = min(Rating),
    max = max(Rating),
    mad = mad(Rating)
  )

```

```

## # A tibble: 1 x 4
##   med   min   max   mad
##   <dbl> <dbl> <dbl> <dbl>
## 1   4.3   3.7   4.5  0.15

```

By itself, the `summarise()` function is not very useful; we can also simply use the `summary()` function or directly enter the vector we are interested in as an argument to the functions: `mad(students_dataset$grade)` = 0.5908503. The power of `summarise()` is in its combined use with the `group_by()` function, which makes it easy to make grouped summaries:

```

students_dataset %>%
  group_by(prog) %>%
  summarise(
    mean = mean(grade),
    variance = var(grade),
    min = min(grade),
    max = max(grade)
  )

```

```

## # A tibble: 2 x 5
##   prog          mean variance   min   max
##   <chr>        <dbl>    <dbl> <dbl> <dbl>
## 1 Science         7.35     0.736  5.92  8.94
## 2 Social Science  6.96     0.868  5.53  9.29

```

19. Create a grouped summary of the ratings per category in the popular apps dataset.

Final exercise

20. Create an interesting summary based on the Google play store apps dataset. An example could be “do games get higher ratings than communication apps?”

```

read_csv("data/googleplaystore.csv") %>%
  filter(Category == "GAME" | Category == "COMMUNICATION") %>%

```

```

select(App, Category, Rating) %>%
distinct(App, .keep_all = TRUE) %>%
group_by(Category) %>%
summarise(
  mean = mean(Rating, na.rm = TRUE),
  median = median(Rating, na.rm = TRUE)
)

```

```

## # A tibble: 2 x 3
##   Category      mean median
##   <chr>      <dbl>  <dbl>
## 1 COMMUNICATION 4.12    4.2
## 2 GAME          4.25    4.3

```