

Data Literate with R

Nicolas Meseth

Table of contents

Preface	3
Download materials	3
 I Data Transformation	 4
1 Operations	6
2 Select columns	7
2.1 By column names	7
2.2 By name patterns	8
2.2.1 Names starting with a string	8
2.2.2 Names ending with a string	8
2.2.3 Names with a string anywhere	8
2.2.4 Using regular expressions	8
2.3 By data type	8
3 Filter rows	10
4 Add columns	11
5 Summarize rows	12
6 Sort rows	13
 II Data Visualization	 14
7 Overview	16

Preface

Download materials

You can download the ZIP-archive with all material [here](#). This archive includes:

Folder	Content
book	The compiled book in PDF format
data	All data from the chapters
docs	All chapters as single PDF files
exercises	All exercises as PDF files (sometimes with solutions)
scripts	All code from the chapters as plain R-Scripts (.R)
slides	A collection of slide decks in PDF format

Part I

Data Transformation

This part introduces the basic tools for data transformation with R.

1 Operations

Data is the new oil, according to the mathematician [Clive Humby](#):

“Data is the new oil. Like oil, data is valuable, but if unrefined, it cannot really be used. It has to be changed into gas, plastic, chemicals, etc. to create a valuable entity that drives profitable activity. So, must data be broken down, analysed for it to have value.”

If we take this analogy seriously, the data, like oil, needs to be refined to turn it into something of value. Two important tools for refining data into a valuable output are *data transformation* and *data visualization*, both of which are the main focus of this book. In this part of the book, we first need to learn how to transform data so that we can apply visualization later on.

To learn how to transform data, we need to learn how to to the following operations:

- Remove any variables we don’t currently need (or specify those we **do** need)
- Remove any records we don’t currently need (or specify those we **do** need)
- Add new variables that don’t exist yet
- Summarize many records into one or a few numbers
- Change the order of the records

The goal of the following chapters is to introduce means to perform theses five operations with R.

2 Select columns

This chapter introduces tools to remove unnecessary columns from the data set. Or, if stated in a positive manner, we learn how to specify the columns we need for our analysis. As with most data transformation operations, we mostly introduce functions from the `{dplyr}` package.

The function `select()` is the designated tool to select columns with `{dplyr}`. By passing different things to the function, we can efficiently define the set of columns in the resulting data frame.

2.1 By column names

The easiest and intuitive way to specify the columns we want is by listing their names. We can pass one or more column names to the `select()` function. In case of two or more, we use commas to separate the names:

```
# Just one column name
orders %>%
  select(order_id)

#> # A tibble: 2,874 x 1
#>   order_id
#>   <dbl>
#> 1 1130007101519
#> 2 1130014965839
#> 3 1130026958927
#> ...

# A list of column names
orders %>%
  select(order_id, total_price)

#> # A tibble: 2,874 x 2
#>   order_id total_price
#>   <dbl>      <dbl>
#> 1 1130007101519      94.7
```

```
#> 2 1130014965839      32.2
#> 3 1130026958927      30.2
#> ...
```

When we only want a few columns, this approach works fine and is usually a good choice. I expect you apply this method in more than 90% of all cases. However, there are cases when you'd wish there was something more flexible. Luckily, there is.

2.2 By name patterns

2.2.1 Names starting with a string

Sometimes we want to select columns based on a pattern of their names. Take the orders data set as an example. Here, all variables that contain information about the shipping address have the prefix `shipping`. We leverage this with the helper function `starts_with()`:

```
orders %>%
  select(starts_with("shipping")) %>%
  colnames()

#> [1] "shipping_address_city"      "shipping_address_zip"      "shipping_address_country"
#> [4] "shipping_address_latitude"  "shipping_address_longitude"
```

2.2.2 Names ending with a string

2.2.3 Names with a string anywhere

2.2.4 Using regular expressions

2.3 By data type

```
orders %>%
  select(where(is.numeric))

orders %>%
  select(where(is.logical))

orders %>%
```



```
      select(where(is.character))

orders %>%
  select(where(is.factor))

orders %>%
  select(where(is.list))

# The package lubridate provides a function to check for date (without time)
orders %>%
  select(where(lubridate::is.Date))

# Select all date/time columns
orders %>%
  select(where(lubridate::is.POSIXct))
```

3 Filter rows

4 Add columns

5 Summarize rows

6 Sort rows

Part II

Data Visualization

This part introduces the basic tools for data visualization with R.

7 Overview