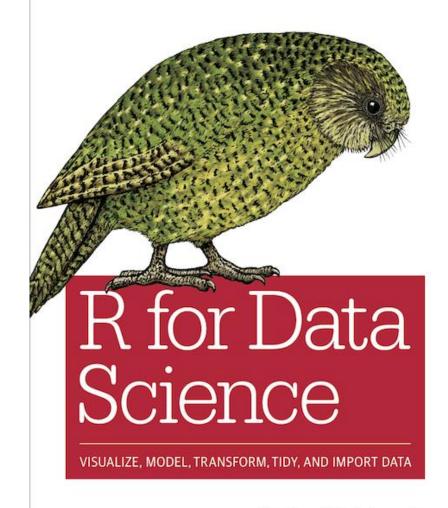


https://r4ds.hadley.nz/

O'REILLY"



Hadley Wickham & Garrett Grolemund

R for Data Science

Search

Table of contents

Welcome

1 Introduction

Explore

- 2 Introduction
- 3 Data visualisation
- 4 Workflow: basics
- 5 Data transformation
- 6 Workflow: scripts
- 7 Exploratory Data Analysis
- 8 Workflow: projects

Wrangle

- 9 Introduction
- 10 Tibbles
- 11 Data import
- 12 Tidy data
- 13 Relational data
- 14 Strings
- 15 Factors
- 16 Dates and times

Program

- 17 Introduction
- 18 Pipes
- 19 Functions
- 20 Vectors
- 21 Iteration

Model

- 22 Introduction
- 23 Model basics
- 24 Model building
- 25 Many models

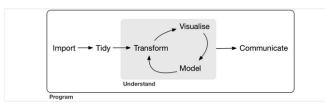
1 Introduction

You're reading the first edition of R4DS; for the latest on this topic see the **Introduction chapter** in the second edition.

Data science is an exciting discipline that allows you to turn raw data into understanding, insight, and knowledge. The goal of "R for Data Science" is to help you learn the most important tools in R that will allow you to do data science. After reading this book, you'll have the tools to tackle a wide variety of data science challenges, using the best parts of R.

1.1 What you will learn

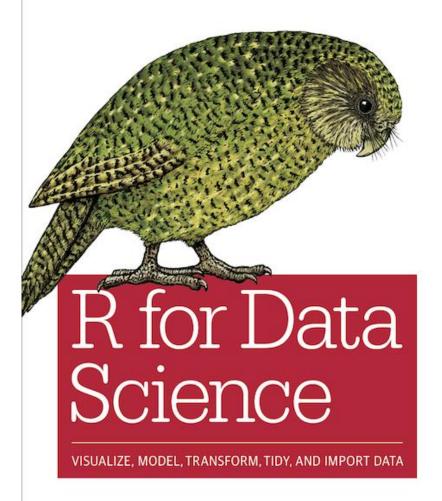
Data science is a huge field, and there's no way you can master it by reading a single book. The goal of this book is to give you a solid foundation in the most important tools. Our model of the tools needed in a typical data science project looks something like this:



First you must **import** your data into R. This typically means that you take data stored in a file, database, or web application programming interface (API), and load it into a data frame in R. If you can't get your data into R, you can't do data science on it!

Once you've imported your data, it is a good idea to **tidy** it. Tidying your data means storing it in a consistent form that matches the semantics of the dataset with the way it is stored. In brief, when your data is tidy, each column is a variable, and each row is an observation. Tidy data is important because the consistent structure lets you focus your struggle on questions about the data, not fighting to get the data into the right form for different functions.

O'REILLY'



Hadley Wickham & Garrett Grolemund



ggplot2



Massimo Bourquin October 2022

Multivariate statistics in R



ggplot2

Wrangle lubridate stringr

Massimo Bourquin October 2022

Multivariate statistics in R



Package #1: dplyr

Data wrangling

A Grammar of Data Manipulation

• A fast, consistent tool for working with data frames

Widely used in ecology and data science in R

• Easy syntax, simplifies complicated tasks

Filtering rows in a data frame

• filter(data, mask1, mask2, etc.)

```
(dec25 <- filter(flights, month == 12, day == 25))
                                                                 Copy
#> # A tibble: 719 x 19
     year month day dep time sched dep time dep delay arr time sched a
    <int> <int> <int>
                        <int>
                                                <dbl>
                                       <int>
                                                         <int>
     2013
             12
                          456
                                         500
                                                           649
                   25
                                                   -4
     2013
                 25
                          524
                                         515
                                                           805
           12
     2013
            12
                 25
                          542
                                         540
                                                           832
#> 3
     2013
            12
                  25
                          546
                                         550
                                                          1022
             12
                  25
                          556
                                         600
                                                           730
     2013
     2013
             12
                   25
                          557
                                         600
                                                   -3
                                                           743
#> # ... with 713 more rows, and 11 more variables: arr_delay <dbl>, carrie
      flight <int>, tailnum <chr>, origin <chr>, dest <chr>, air_time <dl
#> #
      distance <dbl>, hour <dbl>, minute <dbl>, time hour <dttm>
```

Reordering rows in a data frame

arrange(data, col1, col2, etc.)

```
arrange(flights, year, month, day)
                                                                Copy
#> # A tibble: 336,776 x 19
     year month day dep_time sched_dep_time dep_delay arr_time sched_a
    <int> <int> <int> <int>
                                               <dbl>
                                      <int>
                                                       <int>
#> 1
     2013
                          517
                                        515
                                                         830
#> 2
     2013
                          533
                                        529
                                                         850
     2013 1 1
                          542
                                        540
                                                         923
                          544
#> 4
     2013
                                        545
                                                        1004
                          554
                                        600
                                                         812
     2013
                                                 -6
#> 6
     2013
                          554
                                        558
                                                         740
#> # ... with 336,770 more rows, and 11 more variables: arr_delay <dbl>,
      carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest <ch
      air time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time hou
```

Selecting columns

select(data, col1, col2, etc.)

```
# Select columns by name
                                                             Copy
select(flights, year, month, day)
#> # A tibble: 336,776 x 3
     year month day
    <int> <int> <int>
#> 1 2013 1
#> 2 2013 1 1
     2013 1 1
     2013 1 1
     2013 1 1
#> 6 2013
#> # ... with 336,770 more rows
# Select all columns between year and day (inclusive)
```

Creating new columns based on existing ones

mutate(data, new_col = calculation)

```
mutate(flights_sml,
 gain = dep_delay - arr_delay,
 speed = distance / air_time * 60
#> # A tibble: 336,776 x 9
     year month day dep_delay arr_delay distance air_time gain speed
    <int> <int> <int>
                      <dbl>
                               <dbl>
                                       <dbl>
                                              <dbl> <dbl> <dbl>
    2013
                                        1400
                                                227
                                                      -9 370.
#> 2
    2013 1 1
                                 20
                                        1416
                                                227 -16 374
                                        1089
     2013 1 1
                                  33
                                                160
                                                     -31 408.
    2013
                                                183
                                                      17 517.
                         -1
                                 -18
                                        1576
    2013
                                 -25
                                        762
                         -6
                                                116 19 394.
#> 6
     2013
                                 12
                                        719
                                                150
                                                     -16 288.
                         -4
#> # ... with 336,770 more rows
```

Grouping data and summarising values

- group_by(data, col1, col2, etc.)
- summarise(data, new_summary = calculation)

```
by_day <- group_by(flights, year, month, day)</pre>
                                                             Copy
summarise(by_day, delay = mean(dep_delay, na.rm = TRUE))
#> `summarise()` regrouping output by 'year', 'month' (override with `.grangering)
#> # A tibble: 365 x 4
#> # Groups: year, month [12]
#> year month day delay
#> <int> <int> <dbl>
#> 1 2013 1 1 11.5
#> 2 2013 1 2 13.9
#> 3 2013 1 3 11.0
#> 4 2013 1 4 8.95
#> 5 2013 1 5 5.73
#> 6 2013 1 6 7.15
#> # ... with 359 more rows
```

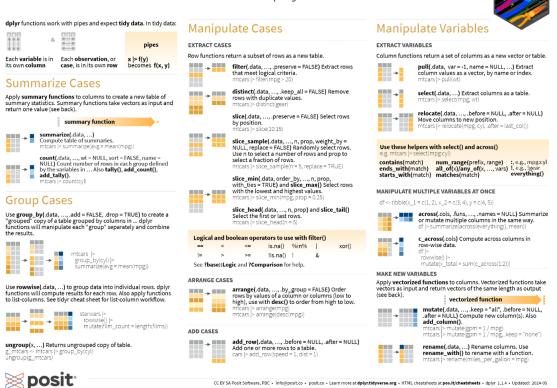
Summary

- filter(data, mask1, mask2, etc.) > filter rows
- arrange(data, col1, col2, etc.) > order rows
- select(data, col1, col2, etc.) > select columns
- mutate(data, new_col = calculation) > create new columns based on existing ones
- group_by(data, col1, col2, etc.) > group a dataframe based on categorical columns
- **summarise(data, new_summary = calculation)** > summarise data e.g. groups means

Cheatsheet

https://github.com/rstudio/cheatsheets/blob/main/data-transformation.pdf

Data transformation with dplyr:: **CHEATSHEET**



Vectorized Functions Summary Functions Combine Tables COMBINE VARIABLES TO USE WITH MUTATE () TO USE WITH SUMMARIZE () COMBINE CASES mutate() applies vectorized functions to summarize() applies summary functions to columns to create new columns. Vectorized columns to create a new table. Summary 2 t 1 X b u 2 functions take vectors as input and return functions take vectors as input and return single vectors of the same length as output. values as output. bind cols(.....name repair) Returns tables vectorized function summary function placed side by side as a single table. Column engths must be equal. Columns will NOT be other as a single table. Set .id to matched by id (to do that look at Relational Data OFFSET COUNT below), so be sure to check that both tables are v d w 4 pictured). dplyr::lag() - offset elements by 1 dplyr::lead() - offset elements by -1 r::n() - number of values/rows ordered the way you want before binding. ::n distinct() - # of uniques RELATIONAL DATA CUMULATIVE AGGREGATE Use a "Filtering Join" to filter one table against dplyr::cumall() - cumulative all() Use a "Mutating Join" to join one table to cumany() - cumulative any() columns from another, matching values with the the rows of another. mean() - mean, also mean(!is.na()) cummax() - cumulative max() rows that they correspond to. Each join retains a median() - mediar dplyr::cummean() - cumulative mean() cumulative min() cumprod() - cumulative prod() msum() - cumulative sum(mean() - proportion of TRUEs left_join(x, y, by = NULL, copy = FALSE, suffix = c(".x", ".y"), ..., keep = FALSE, na_matches = "na") Join matching semi_join(x, y, by = NULL, copy = FALSE, ..., na_matches = "na") Return rows of x /r::cume_dist() - proportion of all values <= /r::dense_rank() - rank w ties = min, no gaps values from v to x ORDER that have a match in y. Use to see what will be included in a join. ::first() - first value right_join(x, y, by = NULL, copy = FALSE, suffix = c(".x", ".y"), ..., keep = FALSE, na_matches = "na") Join matching ::min_rank() - rank with ties = min ::last() - last value :ntile() - bins into n bins anti_join(x, y, by = NULL, copy = FALSE, ..., na_matches = "na") Return rows of x ::nth() - value in nth location of vector percent_rank() - min_rank scaled to [0,1] row number() - rank with ties = "first values from x to v. that do not have a match in v. Use to see RANK what will not be included in a join. inner_join(x, y, by = NULL, copy = FALSE, suffix = c("x", "y"), ..., keep = FALSE, na_matches = "na") Join data. Retain quantile() - nth quantile +, -, *, /, ^, %/%, %% - arithmetic ops log(), log2(), log10() - logs nin() - minimum value Use a "Nest Join" to inner join one table to max() - maximum value another into a nested data frame. <, <=, >, >=, !=, == - logical comparisons ::between() - x >= left & x <= right only rows with matches. full_join(x, y, by = NULL, copy = FALSE, s t 1 suffix = c("x", ".y"), ..., keep = FALSE, b u 2 suffix = c("x", "ly"), loin data. Retain all ::near() - safe == for floating point numbers IQR() - Inter-Quartile Range nad() - median absolute deviation MISCELLANEOUS sd() - standard deviation dplyr::case_when() - multi-case if_else() dw M 1 values all rows var() - variance mutate(type = case when SET OPERATIONS height > 200 | mass > 200 ~ "large" Row Names COLUMN MATCHING FOR JOINS species == "Droid" TRUE intersect(x, y, ...) Rows that appear in both x and y. Tidy data does not use rownames, which store a Use by = c("col1", "col2", ...) to variable outside of the columns. To work with the specify one or more common dplyr::coalesce() - first non-NA values by rownames, first move them into a column. columns to match on. element across a set of vectors if else() - element-wise if() + else(tibble::rownames_to_column() Move row names into col. Rows that appear in x but not y. r::na_if() - replace specific values with NA Use a named vector, by = c("col1" = pmax() - element-wise max() "col2"), to match on columns that

rownames to column(var = "C")

tibble::column_to_rownames()

Move col into row names.

Also tibble::has_rownames() and tibble::remove_rownames(

have different names in each table.

Use suffix to specify the suffix to

give to unmatched columns that

have the same name in both tables

 $left_join(x, y, by = c("C" = "D")$

bind rows(....id = NULL)

Returns tables one on top of the

a column name to add a column

of the original table names (as

FALSE, keep = FALSE, name =

NULL, ...) Join data, nesting

data frame column.

Rows that appear in x or y, duplicates removed). union_all()

Use setequal() to test whether two data sets

contain the exact same rows (in any order)

retains duplicates.

matches from y in a single new

Also includes magrittr: pipelines



• %>% (called pipe) to perform operations sequentially

• Improves the **readability** of the code:

```
data %>% filter(blood_oxygen < 0.9) %>%
  group_by(status, sex) %>%
  summarise(mean_ox = mean(blood_oxygen))
```



ggplot2

purrr



Massimo Bourquin October 2022

Multivariate statistics in R



Package #2: ggplot2

Plotting library

Cheatsheet:

https://github.com/rstudio/cheatsheets/blob/main/datavisualization.pdf

Better visualisations / flexibility

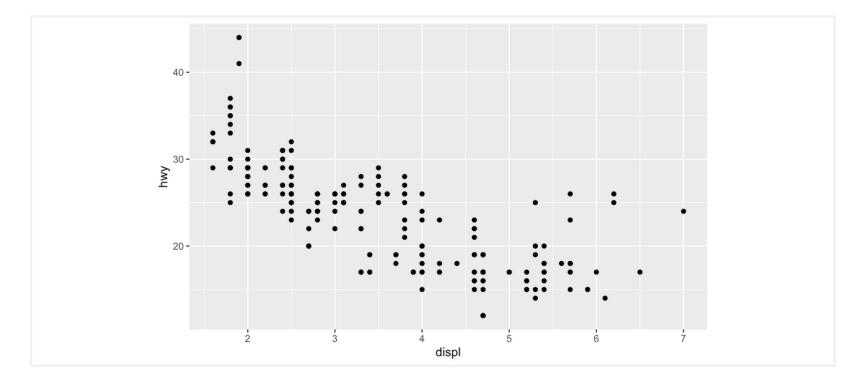
Less code / easier to understand



Basic example

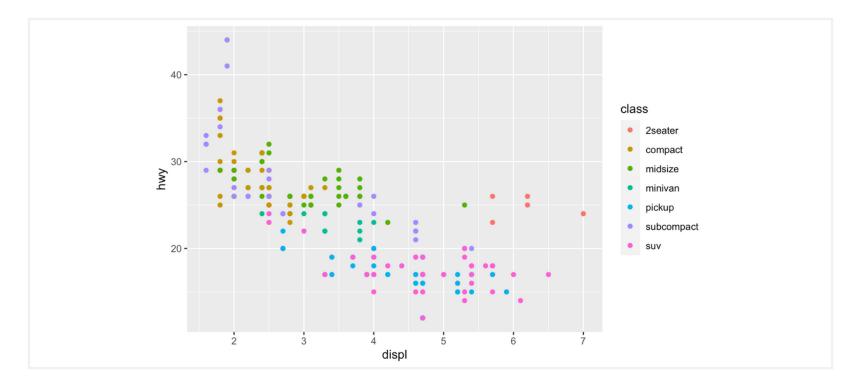
*	manufacturer ‡	model ‡	displ ‡	year ‡	cyl ‡	trans ‡	drv ‡	cty ‡	hwy ‡	fl ‡	class ‡
1	audi	a4	1.8	1999	4	auto(l5)	f	18	29	р	compact
2	audi	a4	1.8	1999	4	manual(m5)	f	21	29	р	compact
3	audi	a4	2.0	2008	4	manual(m6)	f	20	31	р	compact
4	audi	a4	2.0	2008	4	auto(av)	f	21	30	р	compact
5	audi	a4	2.8	1999	6	auto(l5)	f	16	26	р	compact
6	audi	a4	2.8	1999	6	manual(m5)	f	18	26	р	compact
7	audi	a4	3.1	2008	6	auto(av)	f	18	27	р	compact
8	audi	a4 quattro	1.8	1999	4	manual(m5)	4	18	26	р	compact
9	audi	a4 quattro	1.8	1999	4	auto(l5)	4	16	25	р	compact
10	audi	a4 quattro	2.0	2008	4	manual(m6)	4	20	28	р	compact

```
ggplot(data = mpg) +
  geom_point(mapping = aes(x = displ, y = hwy))
Copy
```



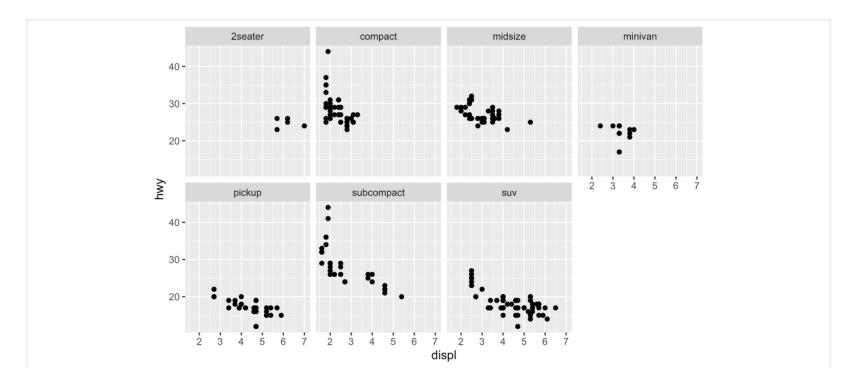
Customisation

```
ggplot(data = mpg) +
  geom_point(mapping = aes(x = displ, y = hwy, color = class))
```



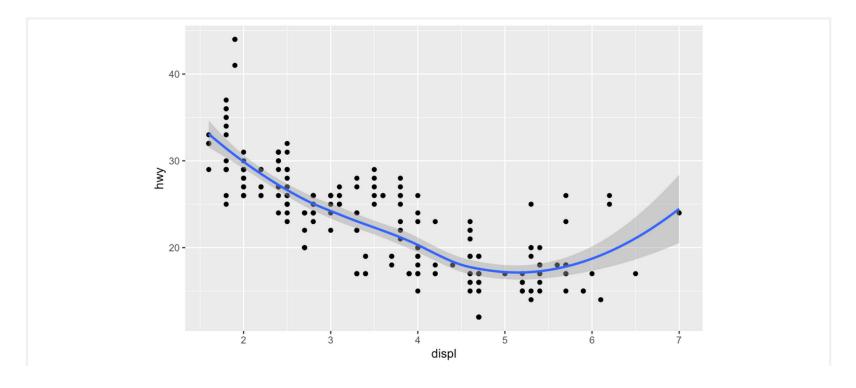
Facet_grid / wrap

```
ggplot(data = mpg) +
  geom_point(mapping = aes(x = displ, y = hwy)) +
  facet_wrap(~ class, nrow = 2)
```



Plotting models

```
ggplot(data = mpg) +
  geom_point(mapping = aes(x = displ, y = hwy)) +
  geom_smooth(mapping = aes(x = displ, y = hwy))
```

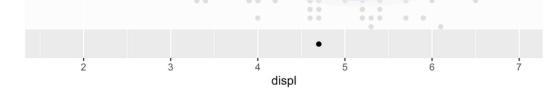


Plotting models

```
ggplot(data = mpg) +
                                                                       Copy
  geom\_point(mapping = aes(x = displ, y = hwy)) +
  geom\_smooth(mapping = aes(x = displ, y = hwy))
```

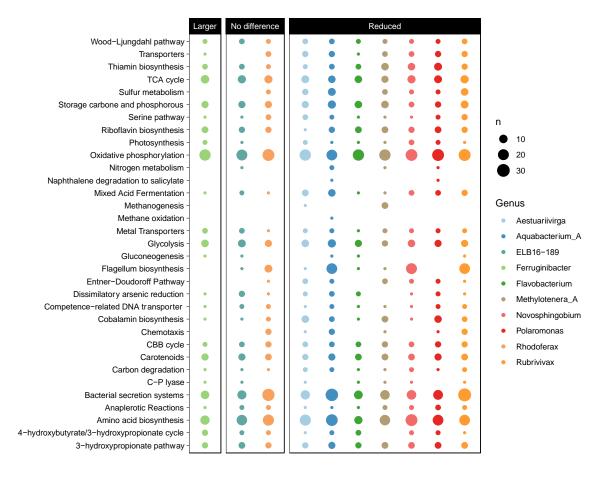


- → Linear models 'lm'
- → Generalised additive models 'gam'
- → Moving averages 'loess'
- → Generalised linear models 'glm'
- →Etc.



A more complicated example

```
ggplot(core_funcs, aes(x=Genus, y=Category)) + geom_count(aes(colour=Genus)) + xlab('') + ylab('') +
    scale_colour_manual(values = cols) + theme_linedraw() +
    theme(axis.text.x = element_blank(), axis.ticks.x = element_blank(), panel.grid.major = element_blank(), panel.grid.minor = element_blank()) +
    facet_grid(.~Genus_category, scales = 'free', space = 'free')
```





ggplot2



Massimo Bourquin October 2022

Multivariate statistics in R



ggplot2

purrr



Massimo Bourquin October 2022

Publication.

Multivariate statistics in R



Package #3: ggpubr

Publication ready plots

ggplot2 but easier

• ggplot2 needs some formatting before publication

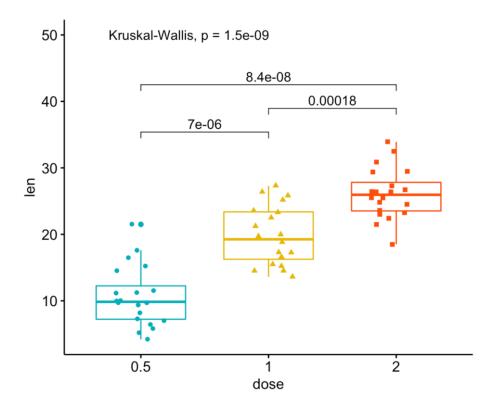
• Syntax can be quite opaque and difficult to assimilate

• ggpubr provides easy-to-use equivalents of ggplot2 functions

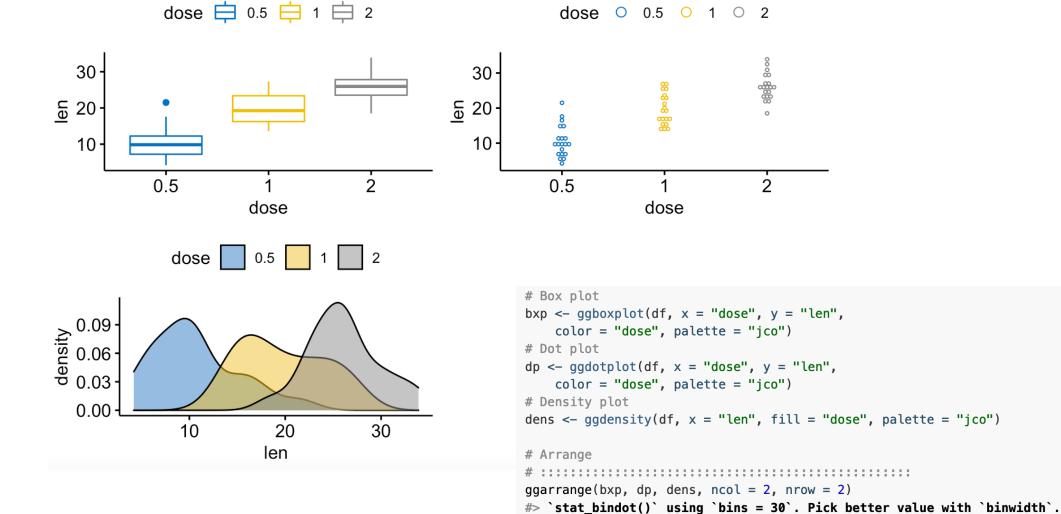
Adding statistics on plots

```
# Add p-values comparing groups
# Specify the comparisons you want
my_comparisons <- list( c("0.5", "1"), c("1", "2"), c("0.5", "2") )
p + stat_compare_means(comparisons = my_comparisons)+ # Add pairwise comparisons p-value
    stat_compare_means(label.y = 50)  # Add global p-value</pre>
```



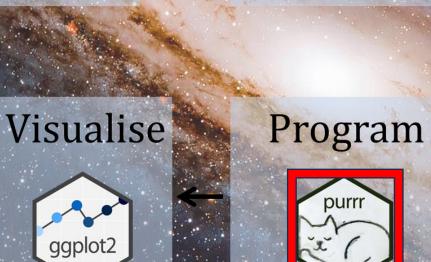


Arranging figures panels



Other useful packages







Massimo Bourquin October 2022

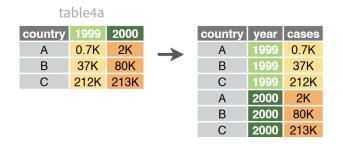
Multivariate statistics in F

Tidyr





Reshape Data - Pivot data to reorganize values into a new layout.



pivot_longer(data, cols, names_to = "name", values_to = "value", values_drop_na = FALSE)

"Lengthen" data by collapsing several columns into two. Column names move to a new names_to column and values to a new values_to column.

pivot longer(table4a, cols = 2:3, names to ="year", values to = "cases")



country	year	type	count		country	year	cases	pop
Α	1999	cases	0.7K	_	Α	1999	0.7K	19M
Α	1999	pop	19M	→	Α	2000	2K	20M
Α	2000	cases	2K		В	1999	37K	172M
Α	2000	рор	20M		В	2000	80K	174M
В	1999	cases	37K		С	1999	212K	1T
В	1999	рор	172M		С	2000	213K	1T
В	2000	cases	80K					
В	2000	рор	174M					
С	1999	cases	212K					
С	1999	рор	1T					
С	2000	cases	213K					

pivot_wider(data, names_from = "name", values from = "value")

The inverse of pivot_longer(). "Widen" data by expanding two columns into several. One column provides the new column names, the other the values.

pivot_wider(table2, names_from = type, values_from = count)

Purrr

Cheatsheet: https://github.com/rstudio/cheatsheets/blob/main/purrr.pdf



Never copy and paste more than twice!

You can use functions or loops

- Purr does the same but with clarity with the map family
- Check https://r4ds.had.co.nz/iteration.html for a nice tutorial

Purrr



Cheatsheet: https://github.com/rstudio/cheatsheets/blob/main/purrr.pdf

```
df <- tibble(
    a = rnorm(10),
    b = rnorm(10),
    c = rnorm(10),
    d = rnorm(10)
)</pre>
```

```
median(df$a)
#> [1] -0.2457625
median(df$b)
#> [1] -0.2873072
median(df$c)
#> [1] -0.05669771
median(df$d)
#> [1] 0.1442633
```

OR

```
output <- vector("double", ncol(df)) # 1. output
for (i in seq_along(df)) { # 2. sequence
  output[[i]] <- median(df[[i]]) # 3. body
}
output
#> [1] -0.24576245 -0.28730721 -0.05669771 0.14426335
```

0

R

```
map_dbl(df, mean)
#> a b c d
#> -0.3260369 0.1356639 0.4291403 -0.2498034
map_dbl(df, median)
#> a b c d
#> -0.51850298 0.02779864 0.17295591 -0.61163819
```

Ask the internet



 Coding is all about searching on the internet, don't be shy to search and copy some code (always copy the link too, for credits and reference) !!!

• The more you work with a package/language, the better you will know the associated **keywords** to search.

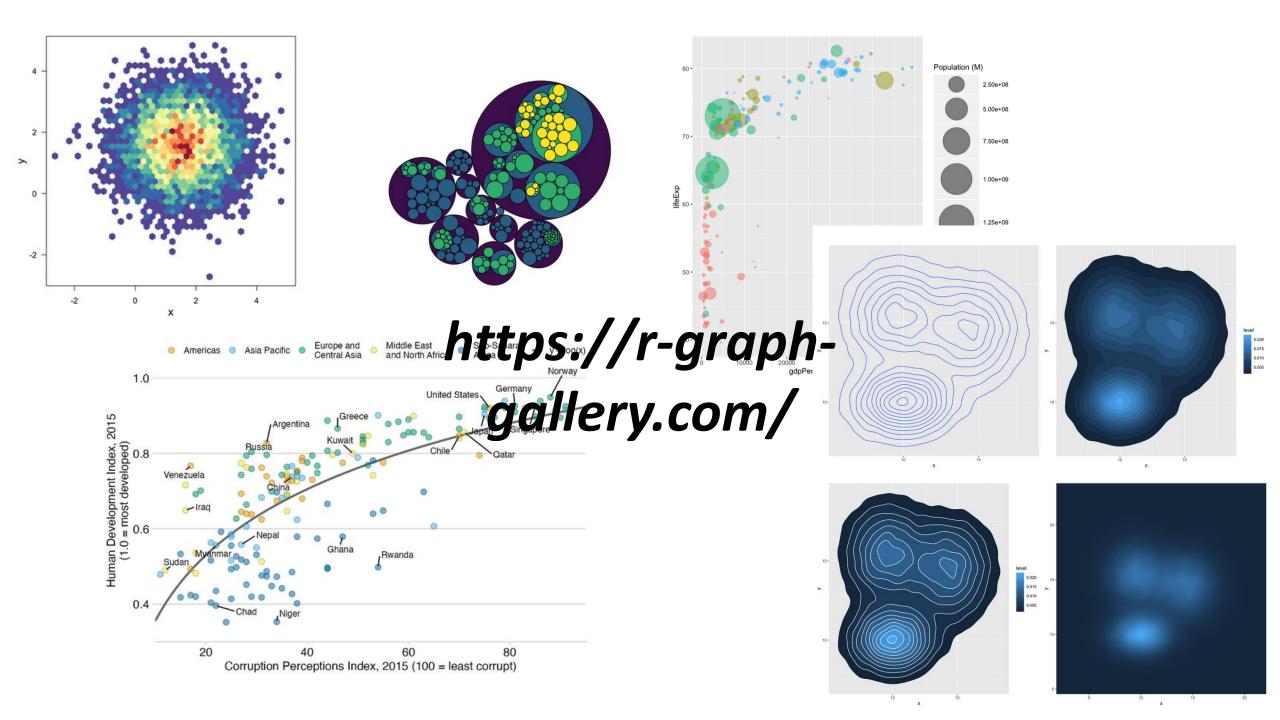
Resources

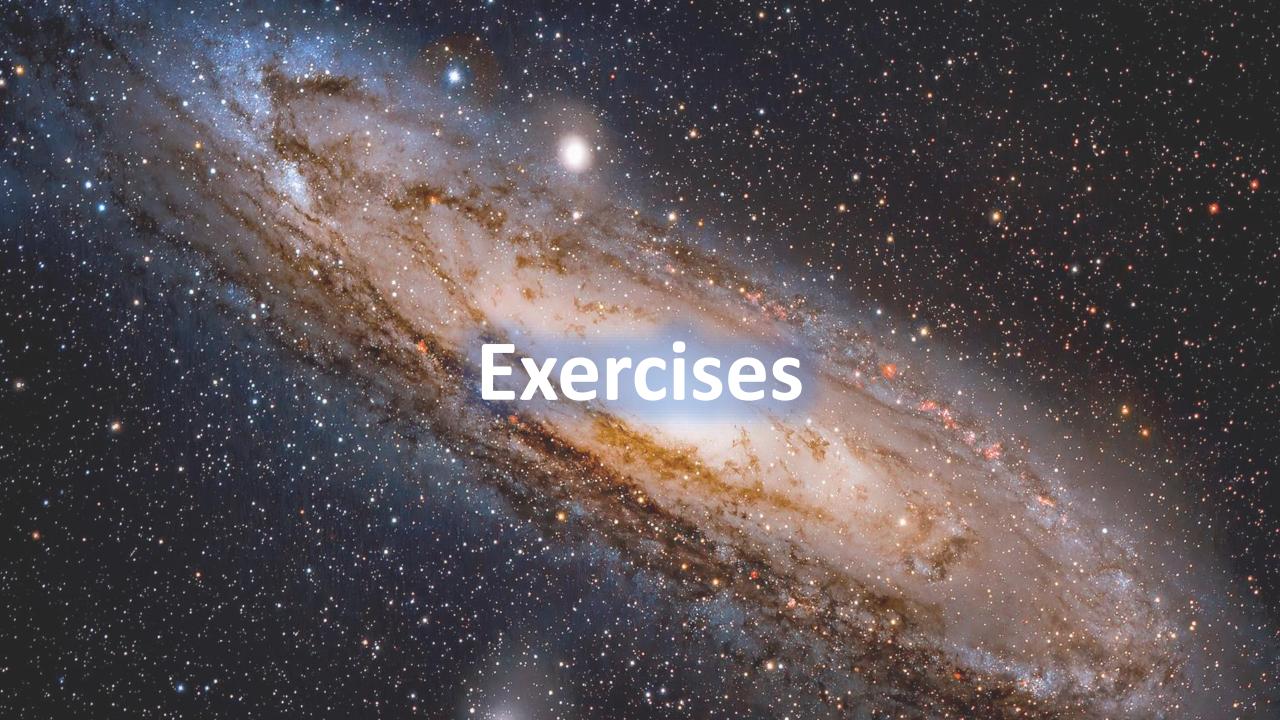
2nd edition of "R for Data Science" (https://r4ds.hadley.nz/):

- Data visualisation
- Data transformation

Other resources:

- https://www.tidyverse.org/
- https://cran.rproject.org/web/packages/magrittr/vignettes/magrittr.html
- https://rpkgs.datanovia.com/ggpubr/





- 1. Load the iris data and do a scatter plot of the sepal.width against the sepal.length with a different colour for each species.
- 2. Add a linear model for each species.
- 3. Compare this plot with another one showing the sepal.length/sepal.width ratio to the petal.length/petal.width ratio. Place them side by side.
- **4. Summarise** the sepal.width (mean, sd) for each species, and plot a **bar** chart with **error bars** (=sd).
- 5. Try to reproduce the figure on the right.

