

E0410 Fundamentals of Statistics for Scientific Data Using R

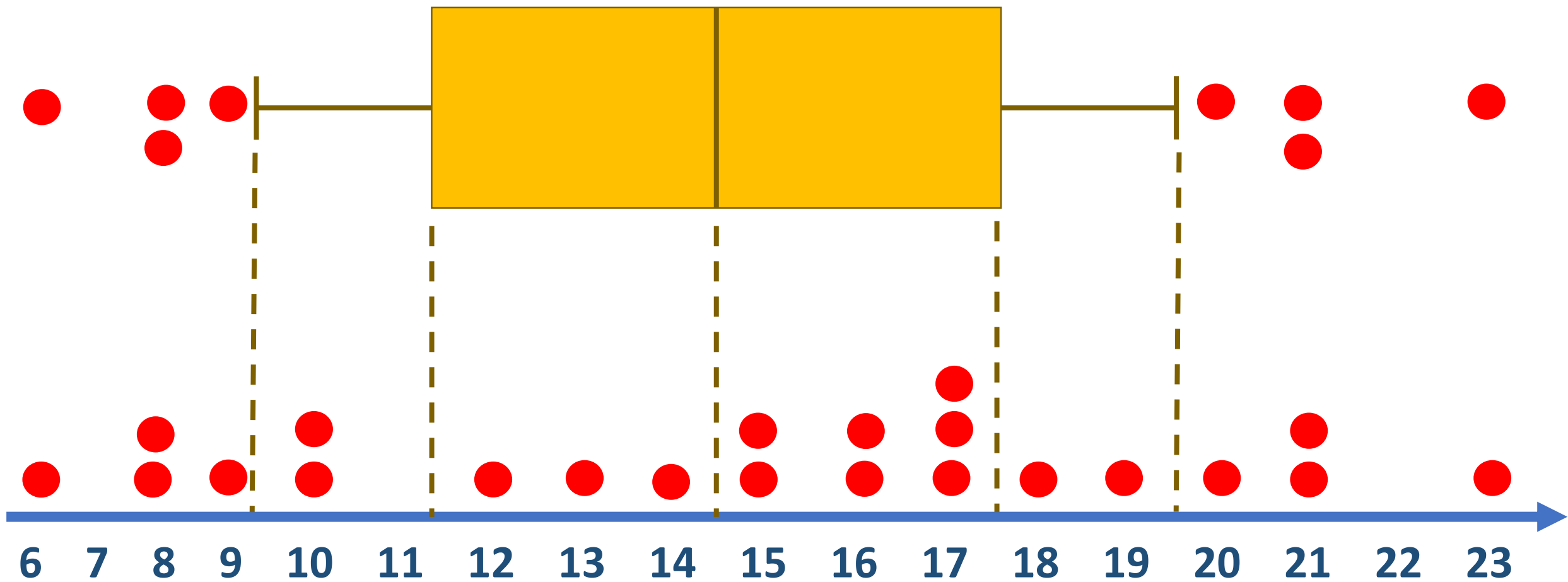
by Daria Sapunova, PhD student, RECETOX

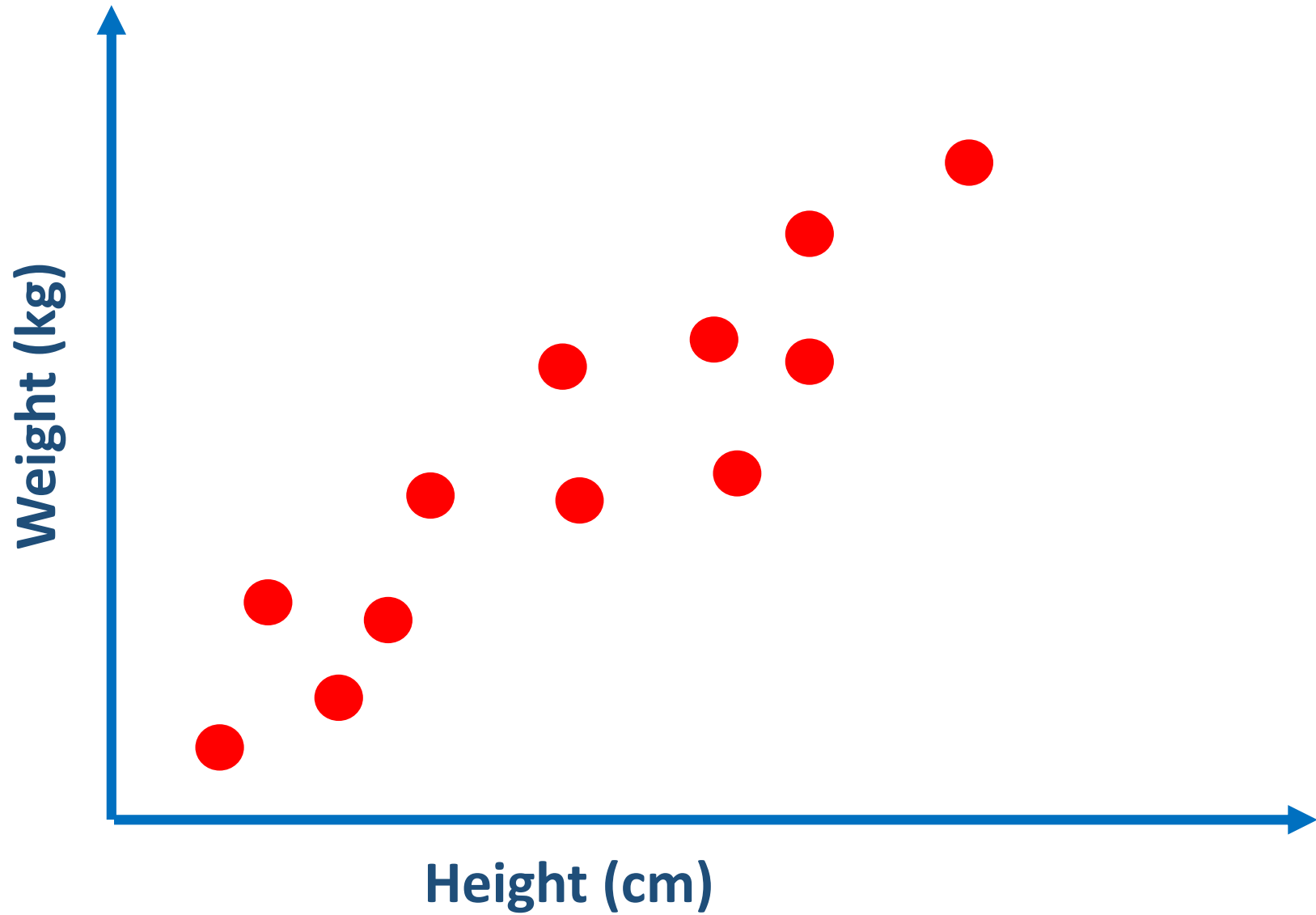
daria.sapunova@recetox.muni.cz

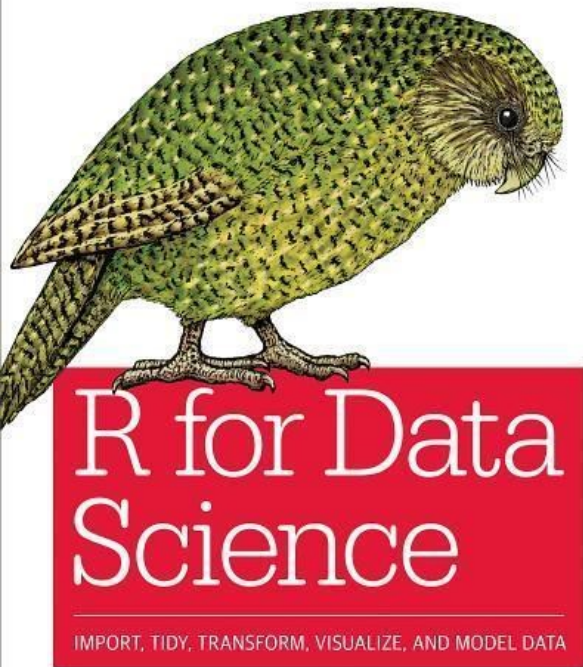
Bohunice, D29, room 123

Theoretical part

The data is only for
visual demonstration!







Hadley Wickham &
Garrett Grolemund

R for Data Science

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

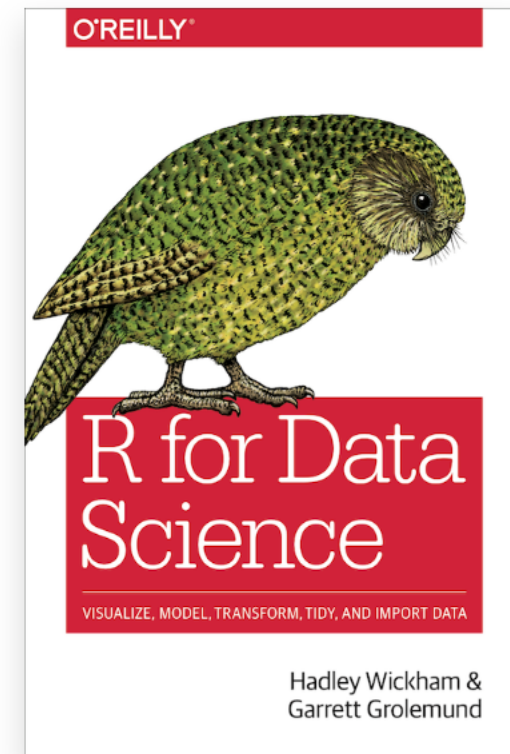
12 Tidy data

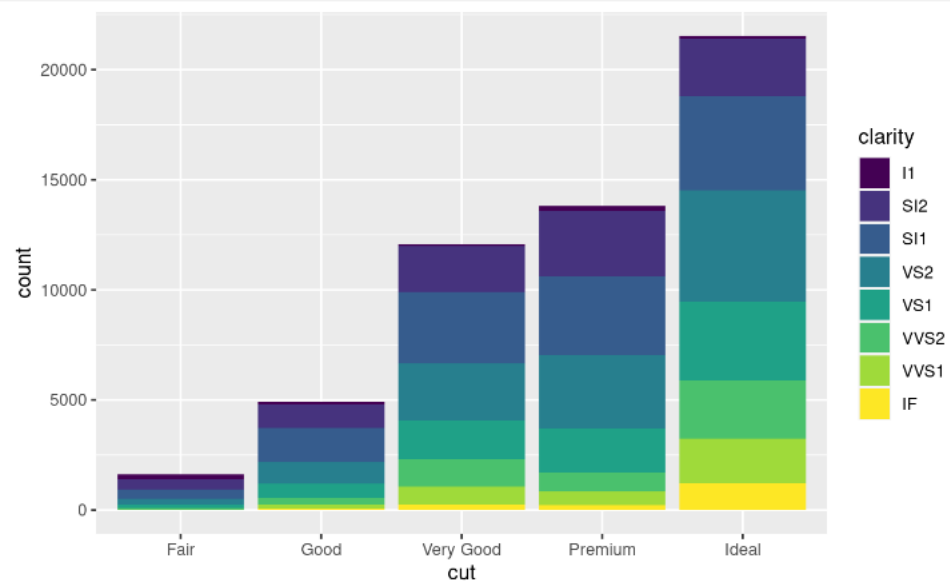
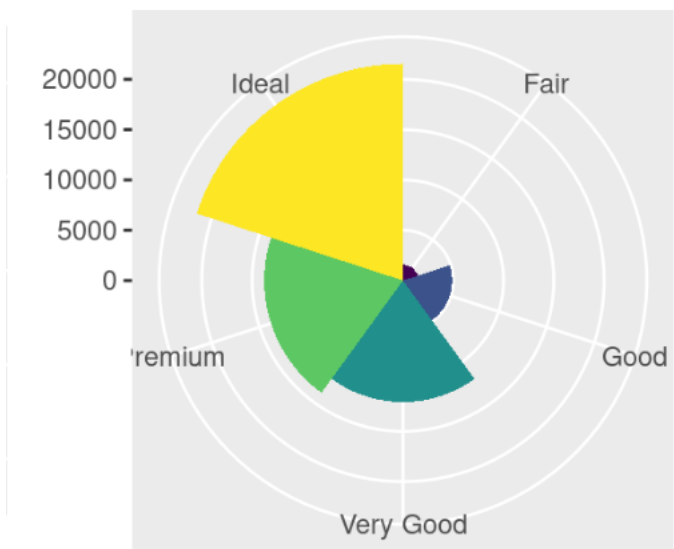
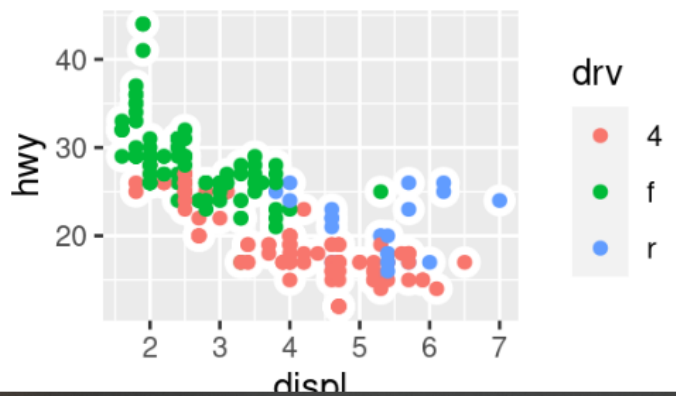
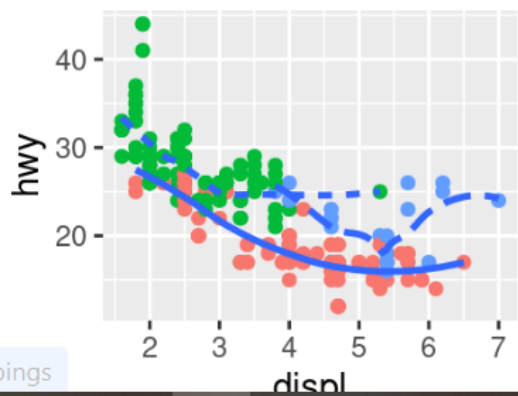
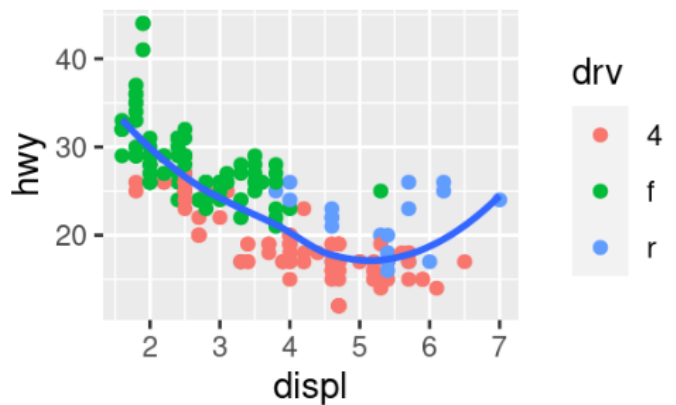
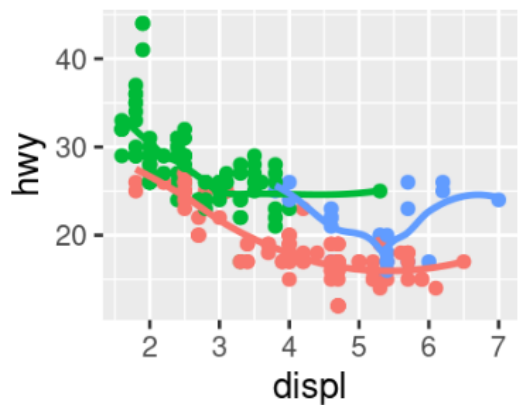
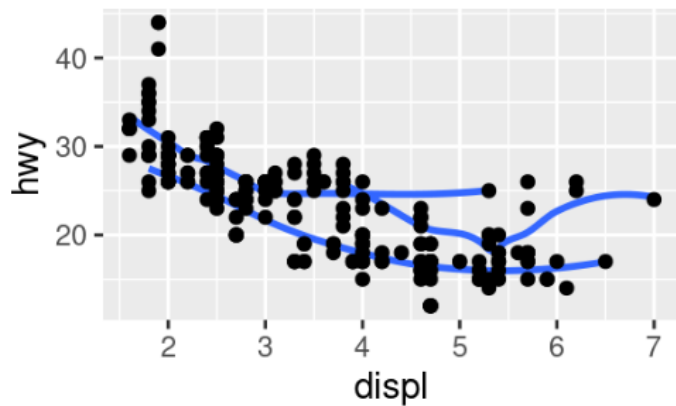
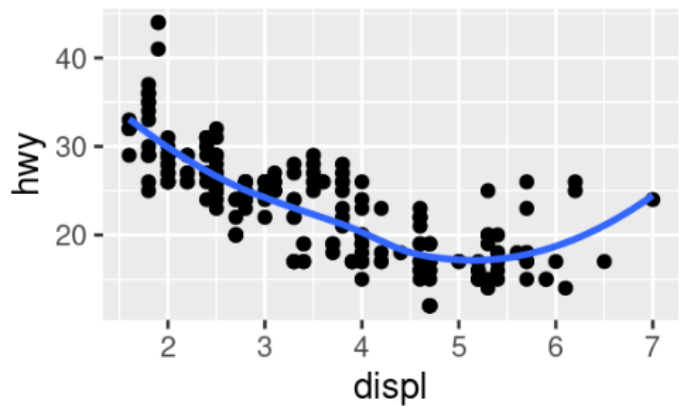
13 Relational data

Welcome

This is the website for “**R for Data Science**”. This book will teach you how to do data science with R: You’ll learn how to get your data into R, get it into the most useful structure, transform it, visualise it and model it. In this book, you will find a practicum of skills for data science. Just as a chemist learns how to clean test tubes and stock a lab, you’ll learn how to clean data and draw plots—and many other things besides. These are the skills that allow data science to happen, and here you will find the best practices for doing each of these things with R. You’ll learn how to use the grammar of graphics, literate programming, and reproducible research to save time. You’ll also learn how to manage cognitive resources to facilitate discoveries when wrangling, visualising, and exploring data.

This website is (and will always be) **free to use**, and is licensed under the Creative








cole nussbaumer knaflic


storytelling with data

a data
visualization
guide for
business
professionals

WILEY



	A	B	C
1	15%	22%	42%
2	40%	36%	20%
3	35%	17%	34%
4	30%	29%	28%
5	55%	30%	58%
6	11%	25%	49%



	A	B	C
Category 1	15%	22%	42%
Category 2	40%	36%	20%
Category 3	35%	17%	34%
Category 4	30%	29%	28%
Category 5	55%	30%	58%
Category 6	11%	25%	49%

010%




cole nussbaumer knaflic

storytelling
with data *let's*
PRACTICE!



	A	B	C
Y1	15%	22%	42%
Y2	40%	36%	20%
Y3	35%	17%	34%
Y4	30%	29%	28%



	A	B	C
CATEGORY 1	15%	22%	42%
CATEGORY 2	40%	36%	20%
CATEGORY 3	35%	17%	34%
CATEGORY 4	30%	29%	28%

010%

WILEY



*Real world
data*

iris &
mtcars

How your data should look before importing into R

Food c	Food - description	Chemical	units	n	Mean	Stdev	Food origin	Food origin sp	Production	Date be	Date er
A031G	Eggs	PFOS	ng/g fw	3	52.8	30.9	Belgium	Gavere, Belgium	Home raised eg	2011	2011
A031G	Eggs	PFOS	ng/g fw	3	5.3	3.5	Belgium	Mechelen, Belgium	Home raised eg	2011	2011
A031G	Eggs	PFOS	ng/g fw	3	7	4.7	Belgium	Grimbergen, Belgiu	Home raised eg	2011	2011
A031G	Eggs	PFOS	ng/g fw	3	386.1	307.9	Belgium	Zwijndrecht, Belgiu	Home raised eg	2011	2011
A031G	Eggs	PFOS	ng/g fw	3	2.5	1	Belgium	Kessel, Belgium	Home raised eg	2011	2011
A031G	Eggs	PFOS	ng/g fw	3	0.4	0.5	Belgium	Westmalle, Belgiu	Home raised eg	2011	2011
A031G	Eggs	PFOS	ng/g fw	3	4.5	5.3	Belgium	Nijlen, Belgium	Home raised eg	2011	2011
A031G	Eggs	PFOS	ng/g fw	3	13	9	Belgium	Hoeneve, Belgium	Home raised eg	2011	2011
A031G	Eggs	PFOS	ng/g fw	3	9.2	4.6	Belgium	Arendonk, Belgium	Home raised eg	2011	2011
A031G	Eggs	PFOS	ng/g fw	3	3.1	1.8	Belgium	Olmen, Belgium	Home raised eg	2011	2011
A031G	Eggs	PFOS	ng/g fw	3	3.3	1.7	Belgium	Olmen, Belgium	Home raised eg	2011	2011
A031G	Eggs	PFOS	ng/g fw	3	3.4	1.8	Belgium	Lille, Belgium	Home raised eg	2011	2011
A031G	Eggs	PFOS	ng/g fw	3	109.9	59.2	Belgium	Zwijndrecht, Belgiu	Home raised eg	2011	2011
A031G	Eggs	PFOS	ng/g fw	3	7.8	5.9	Belgium	Edegem, Belgium	Home raised eg	2011	2011
A031G	Eggs	PFOS	ng/g fw	3	3	2.9	Belgium	Liedekerke, Belgiu	Home raised eg	2011	2011
A031G	Eggs	PFOS	ng/g fw	3	0.4	0.7	Belgium	Meerhout, Belgium	Home raised eg	2011	2011
A031G	Eggs	PFOS	ng/g fw	3	4.7	3.1	Belgium	Oelegem, Belgium	Home raised eg	2011	2011

How your data should look before importing into R

ID	Age	Gender	Height	Weight	Education	Freq.of fish consumption (times/month)	Source of fish (majority)	Chronical disease	Last med.exam. (years ago)
Person 1	45	Male	170.4	76.4	PhD	21	Market	Diabetes	10
Person 2	26	Male	168.3	65.3	BS	5	Self-fishing	NA	NA
Person 3	54	Female	168.6	75.3	MS	23	Grocery store	Bladder infection	10
Person 4	65	Male	156.6	44.2	MS	10	Market	Diabetes	5
Person 5	21	Female	170.1	69.9	HS	8	Market	No	NA
Person 6	44	Male	176.4	84.3	MS	43	Self-fishing	No	2
...									
Person 400	6	Male	121.2	21.9	NA	15	Market	NA	1

Data preparation – cleaning, harmonizing and structuring

- ✓ Upload your data.
- ✓ Choose variables you are going to work with.
- ✓ Inspect the variables (do descriptive statistics, check unique values, check data type, check amount of missing data).
- ✓ Harmonize your variables (esp. categorical variables).
- ✓ Structure your data.

Preliminary analysis

- ✓ Build histograms in case of numerical variables.
- ✓ Build box plots in case of categorical variables.
- ✓ Build scatterplots to see associations between numerical variables.
- ✓ Check data distributions in case of numerical variables (normality tests).
- ✓ Transform the data if necessary.

You may proceed to the analysis

Task!



Data preparation – cleaning, harmonizing and structuring

✓ Structure your data.

✓ Upload your data.

✓ Choose variables you are going to work with.

✓ Inspect the variables (do descriptive statistics, check unique values, check data type, check amount of missing data).

✓ Harmonize your variables (esp. categorical variables).

Preliminary analysis

- ✓ Build histograms in case of numerical variables.
 - ✓ Build box plots in case of categorical variables.
 - ✓ Build scatterplots to see associations between numerical variables.
-
- ✓ Check data distributions in case of numerical variables (normality tests).
-
- ✓ Transform the data if necessary.

You may proceed to the analysis

Practical part