Data Storytelling Dashboard for Exploring Auckland Air Quality



ZHAOMING SU

DEPARTMENT OF STATISTICS THE UNIVERSITY OF AUCKLAND

supervised by

DR. EARO WANG

PROF. CHRIS WILD

A dissertation submitted in partial fulfilment of the requirements for the degree of Bachelor of Science (Honours) in Statistics, The University of Auckland, 2021.

Contents

A۱	bstract	v
A	cknowledgements	vii
Co	opyright notice	ix
D	eclaration	xi
1	Introduction (Template Demo)	xiii
	1.1 Rmarkdown	. xiii
	1.2 Data	. xiv
	1.3 Figures	. xiv
	1.4 Results from analyses	. xv
	1.5 Tables	. xv
2	Background and related works	xvii
	2.1 Tidy time series data-wrangling toolbox	. xvii
	2.2 Time series graphics toolbox	. xviii
	2.3 HTML widgets for interactive graphics	. xxi
	2.4 Visual analysis for air quality data	. xxiii
3	Auckland air quality data	xxv
	3.1 Introduction	. xxv
	3.2 Data quality and cleaning	. xxvi
	3.3 Data enrichment	.xxviii
4	Design layout and philosophy	xxix
5	Linked interactive graphics	xxxi
6	Modelling	xxxiii
7	Conclusion and future works	xxxv
A	Additional stuff	xxxvii
Bi	bliography	xxxix

Abstract

The abstract should outline the main approach and findings of the thesis and must not be more than 500 words.

Acknowledgements

I would like to thank my pet goldfish for ...

Copyright notice

© ZHAOMING SU (2021).

I certify that I have made all reasonable efforts to secure copyright permissions for third-party content included in this thesis and have not knowingly added copyright content to my work without the owner's permission.

Declaration

This dissertation is an original work of my research and contains no material which has been accepted for the award of any other degree or diploma at any university or equivalent institution and that, to the best of my knowledge and belief, this dissertation contains no material previously published or written by another person, except where due reference is made in the text of the dissertation.

Introduction (Template Demo)

This is where you introduce the main ideas of your thesis, and an overview of the context and background.

In a PhD, Chapter 2 would normally contain a literature review. Typically, Chapters 3–5 would contain your own contributions. Think of each of these as potential papers to be submitted to journals. Finally, Chapter 6 provides some concluding remarks, discussion, ideas for future research, and so on. Appendixes can contain additional material that don't fit into any chapters, but that you want to put on record. For example, additional tables, output, etc.

1.1 Rmarkdown

In this template, the rest of the chapter shows how to use Rmarkdown. The big advantage of using Rmarkdown is that it allows you to include your R code directly into your thesis, to ensure there are no errors in copying and pasting, and that everything is reproducible. It also helps you stay better organized.

For details on using R Markdown see http://rmarkdown.rstudio.com.

1.2 Data

Included in this template is a file called sales.csv. This contains quarterly data on Sales and Advertising budget for a small company over the period 1981–2005. It also contains the GDP (gross domestic product) over the same period. All series have been adjusted for inflation. We can load in this data set using the following command:

```
sales <- ts(read.csv("data/sales.csv")[, -1], start = 1981, frequency = 4)</pre>
```

Any data you use in your thesis can go into the data directory. The data should be in exactly the format you obtained it. Do no editing or manipulation of the data outside of R. Any data munging should be scripted in R and form part of your thesis files (possibly hidden in the output).

1.3 Figures

Figure 1.1 shows time plots of the data we just loaded. Notice how figure captions and references work. Chunk names can be used as figure labels with fig: prefixed. Never manually type figure numbers, as they can change when you add or delete figures. This way, the figure numbering is always correct.

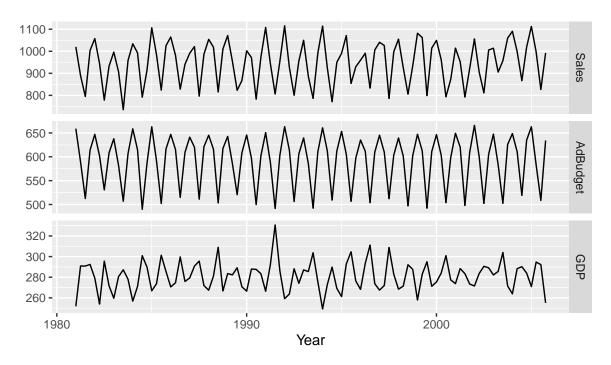


Figure 1.1: *Quarterly sales, advertising and GDP data.*

1.4 Results from analyses

We can fit a dynamic regression model to the sales data.

If y_t denotes the sales in quarter t, x_t denotes the corresponding advertising budget and z_t denotes the GDP, then the resulting model is:

$$y_t - y_{t-4} = \beta(x_t - x_{t-4}) + \gamma(z_t - z_{t-4}) + \theta_1 \varepsilon_{t-1} + \Theta_1 \varepsilon_{t-4} + \varepsilon_t$$
 (1.1)

where $\beta = 2.28$, $\gamma = 0.97$, $\theta_1 = NA$, and $\Theta_1 = -0.90$.

1.5 Tables

Let's assume future advertising spend and GDP are at the current levels. Then forecasts for the next year are given in Table 1.1.

Again, notice the use of labels and references to automatically generate Table numbers. In this case, we need to generate the label ourselves.

Point Forecast	Lo 80	Hi 80	Lo 95	Hi 95
1000.2	947.7	1052.7	919.9	1080.5
1013.1	959.3	1066.8	930.9	1095.3
1076.7	1022.9	1130.6	994.4	1159.0
1003.5	949.7	1057.4	921.2	1085.8

Table 1.1: Forecasts for the next year assuming Advertising budget and GDP are unchanged.

The knitLatex package is useful for generating tables from R output. Other packages can do similar things including the kable function in knitr which is somewhat simpler but you have less control over the result. If you use knitLatex to generate tables, don't forget to include results="asis" in the chunk settings.

Background and related works

2.1 Tidy time series data-wrangling toolbox

The **tsibble** package offers a data infrastructure for wrangling time series data (Wang, Cook, and Hyndman, 2020). A time series data set consists of one or more sequences indexed by time, often with a regular interval. As such, data-wrangling processes of time series data need to account for the special requirements of time series data analysis, including the explicit identification of time gaps and a method of handling multiple time series in a single data set for identifying duplicate records.

Analyses of fixed-interval time series require the data to be free from missing value, especially when the series is self-dependent. Whilst the explicit missing values can be easily handled by the substitution with interpolated values, the implicit gaps with missing index values are often neglected. In the case of multiple time series, locations of implicit time gaps may be different in each sequence; filling the gaps with traditional loops can be time-consuming and inefficient. **tsibble** identifies implicit time gaps with the *index* and *key* variables, such that each variable in the tsibble object is uniquely identified by the index and the interaction of all keys. As such, each time series is uniquely identified by the keys, allowing efficient identification of implicit time gaps, which is achieved by **tsibble** with a range of wrangling verbs.

Duplicates exist in different forms in cross-sectional and time series data. Typically, duplicates are identical observations exhibited as rows in a data frame, yet such definition is inadequate in identifying duplicates in time series data. There exists only one true value at any given point in time for each time series, meaning that there may be duplicate values that are non-identical observations with identical key-index pairs yet different in values. Instead of searching merely for duplicate rows, **tsibble** checks for duplicate key-index pairs. To avoid negligence, the creation of tsibble will fail upon detected duplicates.

2.2 Time series graphics toolbox

2.2.1 Calendar graphics

Calendars are the systematic partition of time from the observed solar-lunar phenomena and cultural custom, which is usable as graphics for temporal representations of societal activities and natural events. Calendar graphics are the method for the aggregated visualisation of time series data at sub-daily intervals, depicting the temporal dimension of time series data as the spatial layout in the calendar grid. The motivation of utilising calendars for data visualisations arises from the convenience of displaying observations in association with exact dates.

Air quality data are conventionally collected at hourly intervals, from which a time series plot becomes overcrowded, impeding the visual detections of abnormalities. Rahman and Lee (2020) depicted a method of non-cartesian heatmaps on a calendar coordinate for visualising air pollution data. Prior to the paper, Liu, Li, and Li (2016) used calendar heatmaps to analyse the correlation of particulate matter temporally. Calendar graphics highlight the abnormalities and allow the association of the detected abnormalities to events with dates, providing insights and directions for analyses.

Calendar graphics is an application of trellis displays (Becker, Cleveland, and Shyu, 1996), which spatially modularise the temporal dimension into conditional groups of small multiples (Tufte, 1983) using calendar period (i.e., month and year) as an integrated and aligned plot. The expanded layout of the temporal dimension on a plane eases the

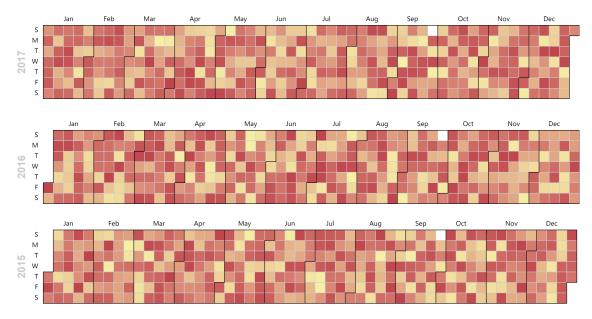


Figure 2.1: A demonstration of a basic calendar heatmap with **Echarts.js** (The Apache Software Foundation, 2020) faceted (unaligned) by **year** ~ **month**. Each tile corresponds to the value of an observed day or aggregated sub-daily observations, whose exact date and day of the week are easily identifiable.

cognitive load (Tufte, 1983) in temporally locating the date of the events and extracting the date components, contrary to conventional time series plots.

2.2.2 Time series plots

In most scenarios, visualising time series relies on connecting points of observations with lines, curves or splines (Wilke, 2019), such that the sequences are plotted against the time index by positions along common xy-scales on the Cartesian coordinate (Hyndman, 2021a). Connected time plots are the most elementary visual method for spotting extrinsic features in time series, including trends, seasons, cycles, clustering and oscillations.

Based on connected time plots, O'Hara-Wild, Hyndman, and Wang (2021) proposed the seasonal plot as a method of visualising seasonal patterns in the **feasts** package. The method conditionally subsets the complete time series into partitions of seasonal periods, each to be plotted in a homogeneous time plot and distinguished using a gradient colour scale for longitudinal comparison between periods.

Nevertheless, interpreting connected time plots relies on the visual alignment of positions to the xy-scales in the Cartesian coordinates. Such visual alignments can be challenging

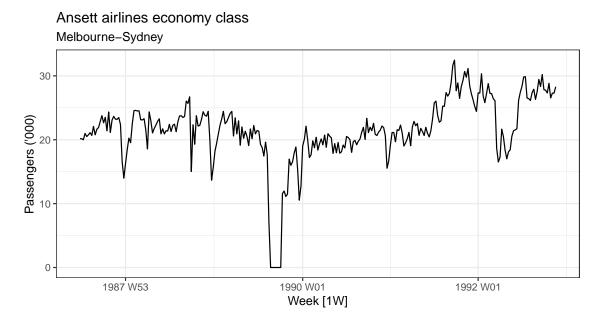


Figure 2.2: A demonstration of a basic line plot with the **ggplot2** package (Wickham, 2016) for the weekly economy passenger load on Ansett Airlines (Hyndman, 2021b). Visual analysis shows weak trend and cycle and abnormal zero values to be investigated. The presence of clustering also indicates a positive temporal dependence in the time series. It is nonetheless uneasy to align any observation to a date accurately.

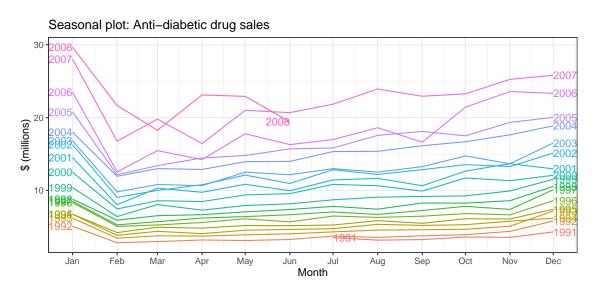


Figure 2.3: A ggplot2-based seasonal time series plot with the feasts package (O'Hara-Wild, Hyndman, and Wang, 2021) for the monthly anti-diabetic drug sales in Australia (Hyndman, 2021b). The plot allows for both visualisation of intra-seasonal patterns and inter-seasonal variations shown as a positive trend in drug sales from year to year.

upon the absence of explicit gridlines and axes, such as when the plot is fitted as a part of trellis displays (Becker, Cleveland, and Shyu, 1996), in which the trend and scale of variations become ambiguous. As such, it is common to fill the area under the curve to emphasise the temporal variation in the plot (Wilke, 2019).

2.3 HTML widgets for interactive graphics

2.3.1 Developing web applications with R

The **shiny** package (Chang et al., 2021) provides a framework for developing web applications with **R** code (R Core Team, 2021), both user and developer-friendly. It enables **R** users with no prior knowledge of HTML, CSS and JavaScript to create custom web applications with sophisticated functionality with template UI components and a server powered with reactive programming.

Reactive programming is the core of computation logic behind **shiny** (Wickham, 2021), greatly simplifies the design of workflow, focusing only on evaluating the changes of values over time. Each change in reactive values is observed as an event by pre-defined callback functions, and workflows are executed as responses to events observed. The lazy nature of reactivity avoids repeated evaluations of expressions leading to wastage in computational resources. Reactive programming also allows users to define abstract workflows without conceiving the low-level data and programming logic by restricting evaluations to merely reactions to events, including user actions and internal value changes. A single reactive value can be observed and called by several callback functions, which can be shared across different functionalities. As the reactive value always keeps the previous evaluated result, conflicts among functionalities are avoided.

The user interface of **shiny** applications provides the front-end inputs and output display from the back end server logic. The collection of user input is the primary source of change of reactive values, events that trigger the evaluation of expressions in the back end server logic. The results of the evaluated expressions are rendered as the outputs, which may be as simple as prints of R objects or as sophisticated as HTML interactive graphics.

2.3.2 Interactive graphics with Echarts JavaScript in R

The echarts4r package (Coene, 2020), the implementation of Echarts JavaScript in **R** (The Apache Software Foundation, 2020), offers access to a powerful open-source JavaScript library for interactive graphics. The creation of Echarts graphics with echarts4r adapts the pipe workflow with dplyr %>% (Wickham et al., 2020; Bache and Wickham, 2014). The wrapper functions of echarts4r enable the configuration of graphic options and parameters without the need of understanding JavaScript syntax.

The creation of echarts4r objects follows the layer-to-layer approach, by which each function either (re-)initiates a chart, adds a single layer of graphical primitive or configures the options. It follows that the number of **echarts4r** function calls needed will be at least the number of graphic primitives to be mapped. In the case of mapping multiple time series as in seasonal plots (Section ref(sec:ts-plot)) with **echarts4r**, loops may be needed.

The underlying data passed to Echarts is treated slightly differently from the tidy data (Wickham and Grolemund, 2016) underpinning **ggplot2**: each row represents an observation, and each column represents a variable. **echarts4r** converts the data frame into a JSON object upon initialisation, with each variable treated as a *serie*. Such a difference implies a different logic for mapping the variables to graphical primitives, especially when mapping unordered categorical data.

It is worth noting that the mapping logic of Echarts in some circumstances is inconsistent with the conventional data organisation of tsibble upon plotting multivariate time series. The latter uniquely maps each observation to a single data point on the plot in a one-to-one relationship between the variables and graphical primitives, yet Echarts maps each cell value to a separate graphical primitive such that each time series is plotted with an x-y pair of two *serie* instead of grouping the observations with tsibble keys. Thus, it is often necessary to "widen" the data frame, breaking the time series into separate columns of variables named by the tsibble keys, achievable with the **tidyr** package (Wickham, 2020).

- 2.3.3 Drill-down in interactive graphics
- 2.3.4 Interactive data tables
- 2.3.5 Interactive maps
- 2.4 Visual analysis for air quality data
 - Wind rose and pollution rose
 - Trend estimate fit plot and ACF

Auckland air quality data

3.1 Introduction

Air quality index (AQI) is a critical indicator of overall air quality by measuring key air pollutant concentrations at a given time. The constitution of AQI consists of ambient air pollutants listed in the National Environmental Standards for Ambient Air Quality which defines the threshold target for calculating AQI (Auckland Regional Council, 2020a). The national standard defines AQI as the maximum ambient air pollutant measurement ratio to the national target as a percentage (Auckland Regional Council, 2020b). Over 10 stations in Auckland monitor a subset of the standard-listed pollutants in an hourly interval.

It is of interest to explore the variation and relationships of AQI and its constituent pollutants with other environmental and meteorological parameters over time. The data, provided by Auckland Regional Council (2021), includes 14 parameters from 10 monitoring stations in Auckland from as North as Takapuna to as South as Patumahoe. Available parameters consist of air quality index (AQI), 10 pollutant levels and four other meteorological variables as per Table 3.1, with various starting dates (since as early as 2003 in Takapuna) until April 2021.

Only six of the standard-listed air pollutants are monitored and available in the data, and each station independently monitors a subset of the six pollutants. As such, the calculation of AQI, based on available data, may be simplified to

$$AQI = 100 \times \max\{\frac{PM_{2.5}}{25}, \frac{PM_{10}}{50}, \frac{NO_2}{200}, \frac{SO_2}{350}, \frac{CO}{10}, \frac{O_3}{150}\}$$
(3.1)

Parameter	Unit	Note
AQI		Air quality index
BC(370)	ngm ⁻³	Black carbon at 370nm wavelength
BC(880)	ngm ⁻³	Black carbon at 880nm wavelength
CO*	mgm ⁻³	Carbon monoxide concentration
NO	μgm ⁻³	Nitrogen monoxide concentration
NO_2^*	μgm ⁻³	Nitrogen dioxide concentration
NOx	μgm ⁻³	Nitrogen oxides concentration
O_3^*	μgm ⁻³	Ozone concentration
PM2.5*	μgm ⁻³	Particulate matter with diameter <2.5µm
$PM10^*$	μgm ⁻³	Particulate matter with diameter <10µm
SO_2^*	μgm ⁻³	Sulphur dioxide concentration
Relative Humidity	%	-
Temperature	$^{\circ}C$	
Wind Speed	ms ⁻¹	
Wind Direction	0	

Table 3.1: *Parameters available in raw data.** *AQI-related ambient air pollutants*

It is noteworthy that the availability of air quality parameters in each monitoring station varies from year to year. Besides, the extreme values addressed in Section 3.2.1 are more frequent in earlier years. The final data set is subsetted from the year 2016.

3.2 Data quality and cleaning

The raw data consists of two separate data sets, each with a different data structure. Cleaning and manipulation are needed to ensure that the two data sets are consistent in structure and free from error. The raw data sets are individually inspected and cleaned before combination. This section outlines the issues found and methods to address them.

3.2.1 Abnormal and missing values

Abnormal or missing values arise from instrumental or input errors. Upon inspection, 104,332 records were found to have a negative value. Nevertheless, all pollutants are reported in units in the form of mass per unit volume, and other parameters, except for

temperature, are only sensible if positive as of Table 3.1. Therefore, 104,257 records of insensible negative values are removed. Besides, conspicuously anomalous records of AQI are found in data, including consecutive hours of >1,000 AQI in Takapuna and numerous AQI values being inconsistent with Formula 3.1 based on available pollutants in the same data set. The anomalous records are nonetheless kept as-is for further verification.

In addition, preliminary inspection finds that 0.81% of records are explicitly missing. Yet after filling the implicit time gaps in the data, 53.71% of records are implied to be missing.

3.2.2 Date and time

A consistent format in date and time is crucial to the accuracy of temporal data. Observations with inconsistent time format are present in the data, where some are recorded in hh:mm:ss whilst others in hh:mm. The inconsistency in the time format is correctable due to the hourly nature of the data. 0.06% of records with missing time are removed.

The time zone of New Zealand changes by +1 during daylight saving. To avoid duplicated index upon boundaries of daylight saving upon data visualisation, all time-stamps are presented in NZST (UTC+12). On the other hand, the date and time in the cleaned data file are stored as a single variable, with its format in compliance with ISO 8601 (International Organization for Standardization, 2019; Wickham, Hester, and Francois, 2018).

3.2.3 Duplicate records

Temporal data should not present duplicate records. Of the 7,292,038 valid records, 239,374 (3.28%) are duplicate with 120,207 redundant records. Further checking reveals that 230,822 of the duplicates have inconsistent values. However, as the scale of the inconsistency of most duplicate records is reasonably small, the first-appearing records of each duplicate are kept.

3.2.4 Structural difference in raw data sets

The primary data set, which records all parameters except for wind direction, is in long format, with each observation consisting of a single record of one parameter for one station at a given hour. Nevertheless, each observation of the wind direction data set consists

of wind direction records of all stations at a given hour. Each data set is pivoted to the structure such that each observation is uniquely identified by the date-time and station with records of all parameters before combination to ensure structural consistency.

3.3 Data enrichment

- Categorisation of AQI
- Categorisation of wind direction

Design layout and philosophy

- Overview of AQI
 - Spatial
 - Temporal (AQI and its constituent)
 - * Calendar
 - * Drill-down line plot
- Data enrichment
 - Explore relationship AQI with wind speed and direction
 - Meteorological data
- Trend analysis

Linked interactive graphics

- Introduction
- Implementation of interactive linking
- Modularisation of Shiny App

Modelling

Conclusion and future works

Appendix A

Additional stuff

You might put some computer output here, or maybe additional tables.

Note that line 5 must appear before your first appendix. But other appendices can just start like any other chapter.

Bibliography

- Auckland Regional Council (2020a). *Auckland Ambient Air Quality Guidelines*. Accessed: 2021-09-22. Auckland, New Zealand. https://drive.google.com/drive/folders/1D-WsQ3ISmYWjJlSI29aA7rlgYPrz0v6F.
- Auckland Regional Council (2020b). *Auckland Ambient Air Quality Targets*. Accessed: 2021-09-22. Auckland, New Zealand. https://unitaryplan.aucklandcouncil.govt.nz/Images/Auckland%20Unitary%20Plan%20Operative/Chapter%20E%20Auckland-wide/1.%20Natural%20Resources/E14%20Air%20quality.pdf.
- Auckland Regional Council (2021). *The Auckland Environmental Data*. Accessed: 2021-05-25. Auckland, New Zealand. https://environmentauckland.org.nz/Data/DataSet.
- Bache, SM and H Wickham (2014). *magrittr: A Forward-Pipe Operator for R*. R package version 1.5. https://CRAN.R-project.org/package=magrittr.
- Becker, RA, WS Cleveland, and MJ Shyu (1996). The Visual Design and Control of Trellis Display. *Journal of Computational and Graphical Statistics* **5**(2), 123–155.
- Chang, W, J Cheng, J Allaire, C Sievert, B Schloerke, Y Xie, J Allen, J McPherson, A Dipert, and B Borges (2021). *shiny: Web Application Framework for R*. R package version 1.6.0. https://CRAN.R-project.org/package=shiny.
- Coene, J (2020). *echarts4r: Create Interactive Graphs with 'Echarts JavaScript' Version 4*. R package version 0.3.3. https://CRAN.R-project.org/package=echarts4r.
- Hyndman, R (2021a). *Forecasting: Principles and Practice*. 3rd ed. Melbourne, VIC, Australia: Monash University. https://otexts.com/fpp3.
- Hyndman, R (2021b). *fpp3: Data for "Forecasting: Principles and Practice"* (3rd Edition). R package version 0.4.0. https://CRAN.R-project.org/package=fpp3.

- International Organization for Standardization (2019). *ISO 8601: Date and Time Format*. https://www.iso.org/iso-8601-date-and-time-format.html.
- Liu, J, J Li, and W Li (2016). Temporal Patterns in Fine Particulate Matter Time Series in Beijing: A Calendar View. *Scientific Reports* **6**(32221).
- O'Hara-Wild, M, R Hyndman, and E Wang (2021). *feasts: Feature Extraction and Statistics for Time Series*. R package version 0.1.7. https://CRAN.R-project.org/package=feasts.
- R Core Team (2021). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing. Vienna, Austria. https://www.R-project.org/.
- Rahman, NHA and MH Lee (2020). Air Pollutant Index Calendar-Based Graphics for Visualizing Trends Profiling and Analysis. *Sains Malaysiana* **49**(1), 201–209.
- The Apache Software Foundation (2020). *Apache ECharts: An Open Source JavaScript Visualization Library*. https://echarts.apache.org/handbook.
- Tufte, ER (1983). The Visual Display of Quantitative Information. Graphics Press.
- Wang, E, D Cook, and RJ Hyndman (2020). A new tidy data structure to support exploration and modeling of temporal data. *Journal of Computational and Graphical Statistics* **29**(3), 466–478.
- Wickham, H (2016). *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York. https://ggplot2.tidyverse.org.
- Wickham, H (2020). *tidyr: Tidy Messy Data*. R package version 1.1.2. https://CRAN.R-project.org/package=tidyr.
- Wickham, H (2021). *Mastering Shiny*. Sebastopol, CA: O'Reilly Media. https://mastering-shiny.org.
- Wickham, H, R François, L Henry, and K Müller (2020). *dplyr: A Grammar of Data Manipulation*. R package version 1.0.1. https://CRAN.R-project.org/package=dplyr.
- Wickham, H and G Grolemund (2016). *R for Data Science*. Sebastopol, CA: O'Reilly Media. https://r4ds.had.co.nz.
- Wickham, H, J Hester, and R Francois (2018). *readr: Read Rectangular Text Data*. R package version 1.3.1. https://CRAN.R-project.org/package=readr.
- Wilke, CO (2019). Fundamentals of Data Visualization. Sebastopol, CA: O'Reilly Media. https://clauswilke.com/dataviz.