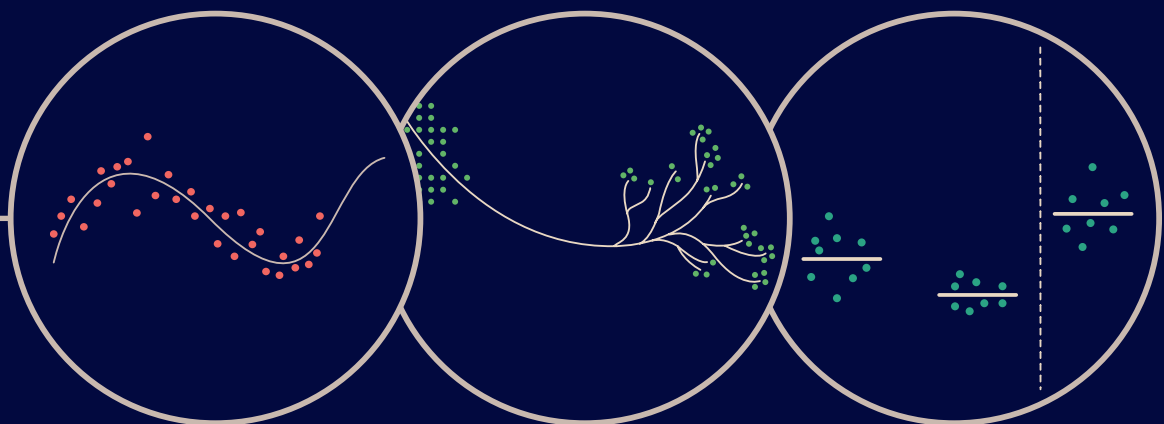# DATA ANALYSIS
## FOR BUSINESS, ECONOMICS, AND POLICY

Gábor Békés | Gábor Kézdi

# DATA ANALYSIS FOR BUSINESS, ECONOMICS, AND POLICY

This textbook provides future data analysts with the tools, methods, and skills needed to answer data-focused, real-life questions; to carry out data analysis; and to visualize and interpret results to support better decisions in business, economics, and public policy.

Data wrangling and exploration, regression analysis, machine learning, and causal analysis are comprehensively covered, as well as when, why, and how the methods work, and how they relate to each other.

As the most effective way to communicate data analysis, running case studies play a central role in this textbook. Each case starts with an industry-relevant question and answers it by using real-world data and applying the tools and methods covered in the textbook. Learning is then consolidated by over 420 practice questions and 117 data exercises.

Extensive online resources, including raw and cleaned data and codes for all analysis in Stata, R, and Python, can be found at www.gabors-data-analysis.com.

**Gábor Békés** is an Assistant Professor at the Department of Economics and Business of the Central European University, and director of the Business Analytics Program. He is a research affiliate at the Center for Economic Policy Research (CEPR). He has published in top economics journals on multinational firm activities and productivity, business clusters, and innovation spillovers. He managed international data collection projects on firm performance and supply chains. He has done both policy advising (the European Commission, ECB) as well as private-sector consultancy (in finance, business intelligence, and real estate). He has taught graduate-level data analysis and economic geography courses since 2012.

**Gábor Kézdi** is a research associate professor at the University of Michigan's Institute for Social Research. He has published in top journals in economics, statistics, and political science on topics including household finances, health, education, demography, and ethnic disadvantages and prejudice. He has managed several data collection projects in Europe; currently, he is co-investigator of the Health and Retirement Study in the US. He has consulted for various governmental and non-governmental institutions on the disadvantage of the Roma minority and the evaluation of social interventions. He has taught data analysis, econometrics, and labor economics from undergraduate to Ph.D. levels since 2002, and supervised a number of MA and Ph.D. students.

# DATA ANALYSIS FOR BUSINESS, ECONOMICS, AND POLICY

## Gábor Békés
Central European University, Budapest

## Gábor Kézdi
University of Michigan, Ann Arbor

CAMBRIDGE
UNIVERSITY PRESS

# CAMBRIDGE
## UNIVERSITY PRESS

University Printing House, Cambridge CB2 8BS, United Kingdom

One Liberty Plaza, 20th Floor, New York, NY 10006, USA

477 Williamstown Road, Port Melbourne, VIC 3207, Australia

314–321, 3rd Floor, Plot 3, Splendor Forum, Jasola District Centre, New Delhi – 110025, India

79 Anson Road, #06–04/06, Singapore 079906

Cambridge University Press is part of the University of Cambridge.

It furthers the University's mission by disseminating knowledge in the pursuit of education, learning, and research at the highest international levels of excellence.

# BRIEF CONTENTS

# CONTENTS

Contents                                                                              **xi**

## Contents                                                                xiii

Contents                                                                                      **xv**

Contents                                                                          **xvii**

Contents     **xix**

# WHY READ THIS BOOK

## An applied data analysis textbook for future professionals

**Data analysis is a process**. It starts with formulating a question and collecting appropriate data, or assessing whether the available data can help answer the question. Then comes cleaning and organizing the data, tedious but essential tasks that affect the results of the analysis as much as any other step in the process. Exploratory data analysis gives context to the eventual results and helps deciding the details of the analytical method to be applied. The main analysis consists of choosing and implementing the method to answer the question, with potential robustness checks. Along the way, correct interpretation and effective presentation of the results are crucial. Carefully crafted data visualization help summarize our findings and convey key messages. The final task is to answer the original question, with potential qualifications and directions for future inquiries.

Our textbook **equips future data analysts with the most important tools, methods and skills** they need through the entire process of data analysis to answer data focused, real life questions. We cover all the fundamental methods that help along the process of data analysis. The textbook is divided into four parts covering **data wrangling and exploration, regression analysis, prediction with machine learning, and causal analysis**. We explain when, why, and how the various methods work, and how they are related to each other.

Our approach has a **different focus compared to the typical textbooks** in econometrics and data science. They are often excellent in teaching many econometric and machine learning methods. But they don't give much guidance about how to carry out an actual data analysis project from beginning to end. Instead, students have to learn all of that when they work through individual projects, guided by their teachers, advisors, and peers – but not their textbooks.

To cover all of the steps that are necessary to carry out an actual data analysis project, we **lean on a set of fully developed case studies**. While each case study focuses on the particular method discussed in the chapter, they illustrate all elements of the process from question through analysis to conclusion. We facilitate individual work by **sharing all data and code in Stata, R, and Python**.

## Curated content and focus for the modern data analyst

Our textbook focuses on the most relevant tools and methods. Instead of dumping many methods on the students, we selected the most widely used methods that tend to work well in many situations. That choice allowed us to discuss each method in detail so students can gain a deep understanding of when, why, and how those methods work. It also allows us to compare the different methods both in general and in the course of our case studies.

The textbook is divided into four parts. The first part starts with data collection and data quality, followed by organizing and cleaning data, **exploratory data analysis** and data visualization, generalizing from the data, and hypothesis testing. The second part gives a thorough introduction to **regression analysis**, including probability models and time series regressions. The third part covers **predictive analytics** and introduces cross-validation, LASSO, tree-based machine learning methods such as random forest, probability prediction, classification, and forecasting from time series data. The fourth part covers **causal analysis**, starting with the potential outcomes framework and causal maps, then discussing experiments, difference-in-differences analysis, various panel data methods and the event study approach.

When deciding on which methods to discuss and in what depth, we drew on our own experience as well as the advice of many people. We have taught Data Analysis and Econometrics to students in Master's programs for years in Europe and the US, and trained experts in business analytics, economics, and economic policy. We used earlier versions of this textbook in many courses with students who differed in background, interest, and career plans. In addition, we have talked to many experts both in academia and in industry: teachers, researchers, analysts, and users of data analysis results. As a result, this textbook offers **a curated content that reflects the views of data analysts with a wide range of experiences**.

## Real life case studies in a central role

A cornerstone of this textbook are 47 case studies spreading over one-third of our material. This reflects our view that working through case studies is the best way to learn data analysis. Each of our case studies starts with a relevant question and answers it in the end, using real life data and applying the tools and methods covered in the particular chapter.

Similarly to other textbooks, our case studies illustrate the methods covered in the textbook. In contrast with other textbooks, though, they are much more than that.

Each of our case studies is a fully developed story linking business or policy questions to decisions in data selection, application of methods and discussion of results. Each case study uses **real-life data** that is messy and often complicated, and it discusses data quality issues and the steps of data cleaning and organization along the way. Then, each case study includes **exploratory data analysis** to clarify the context and help choose the methods for the subsequent analysis. After carrying out the main **analysis**, each case study emphasizes the correct **interpretation** of the results, effective ways to present and visualize the results, and many include robustness checks. Finally, each case study **answers the question** it started with, usually with the necessary qualifications, discussing internal and external validity, and often raising additional questions and directions for further investigation.

Our case studies are different also because they cover a wide range of topics, with a potential appeal to a wide range of students. They cover **consumer decision, economic and social policy, finance, business and management, health, and sport**. Their regional coverage is also wider than usual: one third is from the U.S.A., one third is from Europe and the U.K., and one third is from other countries or includes all countries from Australia to Thailand.

## Support material with data and code shared

We offer a truly comprehensive material with data, code for all case studies, 360 **practice questions**, 112 **data exercises**, derivations for advanced materials and reading suggestions. Each chapter ends with practice questions that help revise the material with a focus on theory. They are followed by data exercises that invite students to carry out analysis on their own, in the form of robustness checks or replicating the analysis using other data.

We share all raw and cleaned data we use in the case studies. We also share the codes that clean the data and produce all results, tables, and graphs in **Stata, R, and Python** so students can tinker with our code and compare the solutions in the different software.

All data and code are available on the textbook website:

`http://gabors-data-analysis.com`

## Who is this book for?

This textbook was written to be a **complete course** in data analysis. It introduces and discusses the most important concepts and methods in exploratory data analysis, regression analysis, machine learning and causal analysis. Thus, readers don't need to have a background in those areas.

The textbook includes formulae to define methods and tools, but it **explains all formulae in plain English**, both when a formula is introduced and, then, when it is used in a case study. Thus, understanding formulae is not necessary to learn data analysis from this textbook. They are of great help, though, and we encourage all students and practitioners to work with formulae whenever possible. The mathematics background required to understand these formulae is quite low, at the the level of basic calculus.

This textbook could be useful for university students in graduate programs as **core text** in applied statistics and econometrics, quantitative methods, or data analysis. The textbook is best used as core text for non-research degree Masters programs or part of the curriculum in a Phd or research Masters programs. It may also **complement online courses** that teach specific methods to give more context and explanation. Undergraduate courses can also make use of this textbook, even though the workload on students exceeds the typical undergraduate workload. Finally, the textbook can serve as a **handbook for practitioners** to guide them through all steps of real-life data analysis.

# SIMPLIFIED NOTATION

Finally, a note for the instructors who plan to use our textbook.

We introduced some new notation in this textbook, to make the formulae simpler and more focused. In particular, our **formula for regressions is slightly different** from the traditional formula. In line with other textbooks, we think that it is good practice to write out the formula for each regression that is analyzed. For this reason, it important to use a notation for the regression formula that is as simple as possible and focuses only on what we care about. Our notation is intuitive, but it's slightly different from traditional practice. Let us explain our reasons.

Our approach starts with the definition of the regression: it is a model for the conditional mean. The formulaic definition of the simple linear regression is $E[y|x] = \alpha + \beta x$. The formulaic definition of a linear regression with three right-hand-side variables is $E[y|x_1, x_2, x_3] = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$.

The regression formula we use in the textbook is a simplified version of this formulaic definition. In particular, we have $y^E$ on the left-hand-side instead of $E[y|...]$. $y^E$ is just a short-hand for the expected value of $y$ conditional on whatever is on the right-hand-side of the regression.

Thus, the formula for the simple linear regression is $y^E = \alpha + \beta x$, and $y^E$ is the expected value of $y$ conditional on $x$. The formula for the linear regression with three right-hand-side variables is $y^E = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$, and here $y^E$ is the expected value of y conditional on $x_1$, $x_2$, and $x_3$. Having $y^E$ on the left-hand-side makes notation much simpler than writing out the conditional expectation formula $E[y|...]$, especially when we have many right-hand-side variables.

In contrast, the traditional regression formula has variable $y$ itself on the left-hand-side not its conditional mean. Thus, it has to involve an additional element, the error term, which is the difference between $y$ and its conditional mean. For example, the traditional formula for the linear regression with three right-hand-side variables is $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + e$.

Our notation is simpler, because it has fewer elements. More importantly, our notation makes it explicit that the regression is a model for the conditional mean. It focuses on what data analysts care about (the right-hand-side variables and their coefficients), without adding anything else.

# ACKNOWLEDGMENTS

Let us first thank our students at Central European University, the University of Michigan, and at the University of Reading. The idea of writing a textbook was born out of teaching and mentoring them. We have learned a lot from teaching them, and many of them helped us writing code, collecting data, reading papers, and hunting for ideas.

Many colleagues helped us with their extremely valuable comments and suggestions. We thank *Eduardo Arino de la Rubia*, *Emily Blanchard*, *Alberto Cairo*, *Gergely Daróczi*, *János Divényi*, *Christian Fons-Rosen*, *Miklós Koren*, *Mike Luca*, *Róbert Lieli*, *László Mátyás*, *Timea Laura Molnar*, *Jenő Pál*, and *Ádám Szeidl* and anonymous reviewers of the first draft of the textbook.

We have received help with our case studies from *Alberto Cavallo*, *Daniella Scur*, *Nick Bloom*, *John van Reenen*, *Anikó Kristof*, *József Keleti*, and *Emily Oster* and *MyChelle Andrews*. We have learned a lot from them.

Several people helped us a great deal with our manuscript. At Cambridge University Press, our commissioning editor, *Phil Good* encouraged us from the day we met. Our editors *Heather Brolly*, *Jane Adams*, *Nicola Chapman*, guided us with kindness and steadfastness from first draft to proofs. We are not native English speakers, and support from *Chris Cartwright*, and *Jon Billam* were very useful. We are grateful for *Sarolta Rózsás* who acted as development editor, reviewing endless versions of chapters, checking consistency, clarity and pushed us to make the text more coherent and accessible.

Creating the code base in Stata, R and Python was a massive endeavour. Both of us are native Stata speakers, and we needed R code that would be fairly consistent with Stata code. Plus, all graphs were produced in R. So we needed help to have all our Stata codes replicated in R, and a great deal of code writing from scratch. *Zsuzsa Holler* and *Kinga Ritter* have provided enormous development support, spearheading this effort for years. Additional code and refactoring in Python and R was created by *Maté Tóth*, *János Bíró*, and *Eszter Pázmándi*. Additional coding, data collection, visualization and editing were done by *Viktória Kónya*, *Zsófia Komuves*, *Dániel Bánki*, *Abuzar Ali*, *Endre Borza* and *Ahmed Al Shaibani*.

Several former students helped us by reading the book and giving feedback as well as editing and checking latex files. We are grateful to *Imre Boda*, *Bonnie Kavoussi*, *Olivér Kiss*.

The wonderful cover design is based on the work by *Ágoston Nagy*, his first, surely not the last. The website was designed by *Imola Csoka*.

Collaborating with many talented young people and bringing them together was one the joys of writing this book.

Let us also shout out to the fantastic R user community – both online and offline – from whom we learnt tremendously. Special thanks to the Rstats and Econ Twitter community – we received wonderful suggestions from tons of people we have never met.

We thank the Central European University for professional and financial support. *Julius Horvath* and *Miklos Koren* as department heads provided massive support form the day we shared our plans.

Finally, let us thank those who were with us throughout the long, and often stressful, process of writing a textbook. Békés thanks *Saci*; Kezdi thanks *Zsuzsanna*, *Anna*, *Dorka*, and *Bori*. We could not have been able to done it without their love and support.

**The Gabors**
*2020 February, Budapest, Hungary and Ann Arbor, Mi, USA*