

Statistics for Social and Behavioral Sciences

Living Standards Analytics

Development through the Lens
of Household Survey Data

 Springer

Statistics for Social and Behavioral Sciences

Advisors:

S.E. Fienberg

W.J. van der Linden

For further volumes:

<http://www.springer.com/series/3463>

Dominique Haughton • Jonathan Haughton

Living Standards Analytics

Development through the Lens
of Household Survey Data

 Springer

Dominique Haughton
Department of Mathematical Sciences
Bentley College
Waltham, MA, USA
dhaughton@bentley.edu

Jonathan Haughton
Department of Economics
Suffolk University
Boston, MA, USA
jhaughto@beaconhill.org

ISBN 978-1-4614-0384-5 e-ISBN 978-1-4614-0385-2
DOI 10.1007/978-1-4614-0385-2
Springer New York Dordrecht Heidelberg London

Library of Congress Control Number: 2011934800

© Springer Science+Business Media, LLC 2011

All rights reserved. This work may not be translated or copied in whole or in part without the written permission of the publisher (Springer Science+Business Media, LLC, 233 Spring Street, New York, NY 10013, USA), except for brief excerpts in connection with reviews or scholarly analysis. Use in connection with any form of information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed is forbidden.

The use in this publication of trade names, trademarks, service marks, and similar terms, even if they are not identified as such, is not to be taken as an expression of opinion as to whether or not they are subject to proprietary rights.

Printed on acid-free paper

Springer is part of Springer Science+Business Media (www.springer.com)

*To our parents
Monique and Paul Boudier
and
Helen and Joe Haughton
and to our daughter Isabelle*

Preface

The Gallup Organization polls a thousand people every day. The Thailand Statistical Office interviews 3,000 households, using detailed surveys, every month. The amount of digital information doubles every 18 months.

We are, to use a headline from *The Economist*, facing a data deluge. What a contrast to the time when Nobel prize winner Wassily Leontief (1971), in his Presidential Address to the American Economic Association, could complain about a plethora of theory and a dearth of data, and call for a shift to “large-scale factual analysis” (p. 5).

The earliest analysis of household survey data – going back at least to the pioneering work of Seebom Rowntree (1901) – was largely confined to tabulations. Starting in 1980, the World Bank’s Living Standards Measurement Survey project boosted the pace and quality of household survey data gathered in Less-Developed Countries; 89 of the surveys may be downloaded from its Web site, but hundreds more such surveys are now available. By 2002 the project had generated 135 technical papers. This second wave emphasized the use of graphical and regression techniques, nicely summed up in the essential volume by Angus Deaton, *The Analysis of Household Surveys: A Microeconomic Approach to Economic Development* (1987).

We are now experiencing a third wave, with the increasing application of an ever-broadening array of analytical tools – such as directed acyclic graphs (DAGs), Kohonen maps, and propensity score matching – in addition to refinements to regression.

The purpose of this book is to introduce, discuss, illustrate, and evaluate the colorful palette of analytical techniques that can be applied to the analysis of household survey data, with an emphasis on the innovations of the past decade or so. It is conceived as an antidote to an overly narrow view of what constitutes legitimate empirical work, and reflects our own preferences as methodological eclectics.

The term “analytics” means the science of analysis, and in the business world – from which we have borrowed the term – it denotes the use of data, often in large quantities, to improve decision making. We use the term in its widest sense, as the harnessing of data, particularly from household surveys, to improve policy

recommendations. It is a large canvas, ranging from the mainstream econometric approach of testing and subsequently revising the sharp lines of model-generated hypotheses – what Deaton (2010, p. 4) calls the hypothetico-deductive method – to the soft brush strokes of some of the atheoretical techniques of data mining and exploratory data analysis. Both painting styles have their place in the gallery of analytics.

This is a gateway book. Most of the chapters begin by introducing a methodological or policy problem, to motivate the subsequent discussion of relevant methods. They then summarize the relevant techniques, and draw on examples – many of them from our own work – and aim to convey a sense of the potential, but also the strengths and weaknesses, of those techniques. The idea is to provide enough detail to allow the reader to take the next steps, but not so much detail as to get bogged down.

To be exhaustive would be too exhausting. For example, we introduce Kohonen maps in Chap. 6, explain how they function, and work through an example. The interested reader will then be well positioned to dig deeper, into a field where more than 5,000 articles have been published.

In writing this book, we have three main audiences in mind. The first is graduate students in statistics, economics, policy analysis, and social sciences, especially, but certainly not exclusively, those interested in the challenges of economic development in the Third World. We would be delighted if this book opens the reader to a handful of new ideas: skim the book, alight on the pages that catch one's fancy, and return to it regularly as a reference and a fount of ideas.

Our second target group is academics, who will likely be very conversant with some of the material in the book, but would appreciate a quick *tour d'horizon* to familiarize them with other interesting, and potentially useful, techniques. This is a book, like Deaton's *Analysis of Household Surveys*, that can serve as a reference work, to be taken down from the shelf and perused from time to time.

Our third audience is practitioners, by whom we mean anyone who works closely with survey data, whether in statistics offices, think tanks, research units, international organizations, central banks, NGOs, businesses – the list is long. We know, from teaching online and internationally, that there are many who, having left the university environment, are not sure how to keep up with new technical developments; we believe the book will help, because it introduces the techniques and ideas without getting too lost in the technical detail.

The Substance

We begin the book with a consideration of graphical methods, because this is often the first step when we are trying to develop a feel for our data. Graphs can be revealing, and they can be helpful in presenting our findings. We start by discussing how to produce a useful histogram, and its continuous-valued cousin, the kernel density. Boxplots are also easy to use and especially helpful when we want to

compare the essential features of two or more distributions side by side. The chapter also includes some discussion of violin plots, scatterplots, and bag plots, before turning to presentational graphics. We agree with Gelman et al. (2002) that graphs could productively be used more often when presenting scientific results: The beautiful bubble plot in Fig. 1.13 contains more information than its apparent simplicity would suggest. The final section of Chap. 1 looks at maps, which can now be produced remarkably quickly and easily; the cartogram in Fig. 1.17 shows the distribution of child mortality worldwide, and instantly conveys the locus of the problem.

After graphics comes regression, which we survey in Chap. 2. Seasoned econometricians and other quantitative researchers can skip this chapter, but it is our experience that regression is sufficiently subtle, and the ideas sufficiently slippery, that one needs a quick review of the material on a regular basis. We note the main problems faced in regression, including measurement error, omitted variable bias, multicollinearity, heteroscedasticity, adjustments for clustered data, outliers, and simultaneity, and suggest ways in which these may be dealt with. Thus the chapter includes a discussion of, among other things, instrumental variables, and quantile regression. It is a whistle stop tour, which is exactly what most of us need.

Household survey data almost never come from simple random samples, and in Chap. 3 we address the issues related to sampling, first reviewing the main types – simple, stratified, cluster – and then presenting the essentials of how to determine an appropriate sample size while recognizing the need to trade off sampling with nonsampling errors. We show how to incorporate sample design into the computation of summary statistics – using Stata, the statistical package that we have used most over the years – and summarize the debate on whether to use weights in regression. The last two sections of the chapter ask how best to survey hard-to-reach groups, such as migrants – the main focus of a recent survey in the two main cities of Vietnam – and groups such as jazz players, or prostitutes, where respondent-driven sampling has been quite successful.

In Chap. 4 we move beyond linear regression, first by making the linear specification more flexible, and then by using nonparametric methods to fit curves. This segues into an explanation of multivariate adaptive regression spline (MARS) models, which we apply to a model of changes in consumption spending in Vietnam between 1993 and 1998. We also discuss classification and regression tree (CART) models; both CART and MARS are particularly good at exploring the data for nonlinearities and interactions. We have used a CART model with some success as a first step in helping us specify the functional form of a parametric model of the determinants of short-term malnutrition in Vietnam.

Much of our interest in working with living standards survey data arises from our desire to say something useful for policy purposes. This requires us to be able to say, “if you do X, then Y will happen,” which is a causal statement. The question of causality, and more specifically how to conceive of and measure causal statements, is the subject of Chap. 5. The experimentalist school focuses on measuring the “effects of causes,” where possible using randomized experiments to try to determine whether microcredit raises spending or flip charts improve exam performance.

The structuralist school worries that the outcomes of experiments leave us with an insufficient understanding of the underlying causal mechanisms, and urge us to pay attention to unearthing the “causes of effects,” which may then be generalized to other situations and applied to policy. Taking its cue from Edward Tufte, who famously wrote that “correlation is not causation but it sure is a hint,” the causal inference school, seeks to measure causality using a combination of correlations and logic. This approach is essentially mechanical, and the results are usually shown in the form of directed acyclic graphs (DAGs). This is unfamiliar terrain for most economists and policy analysts, which is why we devote much of the chapter to explaining how DAGs are constructed and what we might learn from them.

We often group data, for instance looking at income by gender, region, or quintile. In Chap. 6 we explore in more detail how observations may be clustered. This is an exploratory process, traditionally conducted with hierarchical or non-hierarchical clustering, which can produce beautiful graphs. It is also possible to incorporate more statistical structure using latent class models. The second half of the chapter introduces Kohonen maps, which have become very popular: They typically group observations on a two-dimensional grid, and present the results in the form of gorgeous “maps” – all of which we explain and illustrate here.

In approaching any scientific question, or looking at any data, we almost always have at least some idea of what we expect. If the data showed that richer households bought fewer cars, or poorer households eat more caviar, we would be shocked. Bayesian analysis provides a formal framework for incorporating these prior beliefs, in contrast to the more standard frequentist approach that either ignores them entirely, or locks them into rigid models. Chapter 7 provides an introduction to Bayesian analysis, setting out the ideas, the approach, and an example, and then addressing the problems of eliciting priors, applying posterior predictive checking, combining models in the form of Bayesian model averaging, and determining the appropriate sample size for a survey. This is not the easiest chapter in the book – the intrinsic difficulty of the subject helps explain its still-limited spread beyond trained statisticians – but it is likely to be one of the more useful for nonstatistician readers.

We are rediscovering geography, and recognizing once again that what happens in one area can influence what happens nearby. The presence of spatial dependence has implications for how to specify and estimate regression models – most commonly through the use of spatial weights matrices that measure the strength of the contiguity effects. We illustrate the use of these techniques in Chap. 8, drawing on a study of the spatial pattern of unemployment in the Midi-Pyrénées region of France, where we also present an algorithm for choosing among different types of spatial models.

Although it is still comparatively rare, increasing numbers of household surveys are based on panels, where households are surveyed repeatedly over time. In Chap. 9 we show how panel data can allow for more precise inference, and in many cases can help us tackle the knotty problem of unobserved heterogeneity: if households differ in ways we cannot observe, but these differences – in ability or drive, for instance – do not vary over time, then differenced data can sweep

away such effects, laying bare the relationships that we are usually interested in measuring. We illustrate this with an example in which we try to measure the effect on income of loans extended under the Thailand Village Fund, which burst onto the scene in 2002 and by 2004 had become the largest microcredit scheme in the world. Still, panel data are not a panacea; attrition bias can be a problem, and even without attrition, panels become less representative over time.

One of the most important uses of household survey data is to measure poverty, and vulnerability to poverty. Chapter 10 reviews this field, starting with the choice of a measure of well-being, through the construction of a poverty line, to the choice of a summary measure of poverty. We then discuss the robustness of poverty measures, focusing on sampling and measurement error, and explaining the notion of stochastic dominance. After a section in which we consider the problems peculiar to international comparisons of poverty, we consider ways in which vulnerability to poverty – defined as the probability that a household will be poor in the future – may be measured.

We return to an essentially technical issue in Chap. 11, where we look at bootstrapping. This is especially useful when we need to estimate the standard error of a measure – such as the Sen–Shorrocks–Thon index of poverty – and where an analytical formula is not available. The technique can be powerful, especially where the data come from complex samples, and is increasingly straightforward to implement; we illustrate this with an example in which we create a histogram of bootstrapped changes in the poverty rate in Vietnam between 1993 and 1998.

Does a program work? Was a project effective? These are questions addressed by impact evaluation, where we try to compare the actual outcomes, for those who have been “treated,” with a counterfactual, which is our estimate of what would have happened in the absence of the program or project. The traditional gold standard is experimental design, or randomization, but in Chap. 12 we show that even this is not without its limitations. It is much more common to use quasi-experimental methods, of which the most popular are propensity score matching, double differences, and instrumental variables. For each of these we set out the principles, consider an example, and review both the strengths and weaknesses. This is a relatively long and detailed chapter, but it has proven to be effective when teaching impact evaluation to graduate students in economics.

Household survey data mainly come from large, complex questionnaires administered to relatively small samples of perhaps 5,000–10,000. This allows one to conduct the analysis at the level of a country or broad region, but not at the level of a small county or district. Yet we would often like to measure, for instance, poverty rates at a “small-area” level, the better to target spending to alleviate poverty. In Chap. 13 we discuss how to do this, first describing a basic synthetic regression model, and then explaining how one might estimate a two-level, or even multilevel, model with random effects. This chapter applies the methods to Vietnam, and includes two elegant maps that result from the analysis.

Perhaps it is fitting that the last chapter in the book, Chap. 14, looks at duration models. In many cases, the time dimension is central to the analysis, such as

the interval between one birth and the next, or the time spent unemployed. We introduce the Kaplan–Meier estimator, which allows for an exploratory analysis of duration data, and move on to the Cox proportional hazards model, parametric regression models, and mixture models of two Weibull regressions. As always, this chapter is designed to help the reader take the first steps – enough for the first draft of a solid research paper, even if lifting it to the level required for scholarly publication will always call for digging a bit deeper.

Where We Stand

We come to this book with different perspectives – one schooled in economics where the mindset is one of “model first, then test,” the other more comfortable with data mining and letting the numbers speak “for themselves.” The tension between these approaches runs throughout the book, and we see this as a virtue. One of us is skeptical that directed acyclic graphs are useful in helping us understand how the world really works, and thinks that the main virtue of Kohonen maps is that they are pretty. The other has yet to find an instrumental variable that looks compelling, and thinks that a lot of highfalutin theory is “nonsense on stilts.” We do not try to resolve these debates – we are reminded of the observation by George Box that “essentially, all models are wrong, but some are useful” – but instead set out the techniques and ideas, to help the reader develop an informed opinion.

Together, we have over 50 years of experience working with household datasets, and have written over 200 papers, articles, and reports, over 80 of them in scholarly journals. This book is our take on what we find to be most useful, or at least intriguing or innovative; it also contains what we would like our students to know.

We are grateful to all of those who helped us on the way to this book. All of our more than 150 co-authors have at least some claim to intellectual parentage. Glenn Jenkins started the ball rolling in 1979 by interesting one of us in using survey data to address a practical development problem, in this case whether to build small-scale irrigation projects in Malaysia. In 1994, Mark Sidel encouraged us to work with the General Statistics Office in Hanoi; this, and the ongoing support from Nguyen Phong of the GSO, explains why so many of the examples in this book are drawn from the various living standards surveys undertaken in Vietnam.

We would like to thank our institutions – Bentley University and Suffolk University – for providing research support and sabbatical leaves that helped us get the book written. We are grateful to John Kimmel for trusting us with the project, and waiting patiently for it to progress, and to Marc Strauss for taking up the baton; to Dan Westbrook for reviewing an early draft; and to Maria Skaletsky, Sunida Susantud, Bayar Tumennasan, and Jason Wells for very helpful comments.

References

- Deaton, Angus. 1997. *The analysis of household surveys: A microeconomic approach to development policy*. Baltimore: Johns Hopkins University Press.
- The Economist*. 2010. Special report on managing information. February 27.
- Gelman, A., C. Pasarica, and R. Dodhia. 2002. Let's practice what we preach: Turning tables into graphs. *The American Statistician*, 56(2):121–130.
- Leontief, Wassily. 1971. Theoretical assumptions and nonobserved facts. *American Economic Review*, 61(1):1–7.
- Rowntree, Seebohm. 1901. *Poverty: A study of town life*. London: Macmillan.

About the Authors

Dominique Haughton (PhD, MIT 1983) is Professor of Mathematical Sciences at Bentley University in Waltham, Massachusetts, near Boston, and Affiliated Researcher at the Université Toulouse 1, France. Her major areas of interest are applied statistics, statistics and marketing, the analysis of living standards surveys, data mining, and model selection. She is the editor-in-chief of *Case Studies in Business, Industry and Government Statistics (CSBIGS)*, and has published over 50 articles in scholarly journals, including *The American Statistician*, *Annals of Statistics*, *Sankhya*, *Communications in Statistics*, and *Statistica Sinica*. In 2011, she was elected a Fellow of the American Statistical Association.

Jonathan Haughton (PhD, Harvard 1983) is Professor of Economics at Suffolk University, and Senior Economist at the Beacon Hill Institute for Public Policy, both in Boston. A specialist in the areas of economic development, international trade, and taxation, and a prize-winning teacher, he has lectured, taught, or conducted research in over a score of countries on five continents. His *Handbook on Poverty and Inequality* (with Shahidur Khandker) was published by the World Bank in 2009, his articles have appeared in over 30 scholarly journals, and he has written numerous book chapters and over a hundred reports.

Contents

1 Graphical Methods	1
1.1 Introduction	1
1.2 Exploratory Graphical Methods	3
1.2.1 Histograms	3
1.2.2 Kernel Densities.....	6
1.2.3 Boxplots	9
1.2.4 Scatterplots	12
1.2.5 Bagplots	14
1.3 Presentational Graphics	15
1.4 Statistics with Maps.....	18
1.5 Conclusion	21
References.....	21
2 Regression	23
2.1 Introduction	23
2.2 Basics.....	24
2.2.1 Inference	25
2.3 Addressing Regression Problems	28
2.3.1 Measurement Error.....	28
2.3.2 Omitted Variable Bias	30
2.3.3 Multicollinearity	32
2.3.4 Heteroscedasticity.....	33
2.3.5 Clustering	37
2.3.6 Outliers.....	38
2.3.7 Simultaneity.....	43
2.4 Conclusion	47
References.....	48

- 3 Sampling**..... 49
 - 3.1 Introduction 49
 - 3.2 Types of Sampling 51
 - 3.2.1 Simple Random Sampling 51
 - 3.2.2 Stratified Sampling..... 51
 - 3.2.3 Cluster Sampling..... 52
 - 3.3 Sample Size 52
 - 3.3.1 Sampling vs. Nonsampling Errors 53
 - 3.4 Incorporating Sample Design 54
 - 3.5 Design vs. Model-Based Sampling 56
 - 3.5.1 Illustration: Design-Based vs. Model-Based Means..... 56
 - 3.6 Weights or Not?..... 57
 - 3.6.1 Illustration: Weighting Regression 59
 - 3.7 Sampling Hard-to-Reach Groups: Vietnam 60
 - 3.8 Respondent-Driven Sampling: Hard-to-Reach Groups..... 62
 - References..... 65
- 4 Beyond Linear Regression** 67
 - 4.1 Introduction 67
 - 4.2 Flexibility in Linear Regression Models 68
 - 4.3 Nonlinear Models 71
 - 4.4 Nonparametric Models..... 72
 - 4.5 Higher-Dimension Models..... 75
 - 4.5.1 MARS Models 76
 - 4.5.2 MARS Application: Changes in Expenditure
in Vietnam, 1993–1998 79
 - 4.5.3 CART Models..... 82
 - 4.5.4 CART as a Preprocessor: A Nutrition Example..... 84
 - 4.5.5 CART as a Classifier: An Expenditure Example..... 87
 - 4.6 Beyond Regression..... 88
 - References..... 90
- 5 Causality**..... 91
 - 5.1 Introduction 91
 - 5.2 The Experimentalist School 91
 - 5.3 The Structuralist School 93
 - 5.4 The Causal Inference School 94
 - 5.4.1 Establishing Causality 95
 - 5.5 Creating Directed Acyclic Graphs..... 96
 - 5.5.1 Basic Tools: Probability 97
 - 5.5.2 Directed Acyclic Graphs..... 97
 - 5.5.3 d-Separation..... 98
 - 5.5.4 Illustrative Example of Measuring
Causality 100
 - 5.5.5 Tetrad, and the Partial Correlation Algorithm..... 101

- 5.6 A DAG to Explain World Poverty..... 102
 - 5.6.1 DAGs and Theory: Publishing Productivity..... 104
- 5.7 Conclusion 106
- References..... 107
- 6 Grouping Methods..... 109**
 - 6.1 Introduction 109
 - 6.2 Hierarchical Cluster Analysis..... 110
 - 6.3 Nonhierarchical Clustering 111
 - 6.4 Examples of Cluster Analysis 112
 - 6.4.1 Regions of Slovakia..... 112
 - 6.4.2 Households in South Africa..... 113
 - 6.5 Model-Based Clustering: Latent Class Models..... 116
 - 6.5.1 Applications..... 117
 - 6.6 Case Study: Vietnamese Households, 2002 118
 - 6.7 Kohonen Maps 121
 - 6.7.1 Building Kohonen Maps 125
 - 6.7.2 An Illustration..... 126
 - 6.8 Conclusion 127
 - References..... 128
- 7 Bayesian Analysis 129**
 - 7.1 Introduction 129
 - 7.2 A Worked Example..... 131
 - 7.2.1 Assuming i.i.d. Observations 132
 - 7.2.2 Taking Survey Sampling into Account 134
 - 7.2.3 An Illustration..... 135
 - 7.3 Prior Distributions and Implications of Their Choice..... 137
 - 7.3.1 Eliciting Priors 137
 - 7.3.2 Noninformative Priors 138
 - 7.3.3 Priors in a Linear Regression Context..... 138
 - 7.4 Bayes Factors and Posterior Predictive Checking 139
 - 7.4.1 Bayes Factors 139
 - 7.4.2 Posterior Predictive Checking 141
 - 7.4.3 Example: Fixed vs. Random Effects 142
 - 7.4.4 Example: Modeling Diagnostic Tests 143
 - 7.5 Combining Models: Bayesian Model Averaging..... 145
 - 7.5.1 Practical Issues..... 146
 - 7.6 Bayesian Approach to Sample Size Determination 148
 - 7.7 Conclusion 150
 - References..... 152
- 8 Spatial Models 155**
 - 8.1 Introduction..... 155
 - 8.2 The Starting Point: Including Spatial Variables 156
 - 8.2.1 Exploratory Spatial Data Analysis..... 156
 - 8.2.2 Including Spatial Variables 157

8.3	Spatial Models	159
8.3.1	Spatial Dependence	159
8.3.2	Spatial Heterogeneity	160
8.4	Classifying Spatial Models	161
8.4.1	Measuring Spatial Contiguity	163
8.4.2	Types of Spatial Model	164
8.4.3	Illustrating the Choice of Spatial Model	167
8.5	Other Spatial Models	169
8.5.1	Spatial Expansion Models	170
8.5.2	Geographically Weighted Regression	171
8.5.3	Spatial Effects as Random Effects	171
8.6	Conclusion	172
8.6.1	Estimating Spatial Models	173
	References	173
9	Panel Data	175
9.1	Introduction	175
9.2	Types of Panel Data	175
9.3	Why Panel Data?	176
9.4	Why Not Panel Data?	178
9.5	Application: The Birth and Growth of NFHEs	179
9.6	Statistical Analysis of Panel Data	181
9.7	Illustration: Thai Microcredit	183
	References	186
10	Measuring Poverty and Vulnerability	189
10.1	Introduction	189
10.2	What and Why?	189
10.3	Basic Measurement	190
10.3.1	Measuring Well-Being	190
10.3.2	Adult Equivalents	193
10.3.3	Choosing a Poverty Line	194
10.3.4	Summarizing Poverty Information	201
10.4	Robustness	205
10.4.1	Sampling Error	205
10.4.2	Measurement Error	206
10.4.3	Equivalence Scales	209
10.4.4	Choice of Poverty Line	209
10.5	International Poverty Comparisons	211
10.6	Vulnerability to Poverty	214
	References	218

11	Bootstrapping	221
11.1	Introduction	221
11.2	Bootstrap: Mechanics	222
11.2.1	Further Considerations	224
11.3	Applications to Living Standards	225
11.3.1	SST Index for Vietnam	227
11.3.2	Measuring Vulnerability	227
11.4	Bootstrapping Inequality and Regression	230
11.4.1	Regression	230
11.5	Has Poverty Changed?	231
11.6	Conclusion	233
	References	234
12	Impact Evaluation	235
12.1	Introduction	235
12.2	General Principles	236
12.2.1	Case: The Thailand Village Fund	236
12.2.2	A More Formal Treatment	238
12.3	Experimental Design	240
12.3.1	Case Study: Flip Charts in Kenya	241
12.3.2	Partial Randomization	242
12.3.3	Randomization Evaluated	243
12.4	Quasi-Experimental Methods	245
12.4.1	Solution 1. Matching Comparisons	247
12.4.2	Propensity Score Matching	248
12.4.3	Covariate Matching	255
12.4.4	Solution 2. Double Differences	258
12.4.5	Solution 3. Instrumental Variables	264
12.4.6	Other Solutions	267
12.5	Impact Evaluation: Macro Projects	268
12.5.1	Time-Series Data Analysis: Deviations from Trend	268
12.5.2	CGE and Simulation Models	268
12.5.3	Household Panel Impact Analysis	269
12.5.4	Self-Rated Retrospective Evaluation	269
12.6	In Conclusion	270
	References	271
13	Multilevel Models and Small-Area Estimation	273
13.1	Introduction	273
13.2	Simple Small-Area Models	274
13.3	Synthetic Regression Models	275
13.3.1	An Illustration: Poverty Mapping in Vietnam	275

- 13.4 Random Effects and Multilevel Models..... 276
 - 13.4.1 Basic Idea..... 276
 - 13.4.2 Specifying and Estimating a Two-Level Model 277
 - 13.4.3 An Example: Expenditures in Vietnam 278
 - 13.4.4 Rationale for Using Multilevel Models
for Small-Area Estimation..... 279
- 13.5 Conclusion 285
- References..... 286
- 14 Duration Models..... 289**
 - 14.1 Introduction 289
 - 14.2 Basics 289
 - 14.3 An Exploratory Analysis of Duration Data 291
 - 14.3.1 The Kaplan–Meier Estimator..... 292
 - 14.4 Cox Proportional Hazards Model 293
 - 14.5 Parametric Regression Models 296
 - 14.5.1 Weibull Regression Models 296
 - 14.5.2 A Mixture of Two-Weibull Regression Models 300
 - 14.6 Other Applications..... 302
 - References..... 304
- Index..... 307**

Chapter 1

Graphical Methods

1.1 Introduction

It is tempting, but wrong, to believe that graphical techniques have little to offer for serious researchers in economics, statistics, or policy analysis. Their true power comes from the ability of the eye to discern patterns in a graph that are not clearly evident from lists of numbers or tabulated statistics. In Tufte's pithy phrase, "graphics reveal data" (Tufte 2001, p. 13).

We explore this theme in the chapter, beginning with the use of basic exploratory graphical methods in Sect. 1.2, considering presentational graphics in Sect. 1.3, and introducing some more recent techniques, including maps, in Sect. 1.4.

With data in hand, the most productive first step is often to explore the data graphically. These graphs do not have to be especially polished and beautiful; rather, they need to be easy to produce and thoroughly informative, a visual scratch pad where we use the power of graphics to get a sense of the shape of variables and the interactions among them.

Following Tufte (2001), the point can be emphasized elegantly with the help of Anscombe's quartet – four data sets, reproduced in Table 1.1, that may be summarized by the same linear model, and where the mean values of the X and Y variables are the same in each case. Yet a graphical display of the data sets (Fig. 1.1) demonstrates how very different they are. Real data do not usually yield such coherent or clear patterns, but a good initial graphical analysis can easily come up with surprises – showing outliers, suggesting a need to use a mixture of distributions, or raising questions about how variables are related.

Graphical techniques are also exceptionally useful in presenting the results of one's analysis, and we agree with Gelman et al. (2002) that they are typically underutilized for this purpose. But presentational graphics require an approach that is quite different from that of exploratory graphics: they serve to communicate ideas to others, and so they need to be more beautiful and more carefully constructed.

Table 1.1 Anscombe's quartet of data sets

I		II		III		IV	
X	Y	X	Y	X	Y	X	Y
10.0	8.04	10.0	9.14	10.0	7.46	8.0	6.58
8.0	6.95	8.0	8.14	8.0	6.77	8.0	5.76
13.0	7.58	13.0	8.74	13.0	12.74	8.0	7.71
9.0	8.81	9.0	8.77	9.0	7.11	8.0	8.84
11.0	8.33	11.0	9.26	11.0	7.81	8.0	8.47
14.0	9.96	14.0	8.10	14.0	8.84	8.0	7.04
6.0	7.24	6.0	6.13	6.0	6.08	8.0	5.25
4.0	4.26	4.0	3.10	4.0	5.39	19.0	12.50
12.0	10.84	12.0	9.13	12.0	8.15	8.0	5.56
7.0	4.82	7.0	7.26	7.0	6.42	8.0	7.91
5.0	5.68	5.0	4.74	5.0	5.73	8.0	6.89
9.0	7.5	9.0	7.5	9.0	7.5	9.0	7.5

Notes: Each data set has 11 observations; the means are shown in the bottom row. Every regression line is $Y = 3 + 0.5X$; the standard error of the slope coefficient is 0.118 and its t -statistic is 4.24. In every case, $R^2 = 0.67$

Source: Anscombe 1973

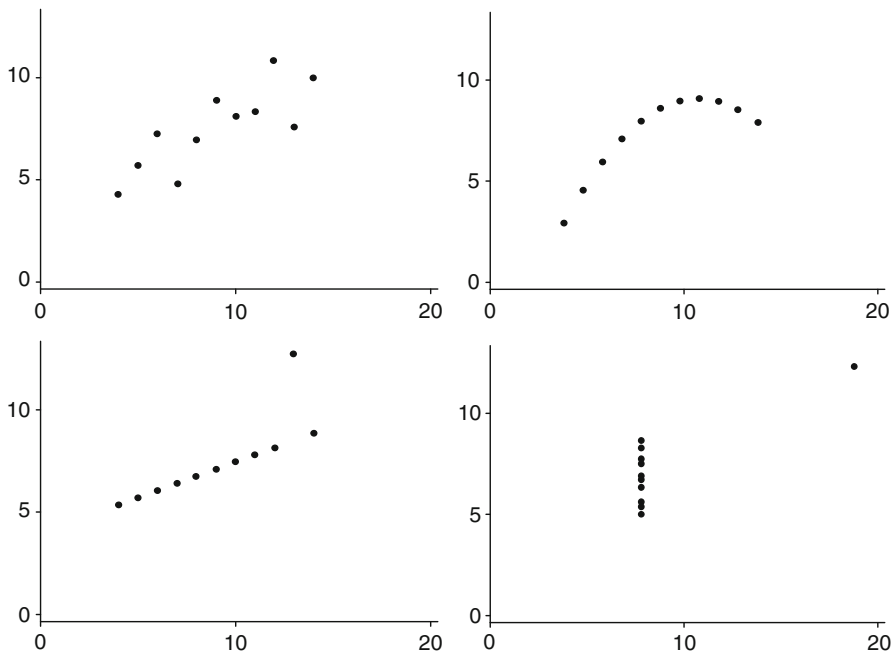


Fig. 1.1 Scatterplot of Anscombe's quartet (Note: Data from Table 1.1)

In Sect. 1.3 we review some of the key principles of graphical design, drawing heavily on the work of Tufte (2001), and suggest a few ways in which graphs could be used to make standard tabular presentations more effective.

1.2 Exploratory Graphical Methods

1.2.1 Histograms

A good place to start any analysis is with the most basic of visual techniques. Consider Fig. 1.2, which shows a simple frequency distribution (histogram) of birth weights of children born in Vietnam in 1992–1993. The data come from the Vietnam Living Standards Survey of 1992–1993, which surveyed 4,800 households nationwide and collected information on birth weights for 1,687 children. The graph represented the first step in an analysis by Sarah Bales (1999) of the determinants of low birth weights, and was generated using Stata.¹

A baby is typically defined as being underweight if he or she weighs less than 2.5 kg at birth. Thus the histogram in Fig. 1.2 alerts us to a problem: an implausibly large number of births are heaped into the 2.5 kg category (and the 3.0, 3.5, and 4.0 kg categories). Indeed, 10.1% of the births were reported as weighing less than 2.5 kg and a further 10.7% as weighing exactly 2.5 kg! The rounding error matters here; the weight of some babies has presumably been rounded up to 2.5 kg, and in other cases the weight has been rounded down to 2.5 kg. So, while it is clear that more than 10.1% of babies are born underweight, but fewer than 20.8%, it is not clear whether it is preferable to define “underweight” as $w < 2.5$ or $w \leq 2.5$ (where w refers to the weight of the baby in kilos). The solution chosen by Bales (1999) was to use both definitions; fortunately, she found that the exact definition of underweight made relatively little difference to the direction and strength of the determinants of low birth weights.

Like a stethoscope, a histogram appears to be a simple tool, but it takes some practice to make it work effectively. The key choice that has to be made is that of the number of classes (“bins”) into which to group the data or, alternatively, the width of each class, and this choice is as much a matter of art as of science.

A histogram aims to lay bare the distribution of the underlying data, and the classification of data into bins serves to filter out some of the noise. This is illustrated in Fig. 1.3, which displays four histograms showing the number of individuals covered by the 1998 Vietnam Living Standards Survey, broken down by age. The bottom right panel of Fig. 1.3 has just ten bins, and hints at a unimodal distribution dominated by the large proportion of individuals in the 10–20 age

¹For a tutorial-based introduction to Stata, with examples that use easily accessible household survey data from Bangladesh, see Appendixes 1 and 2 of Haughton and Khandker (2009).

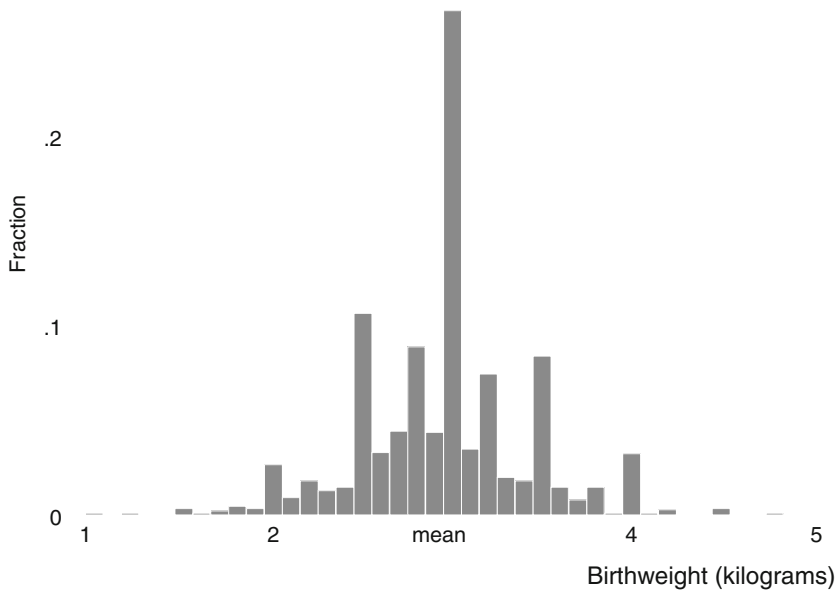


Fig. 1.2 Histogram of birth weights in Vietnam, 1993 (*Source: Vietnam Living Standards Survey, 1993*)

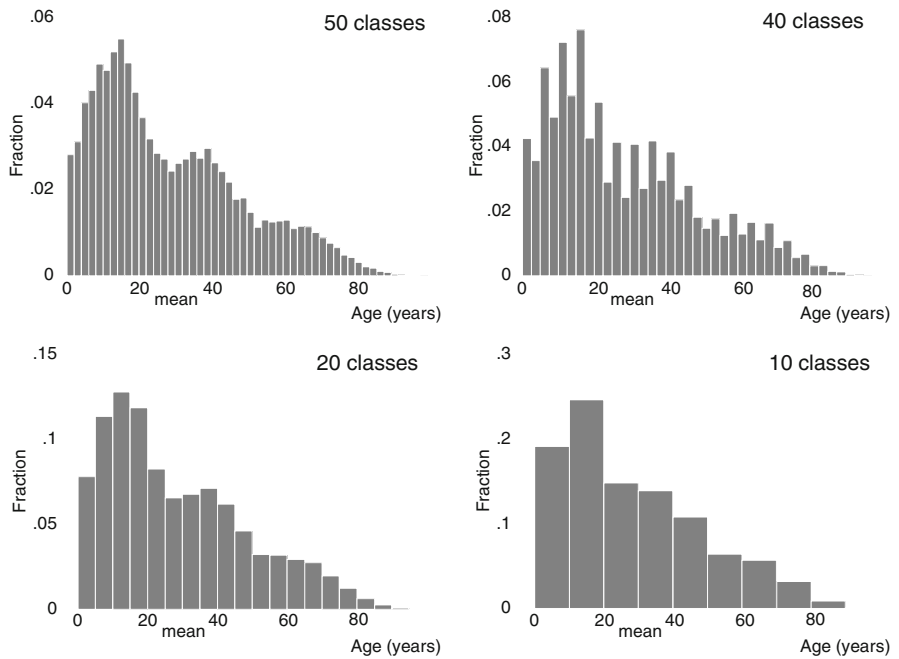


Fig. 1.3 Histograms of ages (in years) for individuals surveyed in the Vietnam Living Standards Survey of 1997–1998, with differing numbers of classes (“bins”)

bracket. The top right panel uses 40 bins, and is clearly unsatisfactory; the age interval 0–100 is divided into 40 equal classes, so those aged 0, 1, or 2 are in the first bracket, those aged 3 or 4 in the next bracket, those aged 5, 6 or 7 in the third bin, and so on, in a saw-toothed fashion. This particular problem arises because the age data are integer rather than continuous, but even continuous data are often subject to heaping, as we saw in Fig. 1.2.

We get closer to a sensible pattern with a 20-bin histogram, which suggests a second peak. This is confirmed by the very informative histogram in the top left panel of Fig. 1.3, which uses 50 classes. Turned on its side, this gives half of a population “pyramid.” Note, in this context, the dramatic reduction in the number of births since they peaked in about 1972 (i.e., 16 years prior to the 1998 survey), and the shortage of those aged roughly 45–60 – casualties of war, as they would have been in the armed forces in the years up to 1975 – and their children (Houghton 2000).

Are there better ways to choose bin widths other than trial and error? Freedman and Diaconis (1981) argue that if the histogram is to serve as a density estimator, then an appropriate rule for determining bin width is

$$BW_{FD}(x) = 2 \times IQR(x) \times n^{-1/3}, \quad (1.1)$$

where $BW_{FD}(x)$ is the Freedman and Diaconis bin width for variable x , $IQR(x)$ is the interquartile range of x , and n is the number of observations. With more observations we can afford to have narrower bins; with greater variation in x the bins need to be wider.

Other rules have been suggested. For instance, Wand (1997) proposes starting with a “zero-stage rule” that sets the bin width as

$$BW_W(x) = 3.49 \times \left[\min \left\{ s, \frac{IQR(x)}{1.349} \right\} \right] \times n^{-1/3}, \quad (1.2)$$

where s is the standard deviation of the sample.

In the example in Fig. 1.3, the Freedman–Diaconis rule gives a bin width of 1.9, implying 52 bins, while the Wand rule generates a bin width of 2.3, implying 43 bins. Neither rule gives results that are as clean as those with 50 bins, but they would filter the data nicely if one were using truly continuous (rather than integer) data.

Statistical software packages try to help the user by starting with sensible guesses of the appropriate number of bins. Microsoft Excel sets the number of bins equal to \sqrt{n} (rounded to the next integer) or 50, whichever is the smallest. Stata uses an only slightly more complex default, which is

$$\text{Number of bins} = \min[50, \min\{\sqrt{n}, 10 \times \ln(n) / \ln(10)\}]. \quad (1.3)$$

This sometimes works well, but usually some further exploration is called for to produce a sensible histogram.

1.2.2 Kernel Densities

A histogram provides a discretized, nonparametric approximation to the underlying density, but it has three drawbacks: it is not smooth, it depends on the bin widths, and it is sensitive to the choice of end points of the bins. So it is often more useful, or at least more elegant, to work with a smoothed version. This is achieved by estimating a kernel density.

Suppose we have a dataset X_1, X_2, \dots, X_n , and array the observations on the horizontal axis of a graph. We are interested in estimating the density, $f(x)$, at any given point x . A natural way to measure the density is by measuring the concentration of observed data points that are in the vicinity of x , say in the interval $x \pm h$, where h is the bandwidth half length. As we move x and its associated interval rightwards along the horizontal axis, we drop points to the left and pick up new observations on the right. The effect on the total number of observations is gradual, hence the smoothing effect.

The process we have described here may be formalized. It generates the naïve (or rectangular) estimate of the density at x , given by

$$\hat{f}_N(x) = \frac{1}{hn} \sum_{i=1}^n W\left(\frac{x - X_i}{h}\right) \equiv \frac{1}{hn} \sum_{i=1}^n W(z), \quad (1.4)$$

where

$$w(z) = \begin{cases} 1/2 & \text{if } |z| < 1 \\ 0 & \text{otherwise} \end{cases} \quad (1.5)$$

so that

$$\hat{f}_N(x) = \frac{1}{nh} \sum_{i=1}^n \frac{1}{2} I(|z| \leq 1), \quad (1.6)$$

where $I(\cdot)$ is an indicator function that takes on a value of 1 if the bracketed expression is true, and zero otherwise.

More generally, we may define a kernel density estimator as

$$\hat{f}(x) = \frac{1}{hn} \left[\sum_{i=1}^n K(z) \right], \quad (1.7)$$

where $K(z)$ is the kernel function, calibrated so that the estimator integrates to 1. The naïve estimator puts an equal weight on all the observations in the interval $x \pm h$ when estimating the density of x , which is why the kernel function in this case is referred to as “rectangular.” However, it is usually more satisfactory to use a symmetric function that puts more weight on values of X_i that are closer to x , and progressively less weight on values further from x . The widely used Epanechnikov

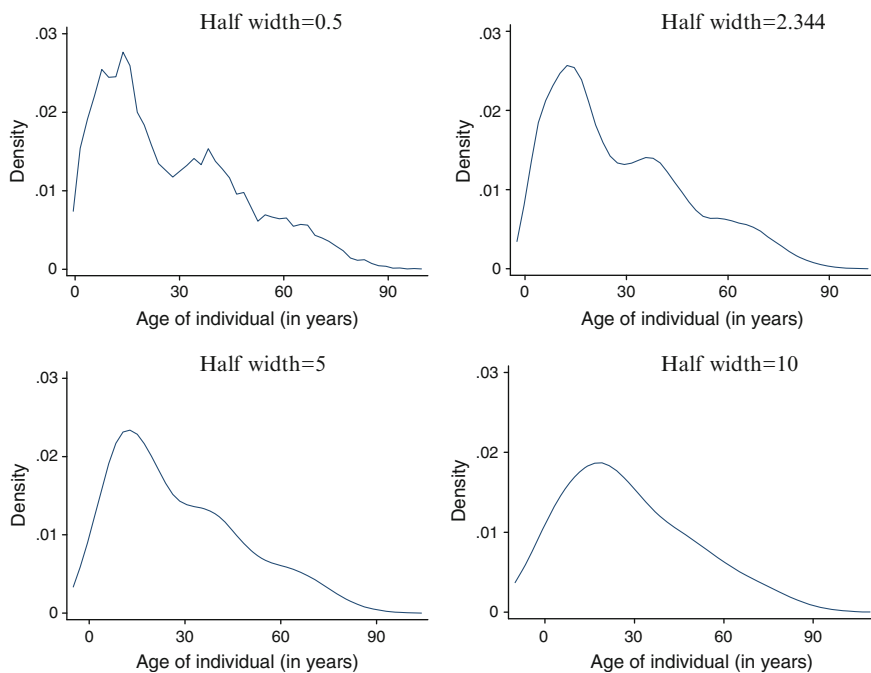


Fig. 1.4 Kernel densities of ages (in years) for individuals surveyed in the Vietnam Living Standards Survey of 1997–1998, with differing half-widths, including the Stata default value of 2.344

kernel is a concave quadratic function with maximum weight at x and zero weights at $x \pm h$. Formally,

$$K_E(z) = \frac{3}{4} \times \left[1 - \frac{1}{5} z^2 \right] / \sqrt{5} \times \mathbf{I}(|z| < \sqrt{5}). \quad (1.8)$$

The principal virtue of the Epanechnikov kernel is that it is the most efficient in minimizing the mean integrated squared error, which is the difference between the true and estimated densities (Stata 2010; Silverman 1986). The Gaussian kernel is also popular, and is given by

$$K_G(z) = \frac{1}{\sqrt{2\pi}} e^{-z^2/2}. \quad (1.9)$$

In fitting a kernel density we have to choose both the kernel function itself and the bandwidth h . It is generally agreed that the important decision concerns the choice of h , just as the choice of bin width is central to the construction of a good histogram. If h is too wide, the kernel density filters the data too much, and potentially valuable information is lost – compare the bottom right panel of Fig. 1.4, which has a wide bandwidth, with the top right panel, where the bandwidth