# Analysing Sociolinguistic Variation

Now in its second edition, this is an invaluable manual for teaching and learning variation analysis, the quantitative study of linguistic variation and change. Written by a leading scholar in the field with over thirty years of experience, it provides an insider's view of the methodology through practical, 'hands-on' advice, including straightforward instructions for conducting analyses using the R programming language, the new gold standard for analysis. It leads readers through each phase of a research study based on data gathered in a sociocultural context, beginning with the selection and sampling of a data source, to hints on successful project design, interview techniques, data management, analysis and interpretation, with systematic procedures provided at each step of the process. This edition has been fully updated, with new insights and explanations in line with recent discoveries in the field, making it essential reading for anyone embarking on their own sociolinguistic research project.

**Sali A. Tagliamonte** is Professor of Linguistics at the University of Toronto, Canada Research Chair in Language Variation and Change, and Fellow of the Royal Society of Canada. Her recent publications include *Making Waves* (2016) and *Teen Talk* (2016). She is also the editor of the book series Studies in Language Variation and Change.

# Analysing Sociolinguistic Variation

## SECOND EDITION

**Sali A. Tagliamonte**

*University of Toronto*

**CAMBRIDGE** UNIVERSITY PRESS

## CAMBRIDGE
### UNIVERSITY PRESS

For Dazzian, Freya, Shaman, Tara, Adrian
And then Emily, Craig, Melissa, Thomas
And now also Aurin, Kieran, Caliope
With love,
*Mum/Salimum/Granna*

# Contents

# Figures

x          **List of Figures**

# Tables

# Preface

The variationist approach to sociolinguistics began during the 1960s, when William Labov, working with Uriel Weinreich and Martin Herzog, developed a theory of language change (Weinreich et al., 1968). Thereafter, Labov continued to advance the method and analysis of language variation and change (e.g. Labov, 1963, 1966, 1969b). In the 1970s, one of Labov's graduate students at the University of Pennsylvania was Shana Poplack. In 1981, Shana became a professor of sociolinguistics at the University of Ottawa's Department of Linguistics, the same year I entered the MA programme. I was fortunate to be Shana's student from 1981 to 1995. We produced many joint publications (e.g. Poplack & Tagliamonte, 1989, 2001), and our work together has had a lasting impact on my research. I also benefited tremendously from the influence of David Sankoff, who always had astute answers to my questions about method and analysis. When I wrote the 2006 edition of this book, everything in it had come directly from what had been passed on from this lineage – training, techniques, insights, knowledge, and sheer passion for the field.

At that time, knowledge and learning in variation analysis had been acquired through word of mouth, from one researcher to the next (see also Guy, 1988:124). It was often noted that the practical details of how to *actually do it* were arcane, largely unwritten, and for the most part, undocumented. That is why the 2006 book was conceived and written. The method needed to be recorded, systematically, thoroughly, and straightforwardly. By 2024 there is a lot of water under the bridge. A variationist approach to language science has changed substantially, and its rigorous, empirical, corpus-based, quantitative methods have spread around the world. Scholars everywhere have been training new generations of variationist linguists, and conferences and workshops focused on many burgeoning new aspects of language variation and change are flourishing.

In completing this second edition, I am indebted to Nathalie Dion, Bridget Jankowski, and Katharina Pabst who worked their way through the text and the code in superb detail, and to Bridget Jankowski and Elena Manzella for checking the final manuscript for me. I have taken everyone's comments into account and then some. My own method has evolved in just this way, changing from one research project to the next, one student or post-doc, one collaborator to the next in my perpetual efforts to do things more efficiently, more usefully, and more transparently (not to mention the joy of working with others!). As far as quantitative savvy is concerned, I have profited from the guidance and wise counsel of many people who are far better at statistics than I am, especially Harald Baayen, John Paolillo, Stefan Gries, Karlien Franco, and Jeremy Needle. Jeremy has been my workshopping

consultant over the past few years and is the expert behind the R-steps for variation linguistics that bring the content up to date with recent developments in the field.

This second edition is an entirely reworked version of the original book. No word, concept or idea has avoided scrutiny and re-evaluation. While the original work enshrined in Weinreich et al. (1968) and then built upon and elaborated by Labov (1982) are fundamental and pervasive, the enrichment and expansion in the field from other sources in the last twenty years should also be apparent and strong, including methods and ideas from the research acumen in variation analysis developing in the EU in historical linguistics, syntax, and cognitive and usage-based linguistics. I have also benefited from the input of many scholars, collaborators, and friends who are too numerous to name. They know who they are. The tracks of their influence are woven indelibly through my life, often in the lines of my CV, but sometimes less conspicuously in citations and acknowledgements, or in the odd social media posts. In sum, this new edition is a thoroughly updated, tried-and-true manual of best practice over the many years I have been working and learning in the ever-blossoming field of language variation and change.