# Multilevel and Longitudinal Modeling Using Stata

## Second Edition

SOPHIA RABE-HESKETH
*University of California, Berkeley*
*Institute of Education, University of London*

ANDERS SKRONDAL
*London School of Econovucs*
*Norwegian Institute of Public Health*

# Contents