# The Workflow of Data Analysis Using Stata

J. SCOTT LONG
*Departments of Sociology and Statistics*
*Indiana University-Bloomington*

# Contents

/