# *Practical Data Science with R*

NINA ZUMEL
JOHN MOUNT

*To our parents*
*Olive and Paul Zumel*
*Peggy and David Mount*

# brief contents

# contents