Jason S. Schwarz • Chris Chapman • Elea McDonnell Feit

# Python for Marketing Research and Analytics

# Contents