2e

# BUSINESS
# ANALYTICS

## COMMUNICATING WITH NUMBERS

Sanjiv Jaggia
Kevin Lertwachara
Alison Kelly
Leida Chen

Mc
Graw
Hill

# CONTENTS