# R

# IN ACTION

Data analysis and graphics with R

Robert I. Kabacoff

# brief contents

# *contents*

appendix B    Customizing the startup environment    406

appendix C    Exporting data from R    408

appendix D    Creating publication-quality output    410

appendix E    Matrix Algebra in R    419

appendix F    Packages used in this book    421

appendix G    Working with large datasets    429

appendix H    Updating an R installation    432

              references    434

              index    438