# Data Mining
## A TUTORIAL-BASED PRIMER

Richard J. Roiger    Michael W. Geatz

# Contents

** *Part IV (Chapters 12, 13, and 14) are available on the CD that accompanies this book.*

# Appendixes