# CONTENTS IN DETAIL

## PART I
## PROLOGUE, AND NEIGHBORHOOD-BASED METHODS

## 1
## REGRESSION MODELS                                                      3

# 6
## TWEAKING THE TREES

# PART III
## METHODS BASED ON LINEAR RELATIONSHIPS

# 8
## PARAMETRIC METHODS     123

# PART IV
# METHODS BASED ON SEPARATING LINES AND PLANES

## 10
## A BOUNDARY APPROACH: SUPPORT VECTOR MACHINES     165

## 11
## LINEAR MODELS ON STEROIDS: NEURAL NETWORKS     185

**PART V
APPLICATIONS**

# 12
# IMAGE CLASSIFICATION                                           199

# 13
# HANDLING TIME SERIES AND TEXT DATA                            211