# DATA MINING WITH R

## *LEARNING WITH CASE STUDIES*

### SECOND EDITION

## Luís Torgo

University of Porto, Portugal

# Contents