# DATA MINING FOR BUSINESS INTELLIGENCE

Concepts, Techniques, and Applications
in Microsoft Office Excel® with XLMiner®

**GALIT SHMUELI**
University of Maryland

**NITIN R. PATEL**
MIT

**PETER C. BRUCE**
statistics.com

# CONTENTS