# DATA MINING

## Concepts, Models, Methods, and Algorithms

### THIRD EDITION

Mehmed Kantardzic

**IEEE** PRESS

WILEY

# CONTENTS