# Data Mining

## Practical Machine Learning Tools and Techniques
## with Java Implementations

## Ian H. Witten

Department of Computer Science
University of Waikato

## Eibe Frank

Department of Computer Science
University of Waikato

**M K**®

# Contents