# Data Mining: Concepts and Techniques

**Second Edition**

Jiawei Han
*University of Illinois at Urbana-Champaign*
Micheline Kamber

# Contents