# PRIVACY-PRESERVING DATA MINING: MODELS AND ALGORITHMS

# PRIVACY-PRESERVING DATA MINING: MODELS AND ALGORITHMS

Edited by

CHARU C. AGGARWAL
IBM T. J. Watson Research Center, Hawthorne, NY 10532

PHILIP S. YU
University of Illinois at Chicago, Chicago, IL 60607

# Contents