# AN INTRODUCTION TO CHEMOINFORMATICS

by

## ANDREW R. LEACH

*GlaxoSmithKline Research and Development,*
*Stevenage, U.K.*

and

## VALERIE J. GILLET

*Department of Information Studies,*
*University of Sheffield, U.K.*

# CONTENTS