# LEARNING TO CLASSIFY TEXT USING SUPPORT VECTOR MACHINES

## Methods, Theory and Algorithms

Thorsten Joachims

# LEARNING TO CLASSIFY TEXT USING SUPPORT VECTOR MACHINES

**Thorsten Joachims**
*Cornell University, U.S.A.*

# Contents