# MACHINE LEARNING IN TRANSLATION

*Peng Wang and David B. Sawyer*

# CONTENTS