# NATURAL LANGUAGE PROCESSING IN THE REAL WORLD

## Text Processing, Analytics, and Classification



## JYOTIKA SINGH
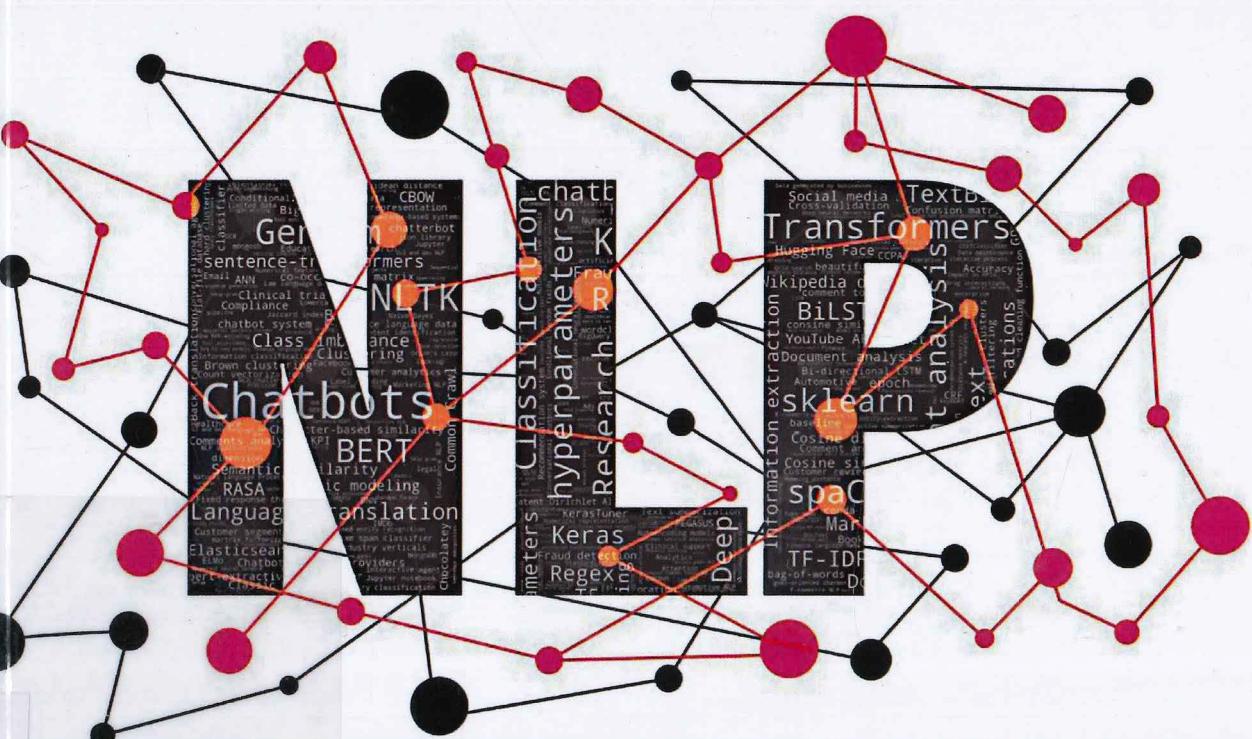
# Natural Language Processing in the Real World

## Text Processing, Analytics, and Classification

Jyotika Singh

# Contents