# Practical Natural Language Processing

## A Comprehensive Guide to Building Real-World NLP Systems

*Sowmya Vajjala, Bodhisattwa Majumder,*
*Anuj Gupta, and Harshit Surana*

# Table of Contents

# Part II.    Essentials

# Part III.   Applied

## Part IV.  Bringing It All Together