# Natural Language Processing in Action

*Understanding, analyzing, and generating text with Python*

HOBSON LANE
COLE HOWARD
HANNES MAX HAPKE

MANNING
SHELTER ISLAND

# Table of Contents

# Part 2. Deeper learning (neural networks)

# 7 Getting words in order with convolutional neural networks (CNNs)

# Part 3. Getting real (real-world NLP challenges)

## *11 Information extraction (named entity extraction and question answering)*

## *12 Getting chatty (dialog engines)*

## Glossary

## Index