# Explainable
# Natural Language Processing

# Synthesis Lectures on Human Language Technologies

### Editor
**Graeme Hirst,** *University of Toronto*

Synthesis Lectures on Human Language Technologies is edited by Graeme Hirst of the University of Toronto. The series consists of 50- to 150-page monographs on topics relating to natural language processing, computational linguistics, information retrieval, and spoken language understanding. Emphasis is on important new techniques, on new applications, and on topics that combine two or more HLT subfields.

### Explainable Natural Language Processing
Anders Søgaard
2021

### Finite-State Text Processing
Kyle Gorman and Richard Sproat
2021

### Semantic Relations Between Nominals, Second Edition
Vivi Nastase, Stan Szpakowicz, Preslav Nakov, and Diarmuid Ó Séagdha
2021

### Embeddings in Natural Language Processing: Theory and Advances in Vector Representations of Meaning
Mohammad Taher Pilehvar and Jose Camacho-Collados
2020

### Conversational AI: Dialogue Systems, Conversational Agents, and Chatbots
Michael McTear
2020

### Natural Language Processing for Social Media, Third Edition
Anna Atefeh Farzindar and Diana Inkpen
2020

### Statistical Significance Testing for Natural Language Processing
Rotem Dror, Lotem Peled, Segev Shlomov, and Roi Reichart
2020

Explainable Natural Language Processing

Anders Søgaard

# Explainable
# Natural Language Processing

Anders Søgaard
University of Copenhagen

## ABSTRACT

This book presents a taxonomy framework and survey of methods relevant to explaining the decisions and analyzing the inner workings of Natural Language Processing (NLP) models. The book is intended to provide a snapshot of Explainable NLP, though the field continues to rapidly grow. The book is intended to be both readable by first-year M.Sc. students and interesting to an expert audience. The book opens by motivating a focus on providing a consistent taxonomy, pointing out inconsistencies and redundancies in previous taxonomies. It goes on to present (i) a taxonomy or framework for thinking about how approaches to explainable NLP relate to one another; (ii) brief surveys of each of the classes in the taxonomy, with a focus on methods that are relevant for NLP; and (iii) a discussion of the inherent limitations of some classes of methods, as well as how to best evaluate them. Finally, the book closes by providing a list of resources for further research on explainability.

## KEYWORDS

# Contents

# Acknowledgments