

Natural Language Processing for Social Media

Second Edition

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.

Natural Language Processing for Social Media, Second Edition

Atefeh Farzindar and Diana Inkpen
2017

Automatic Text Simplification

Horacio Saggion
2017

Neural Network Methods for Natural Language Processing

Yoav Goldberg
2017

Syntax-based Statistical Machine Translation

Philip Williams, Rico Senrich, Matt Post, and Philipp Koehn
2016

Domain-Sensitive Temporal Tagging

Jannik Strötgen and Michael Gertz
2016

Linked Lexical Knowledge Bases: Foundations and Applications

Iryna Gurevych, Judith Eckle-Kohler, and Michael Matuschek
2016

Bayesian Analysis in Natural Language Processing

Shay Cohen
2016

Metaphor: A Computational Perspective

Tony Veale, Ekaterina Shutova, and Beata Beigman Klebanov
2016

Grammatical Inference for Computational Linguistics

Jeffrey Heinz, Colin de la Higuera, and Menno van Zaanen
2015

Automatic Detection of Verbal Deception

Eileen Fitzpatrick, Joan Bachenko, and Tommaso Fornaciari
2015

Natural Language Processing for Social Media

Atefeh Farzindar and Diana Inkpen
2015

Semantic Similarity from Natural Language and Ontology Analysis

Sébastien Harispe, Sylvie Ranwez, Stefan Janaqi, and Jacky Montmain
2015

Learning to Rank for Information Retrieval and Natural Language Processing, Second Edition

Hang Li
2014

Ontology-Based Interpretation of Natural Language

Philipp Cimiano, Christina Unger, and John McCrae
2014

Automated Grammatical Error Detection for Language Learners, Second Edition

Claudia Leacock, Martin Chodorow, Michael Gamon, and Joel Tetreault
2014

Web Corpus Construction

Roland Schäfer and Felix Bildhauer
2013

Recognizing Textual Entailment: Models and Applications

Ido Dagan, Dan Roth, Mark Sammons, and Fabio Massimo Zanzotto
2013

Linguistic Fundamentals for Natural Language Processing: 100 Essentials from Morphology and Syntax

Emily M. Bender
2013

Semi-Supervised Learning and Domain Adaptation in Natural Language Processing

Anders Søgaard

2013

Semantic Relations Between Nominals

Vivi Nastase, Preslav Nakov, Diarmuid Ó Séaghdha, and Stan Szpakowicz

2013

Computational Modeling of Narrative

Inderjeet Mani

2012

Natural Language Processing for Historical Texts

Michael Piotrowski

2012

Sentiment Analysis and Opinion Mining

Bing Liu

2012

Discourse Processing

Manfred Stede

2011

Bitext Alignment

Jörg Tiedemann

2011

Linguistic Structure Prediction

Noah A. Smith

2011

Learning to Rank for Information Retrieval and Natural Language Processing

Hang Li

2011

Computational Modeling of Human Language Acquisition

Afra Alishahi

2010

Introduction to Arabic Natural Language Processing

Nizar Y. Habash

2010

Cross-Language Information Retrieval

Jian-Yun Nie

2010

Automated Grammatical Error Detection for Language Learners

Claudia Leacock, Martin Chodorow, Michael Gamon, and Joel Tetreault
2010

Data-Intensive Text Processing with MapReduce

Jimmy Lin and Chris Dyer
2010

Semantic Role Labeling

Martha Palmer, Daniel Gildea, and Nianwen Xue
2010

Spoken Dialogue Systems

Kristiina Jokinen and Michael McTear
2009

Introduction to Chinese Natural Language Processing

Kam-Fai Wong, Wenjie Li, Ruifeng Xu, and Zheng-sheng Zhang
2009

Introduction to Linguistic Annotation and Text Analytics

Graham Wilcock
2009

Dependency Parsing

Sandra Kübler, Ryan McDonald, and Joakim Nivre
2009

Statistical Language Models for Information Retrieval

ChengXiang Zhai
2008

© Springer Nature Switzerland AG 2022
Reprint of original edition © Morgan & Claypool 2018

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted in any form or by any means—electronic, mechanical, photocopy, recording, or any other except for brief quotations in printed reviews, without the prior permission of the publisher.

Natural Language Processing for Social Media, Second Edition
Atefeh Farzindar and Diana Inkpen

ISBN: 978-3-031-01039-2	paperback
ISBN: 978-3-031-02167-1	ebook
ISBN: 978-3-031-00178-9	hardcover

DOI 10.1007/978-3-031-02167-1

A Publication in the Springer series
SYNTHESIS LECTURES ON HUMAN LANGUAGE TECHNOLOGIES

Lecture #38
Series Editor: Graeme Hirst, *University of Toronto*
Series ISSN
Print 1947-4040 Electronic 1947-4059

Natural Language Processing for Social Media

Second Edition

Atefeh Farzindar
University of Southern California

Diana Inkpen
University of Ottawa

SYNTHESIS LECTURES ON HUMAN LANGUAGE TECHNOLOGIES #38

ABSTRACT TO THE SECOND EDITION

In recent years, online social networking has revolutionized interpersonal communication. The newer research on language analysis in social media has been increasingly focusing on the latter's impact on our daily lives, both on a personal and a professional level. Natural language processing (NLP) is one of the most promising avenues for social media data processing. It is a scientific challenge to develop powerful methods and algorithms which extract relevant information from a large volume of data coming from multiple sources and languages in various formats or in free form. We discuss the challenges in analyzing social media texts in contrast with traditional documents.

Research methods in information extraction, automatic categorization and clustering, automatic summarization and indexing, and statistical machine translation need to be adapted to a new kind of data. This book reviews the current research on NLP tools and methods for processing the non-traditional information from social media data that is available in large amounts (big data), and shows how innovative NLP approaches can integrate appropriate linguistic information in various fields such as social media monitoring, healthcare, business intelligence, industry, marketing, and security and defence.

We review the existing evaluation metrics for NLP and social media applications, and the new efforts in evaluation campaigns or shared tasks on new datasets collected from social media. Such tasks are organized by the Association for Computational Linguistics (such as SemEval tasks) or by the National Institute of Standards and Technology via the Text REtrieval Conference (TREC) and the Text Analysis Conference (TAC). In the concluding chapter, we discuss the importance of this dynamic discipline and its great potential for NLP in the coming decade, in the context of changes in mobile technology, cloud computing, virtual reality, and social networking.

In this second edition, we have added information about recent progress in the tasks and applications presented in the first edition. We discuss new methods and their results. The number of research projects and publications that use social media data is constantly increasing due to continuously growing amounts of social media data and the need to automatically process them. We have added 85 new references to the more than 300 references from the first edition. Besides updating each section, we have added a new application (digital marketing) to the section on media monitoring and we have augmented the section on healthcare applications with an extended discussion of recent research on detecting signs of mental illness from social media.

KEYWORDS

social media, social networking, natural language processing, social computing, big data, semantic analysis

This effort is dedicated to my husband, Massoud, and to my daughters, Tina and Amanda, who are just about the best children a mom could hope for: happy, loving, and fun to be with.

– Atefeh Farzindar

To my wonderful husband, Nicu, with whom I can climb any mountain, and to our sweet baby daughter Nicoleta.

– Diana Inkpen

Contents

	Preface to the Second Edition	xv
	Acknowledgments	xix
1	Introduction to Social Media Analysis	1
1.1	Introduction	1
1.2	Social Media Applications	6
1.2.1	Cross-language Document Analysis in Social Media Data	7
1.2.2	Real-world Applications	7
1.3	Challenges in Social Media Data	8
1.4	Semantic Analysis of Social Media	12
1.5	Summary	13
2	Linguistic Pre-processing of Social Media Texts	15
2.1	Introduction	15
2.2	Generic Adaptation Techniques for NLP Tools	17
2.2.1	Text Normalization	17
2.2.2	Re-training NLP Tools for Social Media Texts	19
2.3	Tokenizers	20
2.4	Part-of-speech Taggers	22
2.5	Chunkers and Parsers	24
2.6	Named Entity Recognizers	27
2.7	Existing NLP Toolkits for English and Their Adaptation	29
2.8	Multi-linguality and Adaptation to Social Media Texts	31
2.8.1	Language Identification	31
2.8.2	Dialect Identification	33
2.9	Summary	39
3	Semantic Analysis of Social Media Texts	41
3.1	Introduction	41
3.2	Geo-location Detection	41

3.2.1	Mapping Social Media Information on Maps	42
3.2.2	Readily Available Geo-location Information	42
3.2.3	Geo-location based on Network Infrastructure	42
3.2.4	Geo-location based on the Social Network Structure	43
3.2.5	Content-based Location Detection	43
3.2.6	Evaluation Measures for Geo-location Detection	47
3.3	Entity Linking and Disambiguation	49
3.3.1	Detecting Entities and Linked Data	50
3.3.2	Evaluation Measures for Entity Linking	52
3.4	Opinion Mining and Emotion Analysis	53
3.4.1	Sentiment Analysis	53
3.4.2	Emotion Analysis	56
3.4.3	Sarcasm Detection	58
3.4.4	Evaluation Measures for Opinion and Emotion Classification	59
3.5	Event and Topic Detection	61
3.5.1	Specified vs. Unspecified Event Detection	61
3.5.2	New vs. Retrospective Events	67
3.5.3	Emergency Situation Awareness	68
3.5.4	Evaluation Measures for Event Detection	69
3.6	Automatic Summarization	69
3.6.1	Update Summarization	71
3.6.2	Network Activity Summarization	71
3.6.3	Event Summarization	72
3.6.4	Opinion Summarization	72
3.6.5	Evaluation Measures for Summarization	73
3.7	Machine Translation	74
3.7.1	Adapting Phrase-based Machine Translation to Normalize Medical Terms	75
3.7.2	Translating Government Agencies' Tweet Feeds	76
3.7.3	Hashtag Occurrence, Layout, and Translation	77
3.7.4	Machine Translation for Arabic Social Media	80
3.7.5	Evaluation Measures for Machine Translation	82
3.8	Summary	82
4	Applications of Social Media Text Analysis	85
4.1	Introduction	85
4.2	Healthcare Applications	85

4.3	Financial Applications	92
4.4	Predicting Voting Intentions	95
4.5	Media Monitoring	97
4.6	Security and Defense Applications	99
4.7	Disaster Response Applications	102
4.8	NLP-based User Modeling	104
4.9	Applications for Entertainment	109
4.10	NLP-based Information Visualization for Social Media	111
4.11	Government Communication	111
4.12	Summary	112
5	Data Collection, Annotation, and Evaluation	113
5.1	Introduction	113
5.2	Discussion on Data Collection and Annotation	113
5.3	Spam and Noise Detection	114
5.4	Privacy and Democracy in Social Media	116
5.5	Evaluation Benchmarks	117
5.6	Summary	119
6	Conclusion and Perspectives	121
6.1	Conclusion	121
6.2	Perspectives	121
A	TRANSLI: a Case Study for Social Media Analytics and Monitoring	125
A.1	TRANSLI architecture	125
A.2	User Interface	126
	Glossary	131
	Bibliography	133
	Authors' Biographies	173
	Index	175

Preface to the Second Edition

This book presents the state-of-the-art in research and empirical studies in the field of Natural Language Processing (NLP) for the semantic analysis of social media data. Because the field is continuously growing, this second edition adds information about recently proposed methods and their results for the tasks and applications that we covered in the first edition.

Over the past few years, online social networking sites have revolutionized the way we communicate with individuals, groups, and communities, and altered everyday practices. The unprecedented volume and variety of user-generated content and the user interaction network constitute new opportunities for understanding social behavior and building socially intelligent systems.

Much of the research on social networks and the mining of the social web is based on graph theory. That is apt because a social structure is made up of a set of social actors and a set of the dyadic ties between these actors. We believe that the graph mining methods for structure and information diffusion or influence spread in social networks need to be combined with the content analysis of social media. This provides the opportunity for new applications that use the information publicly available as a result of social interactions. Adapted classic NLP methods can partially solve the problem of social media content analysis focusing on the posted messages. When we receive a text of less than 10 characters, including an emoticon and a heart, we understand it and even respond to it! It is impossible to use NLP methods to process this type of document, but there is a logical message in social media data based on which two people can communicate. The same logic dominates worldwide, and people from all over the world use it to share and communicate with each other. There is a new and challenging language for NLP.

We believe that we need new theories and algorithms for semantic analysis of social media data, as well as a new way of approaching the big data processing. By semantic analysis, in this book, we mean the linguistic processing of the social media messages enhanced with semantics, and possibly also combining this with the structure of the social networks. We actually use the term in a more general sense to refer to applications that do intelligent processing of social media texts and meta-data. Some applications could access very large amounts of data; therefore, the algorithms need to be adapted to be able process data (big data) in an online fashion and without necessarily storing all the data.

This motivated us to give two tutorials: *Applications of Social Media Text Analysis* at EMNLP 2015¹ and *Natural Language Processing for Social Media* at the 29th Canadian Confer-

¹http://www.emnlp2015.org/tutorials/3/3_OptionalAttachment.pdf

<https://www.cs.cmu.edu/~ark/EMNLP-2015/proceedings/EMNLP-Tutorials/pdf/EMNLP-Tutorials06.pdf>

ence on Artificial Intelligence (AI 2016).² We also organized several workshops on this topic, Semantic Analysis in Social Networks (SASM 2012)³ and Language Analysis in Social Media (LASM 2013⁴ and LASM 2014⁵), in conjunction with conferences organized by the Association for Computational Linguistics⁶ (ACL, EACL, and NAACL-HLT).

Our goal was to reflect a wide range of research and results in the analysis of language with implications for fields such as NLP, computational linguistics, sociolinguistics, and psycholinguistics. Our workshops invited original research on all topics related to the analysis of language in social media, including the following topics.

- What do people talk about on social media?
- How do they express themselves?
- Why do they post on social media?
- How do language and social network properties interact?
- Natural language processing techniques for social media analysis.
- Semantic Web/ontologies/domain models to aid in understanding social data.
- Characterizing participants via linguistic analysis.
- Language, social media, and human behavior.

There were several other workshops on similar topics, for example, the *Making Sense of Microposts* (#Microposts)⁷ workshop series in conjunction with the World Wide Web Conference 2012–2016. These workshops focused in particular on short informal texts that are published without much effort (such as tweets, Facebook shares, Instagram-like shares, and Google+ messages). There has been another series of workshops on Natural Language Processing for Social Media (SocialNLP) since 2013, with SocialNLP 2017 offered in conjunction with EACL 2017⁸ and IEEE BigData 2017.⁹

The **intended audience** of this book is researchers who are interested in developing tools and applications for automatic analysis of social media texts. We assume that the readers have basic knowledge in the area of natural language processing and machine learning. We hope that this book will help the readers better understand computational linguistics and social media

²<http://aigicrv.org/2016/>

³<https://aclweb.org/anthology/W/W12/#2100>

⁴<https://aclweb.org/anthology/W/W13/#1100>

⁵<https://aclweb.org/anthology/W/W14/#1300>

⁶<http://www.aclweb.org/>

⁷<http://microposts2016.seas.upenn.edu/>

⁸<http://eac12017.org/>

⁹<http://cci.drexel.edu/bigdata/bigdata2017/>

analysis, in particular text mining techniques and NLP applications (such as summarization, localization detection, sentiment and emotion analysis, topic detection, and machine translation) designed specifically for social media texts.

Atefeh Farzindar and Diana Inkpen
December 2017

Acknowledgments

This book would not have been possible without the hard work of many people. We would like to thank our colleagues at NLP Technologies Inc., the NLP research group at the University of Ottawa, and our students James Webb and Ruining Liu from the University of Southern California. We would like to thank in particular Prof. Stan Szpakowicz from the University of Ottawa for his comments on the draft of the book, and two anonymous reviewers for their useful suggestions for revisions and additions. We thank Prof. Graeme Hirst at the University of Toronto and Michael Morgan from Morgan & Claypool Publishers for their continuous encouragement.

Atefeh Farzindar and Diana Inkpen
December 2017