



MORGAN & CLAYPOOL PUBLISHERS

# Natural Language Processing for Social Media

Third Edition

Anna Atefeh Farzindar  
Diana Inkpen

***SYNTHESIS LECTURES ON  
HUMAN LANGUAGE TECHNOLOGIES***

Graeme Hirst, *Series Editor*

# Natural Language Processing for Social Media

Third 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, 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

## Deep Learning Approaches to Text Production

Shashi Narayan and Claire Gardent

2020

## Linguistic Fundamentals for Natural Language Processing II: 100 Essentials from Semantics and Pragmatics

Emily M. Bender and Alex Lascarides

2019

## Cross-Lingual Word Embeddings

Anders Søgaard, Ivan Vulić, Sebastian Ruder, Manaal Faruqui

2019

## Bayesian Analysis in Natural Language Processing, Second Edition

Shay Cohen

2019

## Argumentation Mining

Manfred Stede and Jodi Schneider

2018

- [Quality Estimation for Machine Translation](#)  
Lucia Specia, Carolina Scarton, and Gustavo Henrique Paetzold  
2018
- [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 Sennrich, 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



Copyright © 2020 by Morgan & Claypool

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, Third Edition

Anna Atefeh Farzindar and Diana Inkpen

[www.morganclaypool.com](http://www.morganclaypool.com)

ISBN: 9781681738116      paperback

ISBN: 9781681738123      ebook

ISBN: 9781681738147      epub

ISBN: 9781681738130      hardcover

DOI 10.2200/S00999ED3V01Y202003HLT046

A Publication in the Morgan & Claypool Publishers series

*SYNTHESIS LECTURES ON HUMAN LANGUAGE TECHNOLOGIES*

Lecture #46

Series Editor: Grame Hirst, *University of Toronto*

Series ISSN

Print 1947-4040    Electronic 1947-4059

Cover art illustration by Anna Atefeh Farzindar.

# Natural Language Processing for Social Media

Third Edition

Anna Atefeh Farzindar  
University of Southern California

Diana Inkpen  
University of Ottawa

*SYNTHESIS LECTURES ON HUMAN LANGUAGE TECHNOLOGIES #46*



MORGAN & CLAYPOOL PUBLISHERS

## ABSTRACT

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 that extract relevant information from a large volume of data coming from multiple sources and languages in various formats or in free form. This book will 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, and it shows how innovative NLP approaches can integrate appropriate linguistic information in various fields such as social media monitoring, health care, and business intelligence. The book further covers 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), the National Institute of Standards and Technology via the Text REtrieval Conference (TREC) and the Text Analysis Conference (TAC), or the Conference and Labs of the Evaluation Forum (CLEF).

In this third edition of the book, the authors added information about recent progress in NLP for social media applications, including more about the modern techniques provided by deep neural networks (DNNs) for modeling language and analyzing social media data.

## KEYWORDS

social media, social networking, natural language processing, social computing, big data, semantic analysis, artificial intelligence, deep learning

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

*– Anna Atefeh Farzindar*

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

*– Diana Inkpen*



# Contents

<b>List of Figures</b> .....	<b>xvii</b>
<b>List of Tables</b> .....	<b>xix</b>
<b>Preface</b> .....	<b>xxi</b>
<b>Acknowledgments</b> .....	<b>xxv</b>
<b>1 Introduction to Social Media Analysis</b> .....	<b>1</b>
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 Deep Learning techniques for Social Media Data .....	7
1.2.3 Real-world Applications .....	8
1.3 Challenges in Social Media Data .....	9
1.4 Semantic Analysis of Social Media .....	12
1.5 Summary .....	13
<b>2 Linguistic Pre-processing of Social Media Texts</b> .....	<b>15</b>
2.1 Introduction .....	15
2.2 Generic Adaptation Techniques for NLP Tools .....	17
2.2.1 Text Normalization .....	18
2.2.2 Re-training NLP Tools for Social Media Texts .....	20
2.3 Tokenizers .....	21
2.4 Part-of-speech Taggers .....	22
2.5 Chunkers and Parsers .....	25
2.6 Named Entity Recognizers .....	29
2.7 Existing NLP Toolkits for English and Their Adaptation .....	31
2.8 Multi-linguality and Adaptation to Social Media Texts .....	32
2.8.1 Language Identification .....	32
2.8.2 Dialect Identification .....	35
2.9 Summary .....	41

<b>3</b>	<b>Semantic Analysis of Social Media Texts</b>	<b>43</b>
3.1	Introduction	43
3.2	Geo-location Detection	43
3.2.1	Mapping Social Media Information on Maps	44
3.2.2	Readily Available Geo-location Information	44
3.2.3	Geo-location based on Network Infrastructure	44
3.2.4	Geo-location based on the Social Network Structure	45
3.2.5	Content-based Location Detection	45
3.2.6	Evaluation Measures for Geo-location Detection	49
3.3	Entity Linking and Disambiguation	51
3.3.1	Detecting Entities and Linked Data	52
3.3.2	Evaluation Measures for Entity Linking	55
3.4	Opinion Mining and Emotion Analysis	55
3.4.1	Sentiment Analysis	55
3.4.2	Emotion Analysis	58
3.4.3	Sarcasm Detection	60
3.4.4	Evaluation Measures for Opinion and Emotion Classification	61
3.5	Event and Topic Detection	62
3.5.1	Specified vs. Unspecified Event Detection	62
3.5.2	New vs. Retrospective Events	69
3.5.3	Emergency Situation Awareness	70
3.5.4	Evaluation Measures for Event Detection	71
3.6	Automatic Summarization	71
3.6.1	Update Summarization	73
3.6.2	Network Activity Summarization	73
3.6.3	Event Summarization	74
3.6.4	Opinion Summarization	75
3.6.5	Keyphrase Generation	76
3.6.6	Evaluation Measures for Summarization	76
3.7	Machine Translation	77
3.7.1	Neural Machine Translation	78
3.7.2	Adapting Phrase-based Machine Translation to Normalize Medical Terms	79
3.7.3	Translating Government Agencies' Tweet Feeds	79
3.7.4	Hashtag Occurrence, Layout, and Translation	81
3.7.5	Machine Translation for Arabic Social Media	84
3.7.6	Evaluation Measures for Machine Translation	85

3.8	Summary	86
<b>4</b>	<b>Applications of Social Media Text Analysis</b>	<b>87</b>
4.1	Introduction	87
4.2	Healthcare Applications	87
4.3	Financial Applications	96
4.4	Predicting Voting Intentions	99
4.5	Media Monitoring	101
4.6	Security and Defense Applications	104
4.7	Disaster Response Applications	107
4.8	NLP-based User Modeling	109
4.9	Applications for Entertainment	115
4.10	NLP-based Information Visualization for Social Media	117
4.11	Government Communication	117
4.12	Rumor Detection	118
4.13	Recommender systems	119
4.14	Preventing Sexual Harassment	120
4.15	Summary	120
<b>5</b>	<b>Data Collection, Annotation, and Evaluation</b>	<b>121</b>
5.1	Introduction	121
5.2	Discussion on Data Collection and Annotation	121
5.3	Spam and Noise Detection	122
5.4	Privacy and Democracy in Social Media	125
5.5	Evaluation Benchmarks	126
5.6	Summary	128
<b>6</b>	<b>Conclusion and Perspectives</b>	<b>129</b>
6.1	Conclusion	129
6.2	Perspectives	129
<b>A</b>	<b>TRANSLI: a Case Study for Social Media Analytics and Monitoring</b>	<b>133</b>
A.1	TRANSLI architecture	133
A.2	User Interface	134



**Glossary** ..... 139

**Bibliography** ..... 141

**Authors' Biographies** ..... 191

**Index** ..... 193

# List of Figures

1.1	Social networks ranked by the number of active users as of January 2014 (in millions) provided by Statista. . . . .	3
1.2	Number of monthly active Facebook users from the third quarter of 2008 to the first quarter of 2014 (in millions) provided by Statista. . . . .	3
1.3	Number of LinkedIn members from the first quarter of 2009 to the first quarter of 2014 (in millions) provided by Statista. . . . .	4
1.4	A framework for semantic analysis in social media, where NLP tools transform the data into intelligence. . . . .	6
2.1	Methodology for tweet normalization. The dotted horizontal line separates the two steps (detecting the text to be normalized and applying normalization rules) [Akhtar et al., 2015]. . . . .	19
2.2	Taxonomy of normalization edits [Baldwin and Li, 2015]. . . . .	20
2.3	Arabic dialects distribution and variation across Asia and Africa [Sadat et al., 2014a]. . . . .	36
2.4	Division of Arabic dialects in six groups/divisions [Sadat et al., 2014a]. . . . .	36
2.5	Accuracies on the character-based $n$ -gram Markov language models for 18 countries [Sadat et al., 2014a]. . . . .	38
2.6	Accuracies on the character-based $n$ -gram Markov language models for the six divisions/groups [Sadat et al., 2014a]. . . . .	39
2.7	Accuracies on the character-based $n$ -gram Naïve Bayes classifiers for 18 countries [Sadat et al., 2014a]. . . . .	40
2.8	Accuracies on the character-based $n$ -gram Naïve Bayes classifiers for the six divisions/groups [Sadat et al., 2014a]. . . . .	41
3.1	Example of a pair of tweets extracted from the bilingual feed pair Health Canada/Santé Canada, after tokenization. . . . .	81
3.2	An original tweet with hashtags in its three possible regions. . . . .	82
4.1	Examples of annotated social media posts discussing ADRs [Nikfarjam et al., 2015]. . . . .	89

4.2	The DeepHealthMiner neural net architecture [Nikfarjam, 2016]. . . . .	90
4.3	SVM-based text mining procedure for impact management [Schniederjans et al., 2013]. . . . .	99
A.1	TRANSLI Social Media Analytics and monitoring module architecture. . . .	134
A.2	TRANSLI user interface for event creation module. . . . .	135
A.3	TRANSLI user interface for event browsing module. . . . .	135
A.4	TRANSLI user interface to present an event. Components are identified with their IDs. . . . .	137

# List of Tables

1.1	Social media platforms and their characteristics . . . . .	2
2.1	Three examples of Twitter texts . . . . .	18
2.2	Examples of tokenization . . . . .	21
2.3	Penn TreeBank tagset . . . . .	23
2.4	POS tagset from Gimpel et al. [2011] . . . . .	26
2.5	Example of tweet parsed with the TweepoParser . . . . .	28
3.1	An example of annotation with the true location [Inkpen et al., 2015] . . . . .	49
3.2	Classification accuracies for user location detection on the Eisenstein dataset [Liu and Inkpen, 2015] . . . . .	50
3.3	Mean error distance of predictions on the Eisenstein dataset [Liu and Inkpen, 2015] . . . . .	50
3.4	Results for user location prediction on the Roller dataset [Liu and Inkpen, 2015] . . . . .	50
3.5	Performance of the classifiers trained on different features for cities [Inkpen et al., 2015] . . . . .	51
3.6	Classification results for emotion classes and non-emotion by Ghazi et al. [2014] . . . . .	63
3.7	Accuracy of the mood classification by Keshtkar and Inkpen [2012] . . . . .	63
3.8	Statistics on hashtag use in the aligned bilingual corpus [Gotti et al., 2014] . . . . .	82
3.9	Distribution of hashtags in epilogues and prologues [Gotti et al., 2014] . . . . .	82
3.10	Percentage of unknown hashtags to English and French vocabularies of the Hansard corpus [Gotti et al., 2014] . . . . .	83
3.11	Percentage of unknown hashtags to “standard” English and French vocabularies, after automatic segmentation of multiword hashtags into simple words [Gotti et al., 2014] . . . . .	83
3.12	Translation performance obtained by Gotti et al. [2014] . . . . .	86



# Preface

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 third edition adds information about recently proposed methods and their results for the tasks and applications that we covered in the first and second editions.

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 research work 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, 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 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 three tutorials on *Applications of Social Media Text Analysis* at EMNLP 2015<sup>1</sup>, on *Natural Language Processing for Social Media* at the 29th Canadian Con-

<sup>1</sup>[http://www.emnlp2015.org/tutorials/3/3\\_OptionalAttachment.pdf](http://www.emnlp2015.org/tutorials/3/3_OptionalAttachment.pdf)  
<https://www.cs.cmu.edu/~ark/EMNLP-2015/proceedings/EMNLP-Tutorials/pdf/EMNLP-Tutorials06.pdf>

ference on Artificial Intelligence (AI 2016)<sup>2</sup>, and on *How Natural Language Processing Helps Uncover Social Media Insights* at the 33rd International Florida Artificial Intelligence Research Society Conference (FLAIRS 2020). Also on this topic, we organized several workshops (Semantic Analysis in Social Networks (SASM 2012)<sup>3</sup>, Language Analysis in Social Media (LASM 2013<sup>4</sup>, and LASM 2014<sup>5</sup>) in conjunction with conferences organized by the Association for Computational Linguistics<sup>6</sup> (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)<sup>7</sup> workshop series in conjunction with the World Wide Web Conference 2012 to 2016. These workshops focused in particular on short informal texts that are published without much effort (such as tweets, Facebook shares, Instagram-like shares, Google+ messages). There has been another series of Workshops on Natural Language Processing for Social Media (SocialNLP) since 2013. For example, SocialNLP 2017 was in conjunction with EACL 2017<sup>8</sup> and IEEE BigData 2017<sup>9</sup>, and SocialNLP 2020 had two editions, one in conjunction with TheWebConf 2020 and one in conjunction with ACL 2020<sup>10</sup>.

The **intended audience** of this book is researchers that are interested in developing tools and applications for automatic analysis social of 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 analysis, in particular text mining techniques and NLP applications (such as summarization,

<sup>2</sup><http://aigicrv.org/2016/>

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

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

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

<sup>6</sup><http://www.aclweb.org/>

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

<sup>8</sup><http://eacl2017.org/>

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

<sup>10</sup><https://sites.google.com/site/socialnlp2020/>

localization detection, sentiment and emotion analysis, topic detection and machine translation) designed specifically for social media texts.

Besides updating each section in this third edition, we added a new section on keyphrase generation from social media messages and one on neural machine translation in Chapter 3 and three new applications in Chapter 4: rumor detection, recommender systems for social media, and preventing sexual harassment. We discuss the new methods and their results. The number of research projects and publications that use social media data is constantly increasing. Finally, we added more than 50 new references to the approximately 400 references from the second edition.

Anna Atefeh Farzindar and Diana Inkpen  
March 2020





# Acknowledgments

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

Anna Atefeh Farzindar and Diana Inkpen  
March 2020



## CHAPTER 1

# Introduction to Social Media Analysis

## 1.1 INTRODUCTION

Social media is a phenomenon that has recently expanded throughout the world and quickly attracted billions of users. This form of electronic communication through social networking platforms allows users to generate its content and share it in various forms of information, personal words, pictures, audio, and videos. Therefore, social computing is formed as an emerging area of research and development that includes a wide range of topics such as Web semantics, artificial intelligence, natural language processing, network analysis, and Big Data analytics.

Over the past few years, online social networking sites (Facebook, Twitter, YouTube, Flickr, MySpace, LinkedIn, Metacafe, Vimeo, etc.) have revolutionized the way we communicate with individuals, groups, and communities, and have altered everyday practices [Boyd and Ellison, 2007].

The broad categories of social media platforms are: content-sharing sites, forums, blogs, and microblogs. On content sharing sites (such as Facebook, Instagram, Foursquare, Flickr, YouTube) people exchange information, messages, photos, videos, or other types of content. On Web user forums (such as StackOverflow, CNET forums, Apple Support) people post specialized information, questions, or answers. Blogs (such as Gizmodo, Mashable, Boing Boing, and many more) allow people to post messages and other content and to share information and opinions. Micro-blogs (such as Twitter, Sina Weibo, Tumblr) are limited to short texts for sharing information and opinions. The modalities of sharing content in order: posts; comments to posts; explicit or implicit connections to build social networks (friend connections, followers, etc.); cross-posts and user linking; social tagging; likes/favorites/starring/voting/rating/etc.; author information; and linking to user profile features.<sup>1</sup> In Table 1.1, we list more details about social media platforms and their characteristics and types of content shared [Barbier et al., 2013].

Social media statistics for January 2014 have shown that Facebook has grown to more than 1 billion active users, adding more than 200 million users in a single year. Statista,<sup>2</sup> the world's largest statistics portal, announced the ranking for social networks based on the number of active users. As presented in Figure 1.1, the ranking shows that Qzone took second place

<sup>1</sup><http://people.eng.unimelb.edu.au/tbaldwin/pubs/starsem2014.pdf>

<sup>2</sup><http://www.statista.com/>

## 2 1. INTRODUCTION TO SOCIAL MEDIA ANALYSIS

Table 1.1: Social media platforms and their characteristics

Type	Characteristics	Examples
Social Networks	A social networking website allows the user to build a web page and connect with a friend or other acquaintance in order to share user-generated content.	MySpace, Facebook, LinkedIn, Meetup, Google Plus+
Blogs and Blog Comments	A blog is an online journal where the blogger can create the content and display it in reverse chronological order. Blogs are generally maintained by a person or a community. Blog comments are posts by users attached to blogs or online newspaper posts.	Huffington Post, Business Insider, Engadget, and online journals
Microblogs	A microblog is similar to a blog but has a limited content.	Twitter, Tumblr, Sina Weibo, Plurk
Forums	An online forum is a place for members to discuss a topic by posting messages.	Online Discussion Communities, phpBB Developer Forum, Raising Children Forum
Social Bookmarks	Services that allow users to save, organize, and search links to various websites, and to share their bookmarks of Web pages.	Delicious, Pinterest, Google Bookmarks
Wikis	These websites allow people to collaborate and add content or edit the information on a community-based database.	Wikipedia, Wikitravel, Wikihow
Social News	Social news encourage their community to submit news stories, or to vote on the content and share it.	Digg, Slashdot, Reddit
Media Sharing	A website that enables users to capture videos and pictures or upload and share with others.	YouTube, Flickr, Snapchat, Instagram, Vine

with more than 600 million users. Google+, LinkedIn, and Twitter completed the top 5 with 300 million, 259 million, and 232 million active users, respectively.

Statista also provided the growth trend for both Facebook and LinkedIn, illustrated in Figure 1.2 and Figure 1.3, respectively. Figure 1.2 shows that Facebook, by reaching 845 million users at the end of 2011, totaled 1,228 million users by the end of 2013. As depicted in Figure 1.3, LinkedIn also reached 277 million users by the end of 2013, whereas it only had 145 million users

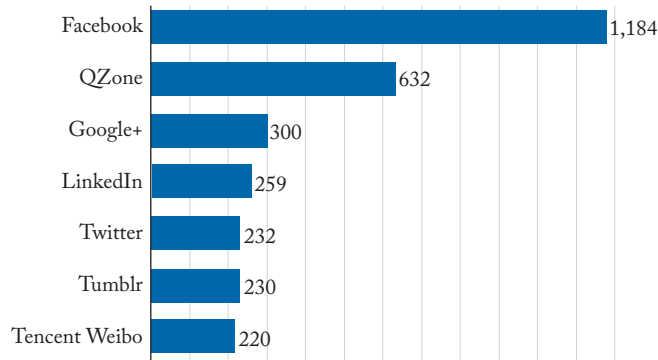


Figure 1.1: Social networks ranked by the number of active users as of January 2014 (in millions) provided by Statista.

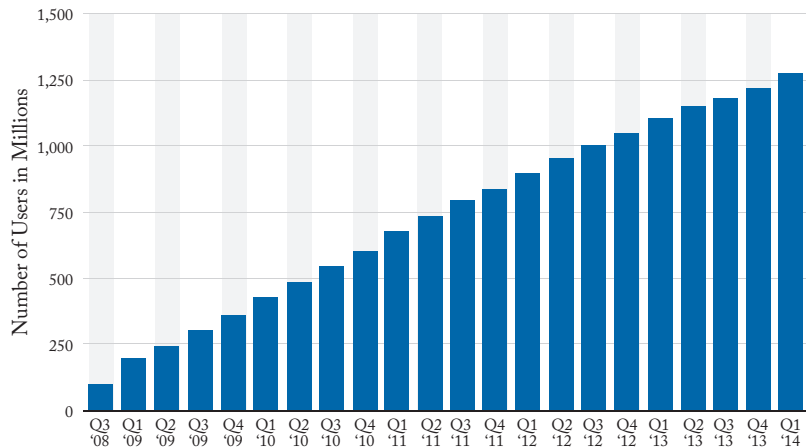
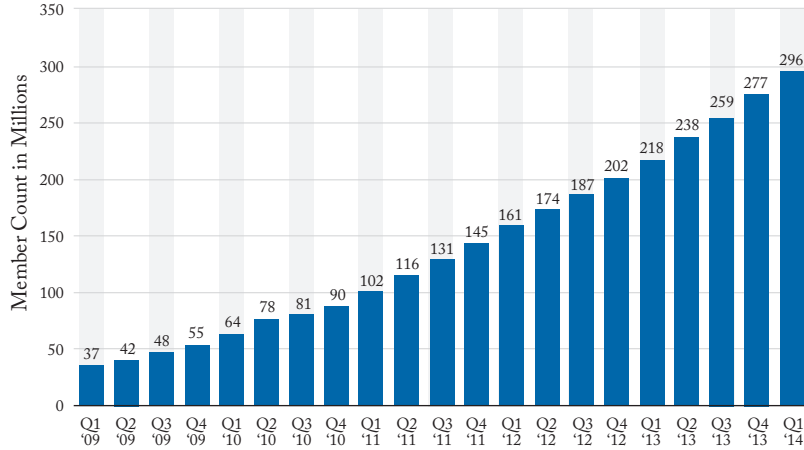


Figure 1.2: Number of monthly active Facebook users from the third quarter of 2008 to the first quarter of 2014 (in millions) provided by Statista.

at the end of 2011. Statista also calculated the annual income for both Facebook and LinkedIn, which in 2013 totalled US\$7,872 and US\$1,528 million, respectively.

Social computing is an emerging field that focuses on modeling, analysis, and monitoring of social behavior on different media and platforms to produce intelligent applications. Social media is the use of electronic and Internet tools for the purpose of sharing and discussing information and experiences with other human beings in efficient ways [Moturu, 2009]. Various social media platforms such as social networks, forums, blogs, and micro-blogs have recently evolved to ensure the connectivity, collaboration, and formation of virtual communities. While traditional media such as newspapers, television, and radio provide unidirectional communica-

## 4 1. INTRODUCTION TO SOCIAL MEDIA ANALYSIS



**Figure 1.3:** Number of LinkedIn members from the first quarter of 2009 to the first quarter of 2014 (in millions) provided by Statista.

tion from business to consumer, social media services have allowed interactions among users across various platforms. Social media have therefore become a primary source of information for business intelligence.

There are several means of interaction in social media platforms. One of the most important is via text posts. The natural language processing (NLP) of traditional media such as written news and articles has been a popular research topic over the past 25 years. NLP typically enables computers to derive meaning from natural language input using the knowledge from computer science, artificial intelligence, and linguistics.

NLP for social media text is a new research area, and it requires adapting the traditional NLP methods to these kinds of texts or developing new methods suitable for information extraction and other tasks in the context of social media.

There are many reasons why the “traditional” NLP are not good enough for social media texts, such as their informal nature, the new type of language, abbreviations, etc. Section 1.3 will discuss these aspects in more detail.

A social network is made up of a set of actors (such as individuals or organizations) and a set of binary relations between these actors (such as relationships, connections, or interactions). From a social network perspective, the goal is to model the structure of a social group to identify how this structure influences other variables and how structures change over time. Semantic analysis in social media (SASM) is the semantic processing of the text messages as well as of the meta-data, in order to build intelligent applications based on social media data.

SASM helps develop automated tools and algorithms to monitor, capture, and analyze the large amounts of data collected from social media in order to predict user behavior or ex-

tract other kinds of information. If the amount of data is very large, techniques for “big data” processing need to be used, such as online algorithms that do not need to store all the data in order to update the models based on the incoming data.

In this book, we focus on the analysis of the textual data from social media, via new NLP techniques and applications. Workshops such as the EACL 2014 Workshop on Language Analysis in Social Media [Farzindar et al., 2014], the NAACL/HLT 2013 workshop on Language Analysis in Social Media [Farzindar et al., 2013], and the EACL 2012 Workshop for Semantic Analysis in Social Media [Farzindar and Inkpen, 2012] have been increasingly focusing on NLP techniques and applications that study the effect of social media messages on our daily lives, both personally and professionally.

Social media textual data is the collection of openly available texts that can be obtained publicly via blogs and micro-blogs, Internet forums, user-generated FAQs, chat, podcasts, online games, tags, ratings, and comments. Social media texts have several properties that make them different than traditional texts, because the nature of the social conversations, posted in real time. Detecting groups of topically related conversations is important for applications, as well as detection emotions, rumors, and incentives. As an example, in order to investigate youths’ experience of grief and mourning, a study applied NLP techniques to their tweets after the death of friends or family members [Patton et al., 2018]. Determining the locations mentioned in the messages or the locations of the users can also add valuable information. The texts are unstructured and are presented in many formats and written by different people in many languages and styles. Also, the typographic errors and chat slang have become increasingly prevalent on social networking sites like Facebook and Twitter. The authors are not professional writers and their postings are spread in many places on the Web, on various social media platforms.

Monitoring and analyzing this rich and continuous flow of user-generated content can yield unprecedentedly valuable information, which would not have been available from traditional media outlets. Semantic analysis of social media has given rise to the emerging discipline of big data analytics, which draws from social network analysis, machine learning, data mining, information retrieval, and natural language processing [Melville et al., 2009].

Figure 1.4 shows a framework for semantic analysis in social media. The first step is to identify issues and opportunities for collecting data from social networks. The data can be in the form of stored textual information (the big data could be stored in large and complex databases or text files), it could be dynamic online data collection processed in real time, or it could be retrospective data collection for particular needs. The next step is the SASM pipeline, which consists of specific NLP tools for the social media analysis and data processing. Social media data is made up of large, noisy, and unstructured datasets. SASM transforms social media data to meaningful and understandable messages through social information and knowledge. Then, SASM analyzes the social media information in order to produce social media intelligence. Social media intelligence can be shared with users or presented to decision-makers to improve



## 6 1. INTRODUCTION TO SOCIAL MEDIA ANALYSIS

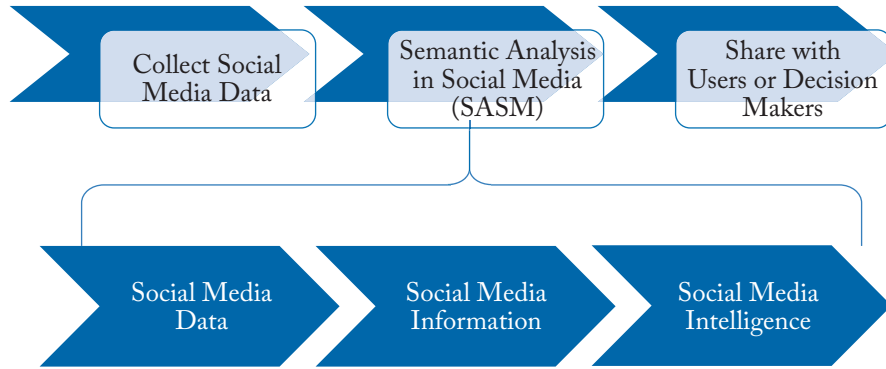


Figure 1.4: A framework for semantic analysis in social media, where NLP tools transform the data into intelligence.

awareness, communication, planning, or problem solving. The presentation of analyzed data by SASM could be completed by data visualization methods.

### 1.2 SOCIAL MEDIA APPLICATIONS

The automatic processing of social media data needs to design appropriate research methods for applications such as information extraction, automatic categorization, clustering, indexing data for information retrieval, and statistical machine translation. The sheer volume of social media data and the incredible rate at which new content is created makes monitoring, or any other meaningful manual analysis, unfeasible. In many applications, the amount of data is too large for effective real-time human evaluation and analysis of the data for a decision maker.

Social media monitoring is one of the major applications in SASM. Traditionally, **media monitoring** is defined as the activity of monitoring and tracking the output of the hard copy, online, and broadcast media which can be performed for a variety of reasons, including political, commercial, and scientific. The huge volume of information provided via social media networks is an important source for open intelligence. Social media make the direct contact with the target public possible. Unlike traditional news, the opinion and sentiment of authors provide an additional dimension for the social media data. The different sizes of source documents—such as a combination of multiple tweets and blogs—and content variability also render the task of analyzing social media documents difficult.

In social media, the real-time event search or event detection The search queries consider multiple dimensions, including spatial and temporal. In this case, some NLP methods such as information retrieval and summarization of social data in the form of various documents from multiple sources become important in order to support the event search and the detection of relevant information.

The semantic analysis of the meaning of a day's or week's worth of conversations in social networks for a group of topically related discussions or about a specific event presents the challenges of cross-language NLP tasks. Social media—related NLP methods that can extract information of interest to the analyst for preferential inclusion also lead us to domain-based applications in computational linguistics.

### 1.2.1 CROSS-LANGUAGE DOCUMENT ANALYSIS IN SOCIAL MEDIA DATA

The application of existing NLP techniques to social media from different languages and multiple resources faces several additional challenges; the tools for text analysis are typically designed for specific languages. The main research issue therefore lies in assessing whether language-independence or language-specificity is to be preferred. Users publish content not only in English, but in a multitude of languages [Blodgett et al., 2018]. This means that due to the language barrier, many users cannot access all available content. The use of machine translation technology can help bridge the language gap in such situations. The integration of machine translation and NLP tools opens opportunities for the semantic analysis of text via cross-language processing.

Natural languages constantly evolve and are adapted based on the environment of their use. Diachronic differences measure the semantic drift for these languages [Jaidka et al., 2018].

### 1.2.2 DEEP LEARNING TECHNIQUES FOR SOCIAL MEDIA DATA

Language-independent NLP tools are very important. They are not only cost- or time-efficient, but can also capture semantic aspects of each language directly. In earlier days, machine learning approaches targeting NLP tasks have mostly relied on shallow models (e.g., Support Vector Machine (SVM) and logistic regression classifiers) that are trained on high-dimensional and sparse features. In the last few years, neural networks based on dense vector representations have produced superior results on various NLP tasks. Deep neural networks [LeCun et al., 2015] enable multi-level automatic feature representation learning. Simple deep learning frameworks were shown to outperform most state-of-the-art approaches in several NLP tasks such as named-entity recognition, semantic role labeling, and part-of-speech tagging [Young et al., 2018]. Then, numerous complex deep learning-based algorithms have been proposed to solve difficult NLP tasks.

Since English is a widely used language, a majority of the research in NLP and deep learning is focused on English. But in multi-lingual countries like India, people generally use words from more than one language in their everyday speech, and on social media sites like Facebook and Twitter. This linguistic behavior is called code-mixing. Deep learning architectures can now be used on such code-mixed tweets, for example for tasks such as humor detection [Sane et al., 2019].

### 1.2.3 REAL-WORLD APPLICATIONS

The huge volume of publicly available information on social networks and on the Web can benefit different areas such as industry, media, healthcare, politics, public safety, and security. Here, we can name a few innovative integrations for social media monitoring, and some model scenarios of government-user applications in coordination and situational awareness. We will show how NLP tools can help governments interpret data in near real-time and provide enhanced command decision at the strategic and operational levels.

#### Industry

There is great interest on the part of industry in social media data monitoring. Social media data can dramatically improve business intelligence (BI). Businesses could achieve several goals by integrating social data into their corporate BI systems, such as branding and awareness, customer/prospect engagement, and improving customer service. Online marketing, stock market prediction, product recommendation, and reputation management are some examples of real-world applications for semantic analysis of social media. Recommender systems are a necessity in the modern era of technology. It is the usual tendency of people to get a review from others before going to a restaurant, watching a movie, or buying any product ranging from furniture to electronics or books. A recommender system is built on a similar approach and aims to give a relevant prediction to the target user based on the user's data, the item's data, and other users' feedback for those items. For example, [Alharthi et al. \[2018\]](#) examined the recommender systems in the field of books. These systems analyze the reading behavior of a user and the kind of books he/she likes, as well as their posting on social media, when available.

#### Media and Journalism

The relationship between journalists and the public became closer thanks to social networking platforms. Statistics published by a 2013 social journalism study show that 25% of major information sources come from social media data.<sup>3</sup> The public relations professionals and journalists use the power of social media to gather the public opinion, perform sentiment analysis, implement crisis monitoring, perform issues- or program-based media analysis, and survey social media.

#### Healthcare

Over time, social media became part of common healthcare. The healthcare industry uses social media tools for building community engagement and fostering better relationships with their clients. The use of Twitter to discuss recommendations for providers and consumers (patients, families, or caregivers), ailments, treatments, and medication is only one example of social media in healthcare. This was initially referred to as social health. Medical forums appeared due to the needs of the patients to discuss their feelings and experiences.

<sup>3</sup><http://www.cision.com/uk/files/2013/10/social-journalism-study-2013.pdf>

This book will discuss how NLP methods on social media data can help develop innovative tools and integrate appropriate linguistic information in order to allow better health monitoring (such as disease spread) or availability of information and support for patients.

### Politics

Online monitoring can help keep track of mentions made by citizens across the country and of international, national, or local opinion about political parties. For a political party, organizing an election campaign and gaining followers is crucial. Opinion mining, awareness of comments and public posts, and understanding statements made on discussion forums can give political parties a chance to get a better idea of the reality of a specific event, and to take the necessary steps to improve their positions.

### Defense and Security

Defense and security organizations are greatly interested in studying these sources of information and summaries to understand situations and perform sentiment analysis of a group of individuals with common interests, and also to be alerted against potential threats to defense and public safety. In this book, we will discuss the issue of information flow from social networks such as MySpace, Facebook, Skyblog, and Twitter. We will present methods for information extraction in Web 2.0 to find links between data entities, and to analyze the characteristics and dynamism of networks through which organizations and discussions evolve. Social data often contain significant information hidden in the texts and network structure. Aggregate social behavior can provide valuable information for the sake of national security.

## 1.3 CHALLENGES IN SOCIAL MEDIA DATA

The information presented in social media, such as online discussion forums, blogs, and Twitter posts, is highly dynamic and involves interaction among various participants. There is a huge amount of text continuously generated by users in informal environments.

Standard NLP methods applied to social media texts are therefore confronted with difficulties due to non-standard spelling, noise, and limited sets of features for automatic clustering and classification. Social media are important because the use of social networks has made everybody a potential author, so the language is now closer to the user than to any prescribed norms [Beverungen and Kalita, 2011, Zhou and Hovy, 2006]. Blogs, tweets, and status updates are written in an informal, conversational tone—often more of a “stream of consciousness” than the carefully thought out and meticulously edited work that might be expected in traditional print media. This informal nature of social media texts presents new challenges to all levels of automatic language processing.

At the surface level, several issues pose challenges to basic NLP tools developed for traditional data. Inconsistent (or absent) punctuation and capitalization can make detection of sentence boundaries quite difficult—sometimes even for human readers, as in the following tweet:

“#qcpoli enjoyed a hearty laugh today with #plq debate audience for @jflisee #notrehome tune was that the intended reaction?” Emoticons, incorrect or non-standard spelling, and rampant abbreviations complicate tokenization and part-of-speech tagging, among other tasks. Traditional tools must be adapted to consider new variations such as letter repetition (“heyyyyyy”), which are different from common spelling errors. Grammaticality, or frequent lack thereof, is another concern for any syntactic analyses of social media texts, where fragments can be as commonplace as actual full sentences, and the choice between “there,” “they are,” “they’re,” and “their” can seem to be made at random.

Social media are also much noisier than traditional print media. Like much else on the Internet, social networks are plagued with spam, ads, and all manner of other unsolicited, irrelevant, or distracting content. Even by ignoring these forms of noise, much of the genuine, legitimate content on social media can be seen as irrelevant with respect to most information needs. André et al. [2012] demonstrate this in a study that assesses user-perceived value of tweets. They collected over 40,000 ratings of tweets from followers, in which only 36% of tweets were rated as “worth reading,” while 25% were rated as “not worth reading.” The least valued tweets were so-called presence maintenance posts (e.g., “Hullo twitter!”). Pre-processing to filter out spam and other irrelevant content, or models that are better capable of coping with noise are essential in any language-processing effort targeting social media.

Several characteristics of social media text are of particular concern to NLP approaches. The particularities of a given medium and the way in which that medium is used can have a profound effect on what constitutes a successful summarization approach. For example, the 140-character limit imposed on Twitter posts makes for individual tweets that are rather contextually impoverished compared to more traditional documents. However, redundancy can become a problem over multiple tweets, due in part to the practice of retweeting posts. Sharifi et al. [2010] note the redundancy of information as a major issue with microblog summarization in their experiments with data mining techniques to automatically create summary posts of Twitter trending topics.

A major challenge facing detection of events of interest from multiple Twitter streams is therefore to separate the mundane and polluted information from interesting real-world events. In practice, highly scalable and efficient approaches are required for handling and processing the increasingly large amount of Twitter data (especially for real-time event detection). Other challenges are inherent to Twitter design and usage. These are mainly due to the shortness of the messages: the frequent use of (dynamically evolving) informal, irregular, and abbreviated words, the large number of spelling and grammatical errors, and the use of improper sentence structure and mixed languages. Such data sparseness, lack of context, and diversity of vocabulary make the traditional text analysis techniques less suitable for tweets [Metzler et al., 2007]. In addition, different events may enjoy different popularity among users, and can differ significantly in content, number of messages and participants, time periods, inherent structure, and causal relationships [Nallapati et al., 2004].

Across all forms of social media, subjectivity is an ever-present trait. While traditional news texts may strive to present an objective, neutral account of factual information, social media texts are much more subjective and opinion-laden. Whether or not the ultimate information need lies directly in opinion mining and sentiment analysis, subjective information plays a much greater role in semantic analysis of social texts.

Topic drift is much more prominent in social media than in other texts, both because of the conversational tone of social texts and the continuously streaming nature of social media. There are also entirely new dimensions to be explored, where new sources of information and types of features need to be assessed and exploited. While traditional texts can be seen as largely static and self-contained, the information presented in social media, such as online discussion forums, blogs, and Twitter posts, is highly dynamic and involves interaction among various participants. This can be seen as an additional source of complexity that may hamper traditional summarization approaches, but it is also an opportunity, making available additional context that can aid in summarization or making possible entirely new forms of summarization. For instance, [Hu et al. \[2007a\]](#) suggest summarizing a blog post by extracting representative sentences using information from user comments. [Chua and Asur \[2012\]](#) exploit temporal correlation in a stream of tweets to extract relevant tweets for event summarization. [Lin et al. \[2009\]](#) address summarization not of the content of posts or messages, but of the social network itself by extracting temporally representative users, actions, and concepts in Flickr data.

As we mentioned, standard NLP approaches applied to social media data are therefore confronted with difficulties due to non-standard spelling, noise, limited sets of features, and errors. Therefore some NLP techniques, including normalization, term expansion, improved feature selection, and noise reduction, have been proposed to improve clustering performance in Twitter news [[Beverungen and Kalita, 2011](#)]. Identifying proper names and language switch in a sentence would require rapid and accurate name entity recognition and language detection techniques. Recent research efforts focus on the analysis of language in social media for understanding social behavior and building socially aware systems. The goal is the analysis of language with implications for fields such as computational linguistics, sociolinguistics, and psycholinguistics. For example, [Eisenstein \[2013a\]](#) studied the phonological variation and factors when transcribed into social media text.

Several workshops organized by the Association for Computational Linguistics (ACL)<sup>4</sup> and special issues in scientific journals dedicated to semantic analysis in social media show how active this research field is.

In this book, we will cite many papers from conferences such as ACL, AAAI, WWW, etc.; many workshop papers; several books; and many journal papers from various relevant journals.

<sup>4</sup>All publications could be found at ACL Anthology <https://www.aclweb.org/anthology/>

## 1.4 SEMANTIC ANALYSIS OF SOCIAL MEDIA

Our goal is to focus on innovative NLP applications (such as opinion mining, information extraction, summarization, and machine translation), tools, and methods that integrate appropriate linguistic information in various fields such as social media monitoring for healthcare, security and defense, business intelligence, and politics. The book contains four major chapters.

- **Chapter 1:** This chapter highlights the need for applications that use social media messages and meta-data. We also discuss the difficulty of processing social media data vs. traditional texts such as news articles and scientific papers.
- **Chapter 2:** This chapter discusses existing linguistic pre-processing tools such as tokenizers, part-of-speech taggers, parsers, and named entity recognizers, with a focus on their adaptation to social media data. We briefly discuss evaluation measures for these tools.
- **Chapter 3:** This chapter is the heart of the book. It presents the methods used in applications for semantic analysis of social network texts, in conjunction with social media analytics as well as methods for information extraction and text classification. We focus on tasks such as: geo-location detection, entity linking, opinion mining and sentiment analysis, emotion and mood analysis, event and topic detection, summarization, machine translation, and other tasks. They tend to pre-process the messages with some of the tools mentioned in Chapter 2 in order to extract the knowledge needed in the next processing levels. For each task, we discuss the evaluation metrics and any existing test datasets.
- **Chapter 4:** This chapter presents higher-level applications that use some of the methods from Chapter 3. We look at: healthcare applications, financial applications, predicting voting intentions, media monitoring, security and defense applications, NLP-based information visualization for social media, disaster response applications, NLP-based user modeling, applications for entertainment, rumor detection, and recommender systems.
- **Chapter 5:** This chapter discusses chapter complementary aspects such as data collection and annotation in social media, privacy issues in social media, spam detection in order to avoid spam in the collected datasets, and we describe some of the existing evaluation benchmarks that make available data collected and annotated for various tasks.
- **Chapter 6:** The last chapter summarizes the methods and applications described in the preceding chapters. We conclude with a discussion of the high potential for research, given the social media analysis needs of end-users.

As mentioned in the Preface, the **intended audience** of this book is researchers that 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. Nonetheless, we will try to define as many notions as we can, in order to



facilitate the understanding for beginners in these two areas. We also assume basic knowledge of computer science in general.

## 1.5 SUMMARY

In this chapter, we reviewed the structure of social network and social media data as the collection of textual information on the Web. We presented semantic analysis in social media as a new opportunity for big data analytics and for intelligent applications. Social media monitoring and analyzing of the continuous flow of user-generated content can be used as an additional dimension which contains valuable information that would not have been available from traditional media and newspapers. In addition, we mentioned the challenges with social media data, which are due to their large size, and to their noisy, dynamic, and unstructured nature.





# Linguistic Pre-processing of Social Media Texts

## 2.1 INTRODUCTION

In this chapter, we discuss current Natural Language Processing (NLP) linguistic pre-processing methods and tools that were adapted for social media texts. We survey the methods used for adaptation to this kind of texts. We briefly define the evaluation measures used for each type of tool in order to be able to mention the state-of-the-art results.

In general, evaluation in NLP can be done in several ways:

- manually, by having humans judge the output of each tool;
- automatically, on test data that humans have annotated with the expected solution ahead of time; and
- task-based, by using the tools in a task and evaluating how much they contribute to the success in the task.

We primarily focus on the second approach here. It is the most convenient since it allows the automatic evaluation of the tools repeatedly after changing/improving their methods, and it allows comparing different tools on the same test data. Care should be taken when human judges annotate data. There should be at least two annotators that are given proper instructions on what and how to annotate (in an annotation manual). There needs to be a reasonable agreement rate between the two or more annotators, to ensure the quality of the obtained data. When there are disagreements, the expected solution will be obtained by resolving the disagreements by taking a vote (if there are three annotators or more, an odd number), or by having the annotators discuss until they reach an agreement (if there are only two annotators, or an even number). When reporting the inter-annotator agreement for a dataset, the kappa statistic also needs to be reported, in order to compensate the obtained agreement for possible agreements due to chance [Artstein and Poesio, 2008, Carletta, 1996].

NLP tools often use supervised machine learning, and the training data are usually annotated by human judges. In such cases, it is convenient to keep aside some of the annotated data for testing and to use the remaining data to train the models. Many of the methods discussed in this book use machine learning algorithms for automatic text classification. That is why we give a very brief introduction here. See, e.g., [Witten and Frank, 2005] for details of the classical

algorithms and [Sebastiani, 2002] for how they can be applied to text data. Also see [Eisenstein, 2019] for more details about deep learning classification techniques for text data.

A supervised text classification model predicts the label  $c$  of an input  $x$ , where  $x$  is a vector of feature values extracted from document  $d$ . The class  $c$  can take two or more possible values from a specified set (or even continuous numeric values, in which case the classifier is called a regression model). The training data contain document vectors for which the classes are provided. The classifier uses the training data to learn associations between features or combinations of features that are strongly associated with one of the classes but not with the other classes. In this way, the trained model can make predictions for unseen test data in the future. There are many classification algorithms. We name here only a few of the classifiers popular in NLP tasks.

Decision trees take one feature at a time, compute its power of discriminating between the classes and build a tree with the most discriminative features in the upper part of the tree; decision trees are useful because the models can be easily understood by humans. Naïve Bayes is a classifier that learns the probabilities of association between features and classes; these models are used because they are known to work well with text data (see a more detailed description in Section 2.8.1). SVMs compute a hyper plane that separates two classes and they can efficiently perform nonlinear classification using what is called a kernel to map the data into a high-dimensional feature space where it become linearly separable [Cortes and Vapnik, 1995]; SVMs are probably the most often used classifiers due to their high performance on many tasks that have small amounts of training data.

Lately, most linguistic tools and applications employ deep neural network classifiers, which were shown to lead to better performance when large amounts of training data are available. These classifiers include Recurrent Neural Networks (RNN) or Convolutional Neural Networks (CNN) [LeCun et al., 2015]. A special type of RNN is Long Short-Term Memory (LSTM) networks, which add some forget gates to allow modeling long-distance context [Hochreiter and Schmidhuber, 1997].

There are many machine learning libraries that can be used. We mention here only a few: Weka<sup>1</sup> and scikit-learn<sup>2</sup> for classical algorithms and PyTorch<sup>3</sup> and TensorFlow<sup>4</sup> for deep learning algorithms. The first one is in Java, while the last three are in Python.

A sequence-tagging model can be seen as a classification model, but fundamentally differs from a conventional one, in the sense that instead of dealing with a single input  $x$  and a single label  $c$  each time, it predicts a sequence of labels  $\mathbf{c} = (c_1, c_2, \dots, c_n)$  based on a sequence of inputs  $\mathbf{x} = (x_1, x_2, \dots, x_n)$  and the predictions from the previous steps. It was applied with success in natural language processing (for sequential data such as sequences of part-of-speech tags, discussed in the previous chapter) and in bioinformatics (for DNA sequences). There exist a number of sequence-tagging models, including Hidden Markov Model (HMM) [Baum and

<sup>1</sup>[urlhttps://www.cs.waikato.ac.nz/ml/weka/](https://www.cs.waikato.ac.nz/ml/weka/)

<sup>2</sup>[urlhttps://scikit-learn.org/stable/](https://scikit-learn.org/stable/)

<sup>3</sup>[urlhttps://pytorch.org/](https://pytorch.org/)

<sup>4</sup>[urlhttps://www.tensorflow.org/](https://www.tensorflow.org/)

Petrie, 1966], Conditional Random Field (CRF) [Lafferty et al., 2001], and Maximum Entropy Markov Model (MEMM) [Berger et al., 1996].

Sequence-to-sequence models based on deep learning can also be used to transform sequences into other sequences, for example a sequence of words into a sequence of part-of-speech tags. They are also being used for the latest machine translation system, to transform a sequence of words in one language into a sequence of words in another language [Gehring et al., 2017].

Before the input texts can be fed into classifiers, each text needs to be transformed into a set of features. For some linguistic tools, extracting the words is sufficient while for semantic tasks and applications the texts need to be transformed into vectors of numeric or discrete values. The simplest way to represent texts is using the Bag-of-Words model (BOW) (a word is present or not, possibility with with frequency information), or more advanced and less sparse vectors called word embeddings [Mikolov et al., 2013]. Linguistic features can be used in addition or instead of the word-based features. These representations are important especially for the semantic tasks and applications discussed in the next two chapters.

The remainder of this chapter is structured as follows. Section 2.2 discusses generic methods of adapting NLP tools to social media texts. The next five sections discuss NLP tools of interest: tokenizers, part-of-speech taggers, chunkers, parsers, and named entity recognizers, as well as adaptation techniques for each. Section 2.7 enumerates some of the existing toolkits that were adapted to social media texts in English. Section 2.8 discusses multi-lingual aspects and language identification issues in social media. Section 2.9 summarizes this chapter.

## 2.2 GENERIC ADAPTATION TECHNIQUES FOR NLP TOOLS

NLP tools are important because they need to be used before we can build any applications that aim to understand texts or extract useful information from texts. Many NLP tools are now available, with acceptable levels of accuracy on texts that are similar to the types of texts used for training the models embedded in these tools. Most of the tools are trained on carefully edited texts, usually newspaper texts, due to the wide availability of these kinds of texts. For example, the Penn TreeBank corpus, consisting of 4.5 million words of American English [Marcus et al., 1993], was manually annotated with part-of-speech tags and parse trees, and it is often the main resource used to train part-of-speech taggers and parsers.

Current NLP tools tend to work poorly on social media texts, because these texts are informal, not carefully edited, and they contain grammatical errors, misspellings, new types of abbreviations, emoticons, etc. They are very different than the types of texts used for training the NLP tools. Therefore, the tools need to be adapted in order to achieve reasonable levels of performance on social media texts.

Table 2.1 shows three examples of Twitter messages, taken from Ritter et al. [2011], just to illustrate how noisy the texts can be.

Table 2.1: Three examples of Twitter texts

No.	Example
1	The Hobbit has FINALLY started filming! I cannot wait!
2	@c@ Yess! Yess! It's official Nintendo announced today that theyWill release the Nintendo 3DS in north America march 27 for \$250
3	Government confirms blast n #nuclear plants n #japan...don't knw wht s gona happen nw...

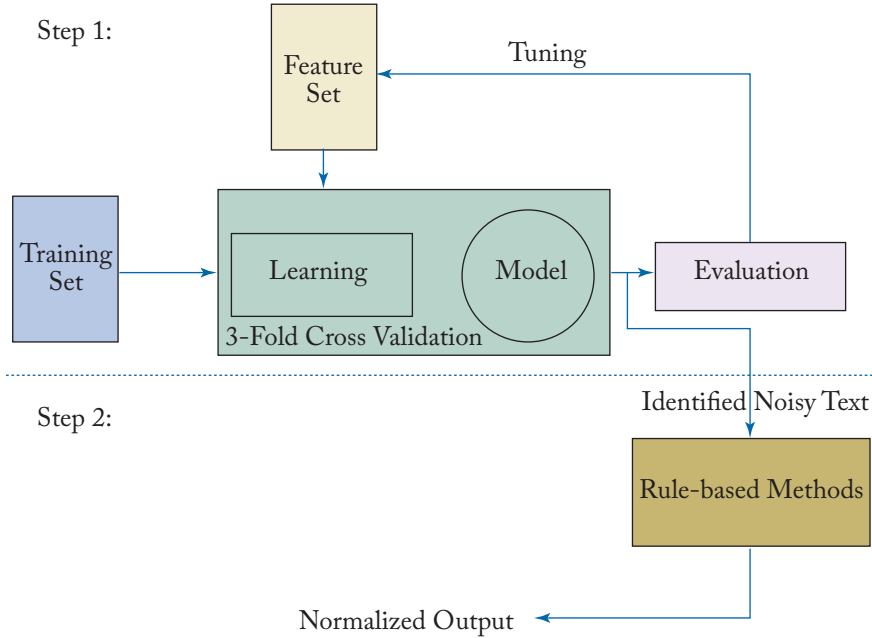
There are two ways to adapt NLP tools to social media texts. The first one is to perform text normalization so that the informal language becomes closer to the type of texts on which the tools were trained. The second one is to re-train the models inside the tool on annotated social media texts. Depending on the goal of the NLP application, a combination of the two techniques could be used, since both have their own limitations, as discussed below (see [Eisenstein \[2013b\]](#) for a more detailed discussion).

### 2.2.1 TEXT NORMALIZATION

Text normalization is a possible solution for overcoming or reducing linguistic noise. The task can be approached in two stages: first, the identification of orthographic errors in an input text, and second, the correction of these errors. Normalization approaches typically include a dictionary of known correctly spelled terms, and detects in-vocabulary and out-of-vocabulary (OOV) terms with respect to this dictionary. The normalization can be basic or more advanced. Basic normalization deals with the errors detected at the POS tagging stage, such as unknown words, misspelled words, etc. Advanced normalization is more flexible, taking a lightly supervised automatic approach trained on an external dataset (annotated with short forms vs. their equivalent long or corrected forms).

For social media texts, the normalization that can be done is rather shallow. Because of its informal and conversational nature, social media text cannot become carefully edited English. Similar issues appear in SMS text messages on phones, where short forms and phonetic abbreviations are often used to save the typing time. According to [Derczynski et al. \[2013b\]](#), text normalization in Twitter messages did not help too much in the named entity recognition task.

Twitter text normalization into traditional written English [[Han and Baldwin, 2011](#)] is not only difficult, but it can be viewed as a “lossy” translation task. For example, many of Twitter’s unique linguistic phenomena are due not only to its informal nature, but also to a set of authors that is heavily skewed toward younger ages and minorities, with heavy usage of dialects that are different than standard English [[Eisenstein, 2013a](#), [Eisenstein et al., 2011](#)].



**Figure 2.1:** Methodology for tweet normalization. The dotted horizontal line separates the two steps (detecting the text to be normalized and applying normalization rules) [Akhtar et al., 2015].

Demir [2016] describes a method of context-tailored text normalization. The method considers contextual and lexical similarities between standard and non-standard words, in order to reduce noise. The non-standard words in the input context in a given sentence are tailored into a direct match, if there are possible shared contexts. A morphological parser is used to analyze all the words in each sentence. Turkish social media texts were used to evaluate the performance of the system. The dataset contains tweets (~11 GB) and clean Turkish texts (~ 6 GB). The system achieved state-of-the-art results on the 715 Turkish tweets.

Akhtar et al. [2015] proposed a hybrid approach for text normalization for tweets. Their methodology proceeds in two phases: the first one detects noisy text, and the second one uses various heuristic-based rules for normalization. The researchers trained a supervised learning model, using 3-fold cross validation to determine the best feature set. Figure 2.1 depicts a schematic diagram of the proposed approach. Their system yielded precision, recall, and F-measure values of 0.90, 0.72, and 0.80, respectively, for their test dataset.

Most practical applications leverage the simpler approach of replacing non-standard words with their standard counterparts as a “one size fits all” task. Baldwin and Li [2015] devised a method that uses a taxonomy of normalization edits. The researchers evaluated this method on

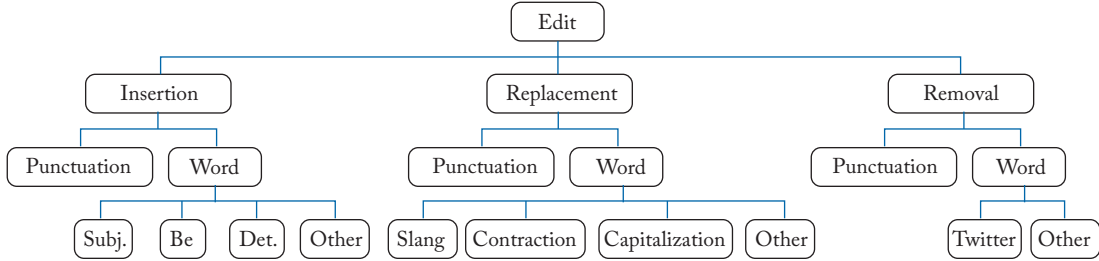


Figure 2.2: Taxonomy of normalization edits [Baldwin and Li, 2015].

three different downstream applications: dependency parsing, named entity recognition, and text-to-speech synthesis. The taxonomy of normalization edits is shown in Figure 2.2. The method categorizes edits at three levels of granularity and its results demonstrate that the targeted application of the taxonomy is an efficient approach to normalization.

The effect of manual vs. automatic lexical normalization for dependency parsing was analyzed by van der Goot [2019]. They showed that for most categories, automatic normalization scores are close to manual normalization but the small differences are important to take into consideration when exploiting normalization in a pipeline setup.

### 2.2.2 RE-TRAINING NLP TOOLS FOR SOCIAL MEDIA TEXTS

Re-training NLP tools for social media texts is relatively easy if annotated training data are available. In general, adapting a tool to a specific domain or a specific type of text requires producing annotated training data for that kind of text. It is easy to collect text of the required kind, but to annotate it can be a difficult and time-consuming process.

Currently, some annotated social media data have become available, but the volume is not high enough. Several NLP tools have been re-trained on newly annotated data, sometimes by also keeping the original annotated training data for newspaper texts, in order to have a large enough training set. Another approach is to use some unlabeled social media text in an unsupervised manner in addition to the small amounts of annotated social media text.

Another question is what kinds of social media texts to use for training. It seems that Twitter messages are more difficult to process than blog posts or messages from forums. Because of the limitation of Twitter messages to 140 characters, more abbreviations and shortened forms of words are used, and more simplified syntax. Therefore, training data should include several kinds of social media texts (unless somebody is building a tool designed for a particular kind of social media text).

We define the tasks accomplished by each kind of tool and we discuss techniques for adapting them to social media texts.

Table 2.2: Examples of tokenization

	(1)	(2)	(3)
Raw	pdf?“<-Wenn	schriftq.Äquivalent	v.14.4
Tokenized	pdf_?“<-_Wenn	schriftq._Äquivalent	v._14._4

## 2.3 TOKENIZERS

The first step in processing a text is to separate the words from punctuation and other symbols. A tool that does this is called a tokenizer. White space is a good indicator of words separation (except in some languages, e.g., Chinese), but even white space is not sufficient. The question of what is a word is not trivial. When doing corpus analysis, there are strings of characters that are clearly words, but there are strings for which this is not clear. Most of the time, punctuation needs to be separated from words, but some abbreviations might contain punctuation characters as part of the word. Take, for example, the sentence: “We bought apples, oranges, etc.” The commas clearly need to be separated from the word “apples” and from the word “oranges,” but the dot is part of the abbreviation “etc.” In this case, the dot also indicates the end of the sentence (two dots were reduced to one). Other examples among the many issues that appear are: how to treat numbers (if they contain commas or dots, these characters should not be separated), or what to do with contractions such as “don’t” (perhaps to expand them into two words “do” and “not”).

While tokenization usually consists of two subtasks (sentence boundary detection and token boundary detection), the Empirist shared task<sup>5</sup> provided sentence boundaries and the participating teams only had to detect token boundaries. Missing whitespace characters presents a major challenge to the task of tokenization. Table 2.2 shows a few examples with their correct tokenization.

### Methods for Tokenizers

Horsmann and Zesch [2016] evaluated a method for dealing with token boundaries consisting of three steps. First, the researchers split the text according to the white space characters. Then they employed regular expressions to refine the splitting of alpha-numerical text segments from punctuation characters in special character sequences such as similes. Finally, these sequences of punctuation are reassembled. They merge the most common combinations of characters into a single token using the training data, and use word lists to merge abbreviations with their following dot character. They increase accuracy in the experiment using more in-domain training data.

<sup>5</sup><https://sites.google.com/site/empirist2015/>



### Evaluation Measures for Tokenizers

Accuracy is a simple measure that calculates how many correct decisions a tool makes. When not all the expected tokens are retrieved, precision and recall are the measure to report. The precision of the tokens recognition measures how many tokens are correct out of how many were found. Recall measures the coverage (from the tokens that should have been retrieved, how many were found). F-measure (or F-score) is often reported when one single number is needed, because F-measure is the harmonic mean of the precision and recall, and it is high only when both the precision and the recall are high.<sup>6</sup> Evaluation measures are rarely reported for tokenizers, one exception being the CleanEval shared task which focused on tokenizing text from web pages [Baroni et al., 2008].

Many NLP projects tend to not mention what kind of tokenization they used, and focus more on higher-level processing. Tokenization, however, can have a large effect on the results obtained at the next levels. For example, Fokkens et al. [2013] replicated two high-level tasks from previous work and obtained very different results, when using the same settings but different tokenization.

### Adapting Tokenizers to Social Media Texts

Tokenizers need to deal with the specifics of social media texts. Emoticons need to be detected as tokens. For Twitter messages, user names (starting with @), hashtags (starting with #), and URLs (links to web pages) should be treated as tokens, without separating punctuation or other symbols that are part of the token. Some shallow normalization can be useful at this stage. Derczynski et al. [2013b] tested a tokenizer on Twitter data, and its F-measure was around 80%. By using regular expressions designed specifically for Twitter messages, they were able to increase the F-measure to 96%. More about such regular expressions can be found in [O'Connor et al., 2010].

## 2.4 PART-OF-SPEECH TAGGERS

Part-of-speech (POS) taggers determine the part of speech of each word in a sentence. They label nouns, verbs, adjectives, adverbs, interjections, conjunctions, etc. Often they use finer-grained tagsets, such as singular nouns, plural nouns, proper nouns, etc. Different tagsets exist, one of the most popular being the Penn TreeBank tagset<sup>7</sup> [Marcus et al., 1993]. See Table 2.3 for one of its more popular lists of the tags. The models embedded in the POS taggers are often complex, based on Hidden Markov Models [Baum and Petrie, 1966], Conditional Random Fields [Lafferty et al., 2001], etc. They need annotated training data in order to learn probabilities and other parameters of the models.

<sup>6</sup>The F-score usually gives the same weight to precision and to recall, but it can weight one of them more when needed for an application.

<sup>7</sup><http://www.comp.leeds.ac.uk/ccalas/tagsets/upenn.html>

Table 2.3: Penn TreeBank tagset

Number	Tag	Description
1	CC	Coordinating conjunction
2	CD	Cardinal number
3	DT	Determiner
4	EX	Existential there
5	FW	Foreign word
6	IN	Preposition or subordinating conjunction
7	JJ	Adjective
8	JJR	Adjective, comparative
9	JJS	Adjective, superlative
10	LS	List item marker
11	MD	Modal
12	NN	Noun, singular or mass
13	NNS	Noun, plural
14	NNP	Proper noun, singular
15	NNPS	Proper noun, plural
16	PDT	Predeterminer
17	POS	Possessive ending
18	PRP	Personal pronoun
19	PRP\$	Possessive pronoun
20	RB	Adverb
21	RBR	Adverb, comparative
22	RBS	Adverb, superlative
23	RP	Particle
24	SYM	Symbol
25	TO	To
26	UH	Interjection
27	VB	Verb, base form
28	VBD	Verb, past tense
29	VBG	Verb, gerund or present participle
30	VCN	Verb, past participle
31	VBP	Verb, non-3rd person singular present
32	VBZ	Verb, 3rd person singular present
33	WDT	Wh-determiner
34	WP	Wh-pronoun
35	WP\$	Possessive wh-pronoun
36	WRB	Wh-adverb

### Methods for Part-of-speech Taggers

Horsmann and Zesch [2016] trained a CRF classifier [Lafferty et al., 2001] using the FlexTag tagger [Zesch and Horsmann, 2016]. There are two adaptations involved in this method. The first is a general domain adaptation. The researchers applied a domain adaptation strategy, which they proposed as a competitive model to improve the accuracy for tagging social media texts. To train their model, they used the CMC and Web corpora subsets from the EmpiriST shared task and some additional 100,000 tokens of newswire text from the Tiger corpus. The second adaptation is specific to the EmpiriST shared task. Because some POS tags are too rare to be learned from training data, the researchers utilized a post-processing step that leveraged heuristics. This step involved the use of regular expressions and word lists from Wikipedia and Wiktionary to improve named entity recognition and case-insensitive matching. Selecting tags from the larger Tiger corpus introduced bias, so the researchers added extra Boolean features to their model.

Deep learning-based POS taggers became easy to build. They directly transform sequences of words into sequences of POS tags. For example, Popov [2016] surveys the techniques that can be applied, starting with word embeddings and enhanced with suffix embeddings.

### Evaluation Measures for Part-of-speech Taggers

The accuracy of the tagging is usually measured as the number of tags correctly assigned out of the total number of words/tokens being tagged.

### Adapting Part-of-speech Taggers

POS taggers clearly need re-training in order to be usable on social media data. Even the set of POS tags used must be extended in order to adapt to the needs of this kind of text. Ritter et al. [2011] used the Penn TreeBank tagset (Table 2.3) to annotate 800 Twitter messages. They added a few new tags for the Twitter-specific phenomena: retweets, @usernames, #hashtags, and URLs. Words in these categories can be tagged with very high accuracy using simple regular expressions, but they still need to be taken into consideration as features in the re-training of the taggers (for example as tags of the previous word to be tagged). In Ritter et al. [2011], the POS tagging accuracy drops from about 97% on newspaper text to 80% on the 800 tweets. These numbers are reported for the Stanford POS tagger [Toutanova et al., 2003]. Their POS tagger T-POS—based on a Conditional Random Field classifier and on the clustering of out-of-vocabulary (OOV) words—also obtained low performance on Twitter data (81%). By re-training the T-POS tagger on the annotated Twitter data (which is rather small), the accuracy increases to 85%. The best accuracy raises to 88% when the size of the training data is increased by adding to the Twitter data the initial Penn TreeBank training data, plus 40,000 tokens of annotated Internet Relay Chat (IRC) data [Forsyth and Martell, 2007], which is similar in style to Twitter data. Similar numbers are reported by Derczynski et al. [2013b] on a part of the same Twitter dataset.

A key reason for the drop in accuracy on Twitter data is that the data contains far more OOV words than grammatical text. Many of these OOV words come from spelling variation, e.g., the use of the word *n* for *in* in Example 3 from Table 2.1. The tag for proper nouns (NNP) is the most frequent tag for OOV words, while in fact only about one third are proper nouns.

Gimpel et al. [2011] developed a new POS tagset for Twitter (see Table 2.4), that is more coarse-grained, and it pays particular attention to punctuation, emoticons, and Twitter-specific tags (@usernames, #hashtags, URLs). They manually tagged 1,827 tweets with the new tagset; then, they trained a POS tagging model that uses features geared toward Twitter text. The experiments conducted to evaluate the model showed 90% accuracy for the POS tagging task. Owoputi et al. [2013] improved on the model by using word clustering techniques and trained the POS tagger on a better dataset of tweets and chat messages.<sup>8</sup> Some of the expressions used in Twitter messages are formal, and some are informal. Therefore, POS tagging for the formal Twitter contexts can be learned together with the exiting news datasets, while POS tagging for the informal Twitter context should be learned separately. Gui et al. [2018] proposed a hypernetwork-based method to generate different parameters to separately model contexts with different expression styles. Experimental results on three test datasets showed that their approach achieves better performance than state-of-the-art methods in most cases.

## 2.5 CHUNKERS AND PARSERS

A **chunker** detects noun phrases, verb phrases, adjectival phrases, and adverbial phrases, by determining the start point and the end point of every such phrase. Chunkers are often referred to as shallow parsers because they do not attempt to connect the phrases in order to detect the syntactic structure of the whole sentence.

A **parser** performs the syntactic analysis of a sentence, and usually produces a parse tree. The trees are often used in future processing stages, toward semantic analysis or information extraction.

A **dependency parser** extracts pairs of words that are in a syntactic dependency relation, rather than a parse tree. Relations can be verb-subject, verb-object, noun-modifier, etc.

### Methods for Parsers

The methods used to build parsers range from early rule-based approaches, to robust probabilistic models and to new types of deep learning-based parsers. For example, Chen and Manning [2014] present a fast and accurate dependency parser using neural networks, trained on newspaper text. Another example is Parsey McParseface<sup>9</sup>, an open-sourced machine learning model-based authored by Google and based on the Tensorflow framework. It contains a globally normalized transition-based neural network model that achieves state-of-the-art part-of-speech tagging, dependency parsing, and sentence compression results.

<sup>8</sup>This data set is available at <http://code.google.com/p/ark-tweet-nlp/downloads/list>.

<sup>9</sup><https://deepai.org/machine-learning-model/parseymcparseface>

Table 2.4: POS tagset from Gimpel et al. [2011]

Tag	Description
N	Common noun
O	Pronoun (personal/WH, not possessive)
^	Proper noun
S	Nominal + possessive
Z	Proper noun + possessive
V	Verb including copula, auxiliaries
L	Nominal + verbal (e.g., i'm), verbal _ nominal (let's)
M	Proper noun + verbal
A	Adjective
R	Adverb
!	Interjection
D	Determiner
P	Pre- or postposition, or subordinating conjunction
&	Coordinating conjunction
T	Verb particle
X	Existential there, predeterminers
Y	X + verbal
#	Hashtag (indicates topic/category for tweet)
@	At-mention (indicates a user as a recipient of a tweet)
~	Discourse marker, indications of continuation across multiple tweets
U	URL or email address
E	Emoticon
\$	Numeral
,	Punctuation
G	Other abbreviations, foreign words, possessive endings, symbols, garbage

### Evaluation Measures for Chunking and Parsing

The Parseval evaluation campaign [Harrison et al., 1991] proposed measures that compare the phrase-structure bracketings<sup>10</sup> produced by the parser with bracketings in the annotated corpus (treebank). One computes the number of bracketing matches  $M$  with respect to the number

<sup>10</sup>A bracketing is a pair of matching opening and closing brackets in a linearized tree structure.

of bracketings  $P$  returned by the parser (expressed as precision  $M/P$ ) and with respect to the number  $C$  of bracketings in the corpus (expressed as recall  $M/C$ ). Their harmonic mean, the  $F$ -measure, is most often reported for parsers. In addition, the mean number of crossing brackets per sentence could be reported, to count the number of cases when a bracketed sequence from the parser overlaps with one from the treebank (i.e., neither is properly contained in the other). For chunking, the accuracy can be reported as the tag correctness for each chunk (labeled accuracy), or separately for each token in each chunk (token-level accuracy). The former is stricter because it does not give credit to a chunk that is partially correct but incomplete, for example one or more words too short or too long.

### Adapting Parsers

Parsing performance also decreases on social media text. Foster et al. [2011] tested four dependency parsers and showed that their performance decreases from 90%  $F$ -score on newspaper text to 70–80% on social media text (70% on Twitter data and 80% on discussion forum texts). After retraining on a small amount of social media training data (1,000 manually corrected parses) plus a large amount of unannotated social media text, the performance increased to 80–83%. Ovreliid and Skjærholt [2012] also show the labeled accuracy of dependency parsers decreasing from newspaper data to Twitter data.

Ritter et al. [2011] also explored shallow parsing and noun phrase chunking for Twitter data. The token-level accuracy for the shallow parsing of tweets was 83% with the OpenNLP chunker and 87% with their shallow parser T-chunk. Both were re-trained on a small amount of annotated Twitter data plus the Conference on Natural Language Learning (CoNLL) 2000 shared task data [Tjong Kim Sang and Buchholz, 2000].

Khan et al. [2013] reported experiments on parser adaptation to social media texts and other kinds of Web texts. They found that text normalization helps increase performance by a few percentage points, and that a tree reviser based on grammar comparison helps to a small degree. A dependency parser named TweepoParser<sup>11</sup> was developed specifically on a recently annotated Twitter treebank for 929 tweets [Kong et al., 2014]. It uses the POS tagset from Gimpel et al. [2011] presented in Table 2.4. Table 2.5 shows an example of output of the parser for the tweet: “They say you are what you eat, but it’s Friday and I don’t care! #TGIF (@ Ogallo Crows Nest) <http://t.co/l3uLuKGk>.”

The columns represent, in order: ID is the token counter, starting at 1 for each new sentence; FORM is the word form or punctuation symbol; CPOSTAG is the coarse-grained part-of-speech tag, where the tagset depends on the language; POSTAG is the fine-grained part-of-speech tag, where the tagset depends on the language, or it is identical to the coarse-grained part-of-speech tag, if not available; HEAD is the head of the current token, which is either an ID (−1 indicates that the word is not included in the parse tree; some treebanks also used zero as ID); and finally, DEPREL is the dependency relation to the HEAD. The set of dependency re-

<sup>11</sup>[http://www.ark.cs.cmu.edu/TweetNLP/#tweepoparser\\_treebank](http://www.ark.cs.cmu.edu/TweetNLP/#tweepoparser_treebank)

Table 2.5: Example of tweet parsed with the TweepoParser

ID	FORM	CPOSTAG	POSTAG	HEAD	DEPREL
1	They	O	O	2	—
2	say	V	V	9	CONJ
3	you	O	O	4	—
4	are	V	V	2	—
5	what	O	O	7	—
6	you	O	O	7	—
7	eat	V	V	4	—
8	,	,	,	—1	—
9	but	&	&	0	—
10	it's	L	L	9	CONJ
11	Friday	^	^	10	—
12	and	&	&	0	—
13	O	O	O	14	—
14	don't	V	V	12	CONJ
15	care	V	V	14	—
16	!	,	,	—1	—
17	#TGIF	#	#	—1	—
18	{@	P	P	0	—
19	Ogalo	^	^	21	MWE
20	Crows	^	^	21	MWE
21	Nest	^	^	18	—
22	)	,	,	—1	—
23	<a href="http://t.co/13uLuKGk">http://t.co/13uLuKGk</a>	U	U	—1	—

lations depends on the particular language. Depending on the original treebank annotation, the dependency relation may be meaningful or simply “ROOT.” So, for this tweet, the dependency relations are MWE (multi-word expression), CONJ (Conjunct), and many other relations between the word IDs, but they are not named (probably due to the limited training data used when the parser was trained). The dependency relations from the Stanford dependency parser are included, if they can be detected in a tweet. If they cannot be named, they are still in the table, but without a label.

## 2.6 NAMED ENTITY RECOGNIZERS

A named entity recognizer (NER) detects names in the texts, as well as dates, currency amounts, and other kinds of entities. NER tools often focus on three types of names: Person, Organization, and Location, by detecting the boundaries of these phrases. There are a few other types of tools that can be useful in the early stages of NLP applications. One example is a **co-reference resolution** tool that can be used to detect the noun that a pronoun refers to or to detect different noun phrases that refer to the same entity. In fact, NER is a semantic task, not a linguistic pre-processing task, but we introduce it this chapter because it became part of many of the recent NLP tools discussed in this chapter. We will talk more about specific kind of entities in Sections 3.2 and 3.3, in the context of integrating more and more semantic knowledge when solving the respective tasks.

### Methods for NER

NER is composed of two sub-tasks: detecting entities (the span of text where a name starts and where it ends) and determining/classifying the type of entity. The methods used in NER are either based on linguistic grammars for each type of entity, either based on statistical methods. Semi-supervised learning techniques were proposed, but supervised learning, especially based on CRFs for sequence learning, are the most prevalent. Hand-crafted grammar-based systems typically obtain good precision, but at the cost of lower recall and months of work by experienced computational linguists. Supervised learning techniques were used more recently due the availability of annotated training datasets, mostly for newspaper texts, such as data from MUC 6, MUC 7, and ACE,<sup>12</sup> and also the CoNLL 2003 English NER dataset [Tjong Kim Sang and De Meulder, 2003].

Tkachenko et al. [2013] described a supervised learning method for named-entity recognition. Feature engineering and learning algorithm selection are critical factors when designing a NER system. Possible features could include word lemmas, part-of-speech tags, and occurrence in some dictionary that encodes characteristic attributes of words relevant for the classification task. Tkachenko et al. [2013] included morphological, dictionary-based, WordNet-based, and global features. For their learning algorithm, the researchers chose CRFs, which have a sequential nature and ability to handle a large number of features. As also mentioned above, CRFs are widely used for the task of NER. For the Estonian dataset, the system produced a gold standard NER corpus, on which their CRF-based model achieved an overall F-score of 0.87.

He and Sun [2017] developed a semi-supervised leaning model based on deep neural networks (B-LSTM). This system combined transition probabilities with deep learning to train the model directly on F-score and label accuracy. The researchers used a modified, labeled corpus which corrected labeling errors in data developed by Peng and Dredze [2016] for NER in Chinese social media. They evaluated their model on NER and nominal mention tasks. The

<sup>12</sup>[http://www.cs.technion.ac.il/~gabr/resources/data/ne\\_datasets.html](http://www.cs.technion.ac.il/~gabr/resources/data/ne_datasets.html)



result for NER on the dataset of Peng and Dredze [2016] is the state-of-the-art NER system in Chinese Social Media. Their Bi-LSTM model achieved an F-score of 0.53.

Approaches based on deep learning were shown to benefit NER systems as well. Aguilar et al. [2017] proposed a multi-task approach by employing the task of Named Entity (NE) segmentation together with the task of fine-grained NE categorization. The multi-task neural network architecture learns higher-order feature representations from word and character sequences along with basic part-of-speech tags and gazetteer information. This neural network acts as a feature extractor to feed a Conditional Random Fields classifier. They obtained the best results in the 3rd Workshop on Noisy User-generated Text (WNUT-2017) with a 0.4186 entity detection F-score. Aguilar et al. [2018] extended the system's architecture and improved the results with 2-3%.

### Evaluation Measures for NER

The precision, recall, and F-measure can be calculated at sequence level (whole span of text) or at token level. The former is stricter because each named entity that is longer than one word has to have an exact start and end point. Once entities have been determined, the accuracy of assigning them to tags such as Person, Organization, etc., can be calculated.

### Adaptation for Named Entity Recognition

Named entity recognition methods typically have 85–90% accuracy on long and carefully edited texts, but their performance decreases to 30–50% on tweets [Li et al., 2012a, Liu et al., 2012b, Ritter et al., 2011].

Ritter et al. [2011] reported that the Stanford NER obtains 44% accuracy on Twitter data. They also presented new NER methods for social media texts based on labeled Latent Dirichlet Allocation (LDA)<sup>13</sup> [Ramage et al., 2009], that allowed their T-Seg NER system to reach an accuracy of 67%.

Derczynski et al. [2013b] reported that NER performance drops from 77% F-score on newspaper text to 60% on Twitter data, and that after adaptation it increases to 80% (with the ANNIE NER system from GATE) [Cunningham et al., 2002]. The performance on newspaper data was computed on the CoNLL 2003 English NER dataset [Tjong Kim Sang and De Meulder, 2003], while the performance on social media data was computed on part of the Ritter dataset [Ritter et al., 2011], which contains of 2,400 tweets comprising 34,000 tokens.

Particular attention is given to microtext normalization, as a way of removing some of the linguistic noise prior to part-of-speech tagging and entity recognition [Derczynski et al., 2013a, Han and Baldwin, 2011]. Some research has focused on named entity recognition algorithms specifically for Twitter messages, training new CRF model on Twitter data [Ritter et al., 2011].

<sup>13</sup>LDA is a method that assumes a number of hidden topics for a corpus, and discovers a cluster of words for each topic, with associated probabilities. Then, for each document, LDA can estimate a probability distribution over the topics. The topics—word clusters—do not have names, but names can be given, for example, by choosing the word with the highest probability in each cluster.