# Neural Network Methods
## for Natural Language Processing

**Yoav Goldberg**

# Neural Network Methods for
# Natural Language Processing

Neural Network Methods for Natural Language Processing

Yoav Goldberg

www.morganclaypool.com

# 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.

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

# Neural Network Methods for Natural Language Processing

Yoav Goldberg
Bar Ilan University

MORGAN & CLAYPOOL PUBLISHERS

## ABSTRACT

Neural networks are a family of powerful machine learning models. This book focuses on the application of neural network models to natural language data. The first half of the book (Parts I and II) covers the basics of supervised machine learning and feed-forward neural networks, the basics of working with machine learning over language data, and the use of vector-based rather than symbolic representations for words. It also covers the computation-graph abstraction, which allows to easily define and train arbitrary neural networks, and is the basis behind the design of contemporary neural network software libraries.

The second part of the book (Parts III and IV) introduces more specialized neural network architectures, including 1D convolutional neural networks, recurrent neural networks, conditioned-generation models, and attention-based models. These architectures and techniques are the driving force behind state-of-the-art algorithms for machine translation, syntactic parsing, and many other applications. Finally, we also discuss tree-shaped networks, structured prediction, and the prospects of multi-task learning.

# Preface

Natural language processing (NLP) is a collective term referring to automatic computational processing of human languages. This includes both algorithms that take human-produced text as input, and algorithms that produce natural looking text as outputs. The need for such algorithms is ever increasing: human produce ever increasing amounts of text each year, and expect computer interfaces to communicate with them in their own language. Natural language processing is also very challenging, as human language is inherently ambiguous, ever changing, and not well defined.

Natural language is symbolic in nature, and the first attempts at processing language were symbolic: based on logic, rules, and ontologies. However, natural language is also highly ambiguous and highly variable, calling for a more statistical algorithmic approach. Indeed, the current-day dominant approaches to language processing are all based on *statistical machine learning*. For over a decade, core NLP techniques were dominated by linear modeling approaches to supervised learning, centered around algorithms such as Perceptrons, linear Support Vector Machines, and Logistic Regression, trained over very high dimensional yet very sparse feature vectors.

Around 2014, the field has started to see some success in switching from such linear models over sparse inputs to *nonlinear neural network models* over dense inputs. Some of the neural-network techniques are simple generalizations of the linear models and can be used as almost drop-in replacements for the linear classifiers. Others are more advanced, require a change of mindset, and provide new modeling opportunities. In particular, a family of approaches based on *recurrent neural networks* (RNNs) alleviates the reliance on the Markov Assumption that was prevalent in sequence models, allowing to condition on arbitrarily long sequences and produce effective feature extractors. These advances led to breakthroughs in language modeling, automatic machine translation, and various other applications.

While powerful, the neural network methods exhibit a rather strong barrier of entry, for various reasons. In this book, I attempt to provide NLP practitioners as well as newcomers with the basic background, jargon, tools, and methodologies that will allow them to understand the principles behind neural network models for language, and apply them in their own work. I also hope to provide machine learning and neural network practitioners with the background, jargon, tools, and mindset that will allow them to effectively work with language data.

Finally, I hope this book can also serve a relatively gentle (if somewhat incomplete) introduction to both NLP and machine learning for people who are newcomers to both fields.

# INTENDED READERSHIP

This book is aimed at readers with a technical background in computer science or a related field, who want to get up to speed with neural network techniques for natural language processing. While the primary audience of the book is graduate students in language processing and machine learning, I made an effort to make it useful also to established researchers in either NLP or machine learning (by including some advanced material), and to people without prior exposure to either machine learning or NLP (by covering the basics from the grounds up). This last group of people will, obviously, need to work harder.

While the book is self contained, I do assume knowledge of mathematics, in particular undergraduate level of probability, algebra, and calculus, as well as basic knowledge of algorithms and data structures. Prior exposure to machine learning is very helpful, but not required.

This book evolved out of a survey paper [Goldberg, 2016], which was greatly expanded and somewhat re-organized to provide a more comprehensive exposition, and more in-depth coverage of some topics that were left out of the survey for various reasons. This book also contains many more concrete examples of applications of neural networks to language data that do not exist in the survey. While this book is intended to be useful also for people without NLP or machine learning backgrounds, the survey paper assumes knowledge in the field. Indeed, readers who are familiar with natural language processing as practiced between roughly 2006 and 2014, with heavy reliance on machine learning and linear models, may find the journal version quicker to read and better organized for their needs. However, such readers may also appreciate reading the chapters on word embeddings (10 and 11), the chapter on conditioned generation with RNNs (17), and the chapters on structured prediction and multi-task learning (MTL) (19 and 20).

# FOCUS OF THIS BOOK

This book is intended to be self-contained, while presenting the different approaches under a unified notation and framework. However, the main purpose of the book is in introducing the neural-networks (deep-learning) machinery and its application to language data, and not in providing an in-depth coverage of the basics of machine learning theory and natural language technology. I refer the reader to external sources when these are needed.

Likewise, the book is not intended as a comprehensive resource for those who will go on and develop the next advances in neural network machinery (although it may serve as a good entry point). Rather, it is aimed at those readers who are interested in taking the existing, useful technology and applying it in useful and creative ways to their favorite language-processing problems.

Further reading    For in-depth, general discussion of neural networks, the theory behind them, advanced optimization methods, and other advanced topics, the reader is referred to other existing resources. In particular, the book by Bengio et al. [2016] is highly recommended.

For a friendly yet rigorous introduction to practical machine learning, the freely available book of Daumé III [2015] is highly recommended. For more theoretical treatment of machine learning, see the freely available textbook of Shalev-Shwartz and Ben-David [2014] and the textbook of Mohri et al. [2012].

For a strong introduction to NLP, see the book of Jurafsky and Martin [2008]. The information retrieval book by Manning et al. [2008] also contains relevant information for working with language data.

Finally, for getting up-to-speed with linguistic background, the book of Bender [2013] in this series provides a concise but comprehensive coverage, directed at computationally minded readers. The first chapters of the introductory grammar book by Sag et al. [2003] are also worth reading.

As of this writing, the progress of research in neural networks and Deep Learning is very fast paced. The state-of-the-art is a moving target, and I cannot hope to stay up-to-date with the latest-and-greatest. The focus is thus with covering the more established and robust techniques, that were proven to work well in several occasions, as well as selected techniques that are not yet fully functional but that I find to be established and/or promising enough for inclusion.

Yoav Goldberg
March 2017

# Acknowledgments

And thanks also to Graeme Hirst, Michael Morgan, Samantha Draper, and C.L. Tondo for orchestrating the effort.

As usual, all mistakes are of course my own. Do let me know if you find any, though, and be listed in the next edition if one is ever made.

Finally, I would like to thank my wife, Noa, who was patient and supportive when I disappeared into writing sprees, my parents Esther and Avner and brother Nadav who were in many cases more excited about the idea of me writing a book than I was, and the staff at The Streets Cafe (King George branch) and Shne'or Cafe who kept me well fed and served me drinks throughout the writing process, with only very minimal distractions.

Yoav Goldberg
March 2017

# Contents