# Natural Language Processing

Jacob Eisenstein

November 13, 2018

# Contents

**IV  Applications**                                                                                  **401**

**17  Information extraction**                                                                         **403**

**18  Machine translation**                                                                           **431**

# Preface

The goal of this text is focus on a core subset of the natural language processing, unified by the concepts of learning and search. A remarkable number of problems in natural language processing can be solved by a compact set of methods:

**Search.** Viterbi, CKY, minimum spanning tree, shift-reduce, integer linear programming, beam search.

**Learning.** Maximum-likelihood estimation, logistic regression, perceptron, expectation-maximization, matrix factorization, backpropagation.

This text explains how these methods work, and how they can be applied to a wide range of tasks: document classification, word sense disambiguation, part-of-speech tagging, named entity recognition, parsing, coreference resolution, relation extraction, discourse analysis, language modeling, and machine translation.

## Background

Because natural language processing draws on many different intellectual traditions, almost everyone who approaches it feels underprepared in one way or another. Here is a summary of what is expected, and where you can learn more:

**Mathematics and machine learning.** The text assumes a background in multivariate calculus and linear algebra: vectors, matrices, derivatives, and partial derivatives. You should also be familiar with probability and statistics. A review of basic probability is found in Appendix A, and a minimal review of numerical optimization is found in Appendix B. For linear algebra, the online course and textbook from Strang (2016) provide an excellent review. Deisenroth et al. (2018) are currently preparing a textbook on *Mathematics for Machine Learning*, a draft can be found online.[1] For an introduction to probabilistic modeling and estimation, see James et al. (2013); for

---

[1] https://mml-book.github.io/

a more advanced and comprehensive discussion of the same material, the classic reference is Hastie et al. (2009).

**Linguistics.** This book assumes no formal training in linguistics, aside from elementary concepts likes nouns and verbs, which you have probably encountered in the study of English grammar. Ideas from linguistics are introduced throughout the text as needed, including discussions of morphology and syntax (chapter 9), semantics (chapters 12 and 13), and discourse (chapter 16). Linguistic issues also arise in the application-focused chapters 4, 8, and 18. A short guide to linguistics for students of natural language processing is offered by Bender (2013); you are encouraged to start there, and then pick up a more comprehensive introductory textbook (e.g., Akmajian et al., 2010; Fromkin et al., 2013).

**Computer science.** The book is targeted at computer scientists, who are assumed to have taken introductory courses on the analysis of algorithms and complexity theory. In particular, you should be familiar with asymptotic analysis of the time and memory costs of algorithms, and with the basics of dynamic programming. The classic text on algorithms is offered by Cormen et al. (2009); for an introduction to the theory of computation, see Arora and Barak (2009) and Sipser (2012).

## How to use this book

After the introduction, the textbook is organized into four main units:

**Learning.** This section builds up a set of machine learning tools that will be used throughout the other sections. Because the focus is on machine learning, the text representations and linguistic phenomena are mostly simple: "bag-of-words" text classification is treated as a model example. Chapter 4 describes some of the more linguistically interesting applications of word-based text analysis.

**Sequences and trees.** This section introduces the treatment of language as a structured phenomena. It describes sequence and tree representations and the algorithms that they facilitate, as well as the limitations that these representations impose. Chapter 9 introduces finite state automata and briefly overviews a context-free account of English syntax.

**Meaning.** This section takes a broad view of efforts to represent and compute meaning from text, ranging from formal logic to neural word embeddings. It also includes two topics that are closely related to semantics: resolution of ambiguous references, and analysis of multi-sentence discourse structure.

**Applications.** The final section offers chapter-length treatments on three of the most prominent applications of natural language processing: information extraction, machine

translation, and text generation. Each of these applications merits a textbook length treatment of its own (Koehn, 2009; Grishman, 2012; Reiter and Dale, 2000); the chapters here explain some of the most well known systems using the formalisms and methods built up earlier in the book, while introducing methods such as neural attention.

Each chapter contains some advanced material, which is marked with an asterisk. This material can be safely omitted without causing misunderstandings later on. But even without these advanced sections, the text is too long for a single semester course, so instructors will have to pick and choose among the chapters.

Chapters 1-3 provide building blocks that will be used throughout the book, and chapter 4 describes some critical aspects of the practice of language technology. Language models (chapter 6), sequence labeling (chapter 7), and parsing (chapter 10 and 11) are canonical topics in natural language processing, and distributed word embeddings (chapter 14) have become ubiquitous. Of the applications, machine translation (chapter 18) is the best choice: it is more cohesive than information extraction, and more mature than text generation. Many students will benefit from the review of probability in Appendix A.

- A course focusing on machine learning should add the chapter on unsupervised learning (chapter 5). The chapters on predicate-argument semantics (chapter 13), reference resolution (chapter 15), and text generation (chapter 19) are particularly influenced by recent progress in machine learning, including deep neural networks and learning to search.

- A course with a more linguistic orientation should add the chapters on applications of sequence labeling (chapter 8), formal language theory (chapter 9), semantics (chapter 12 and 13), and discourse (chapter 16).

- For a course with a more applied focus, I recommend the chapters on applications of sequence labeling (chapter 8), predicate-argument semantics (chapter 13), information extraction (chapter 17), and text generation (chapter 19).

## Acknowledgments

Jacob Eisenstein. Draft of November 13, 2018.

# Notation

As a general rule, words, word counts, and other types of observations are indicated with Roman letters $(a, b, c)$; parameters are indicated with Greek letters $(\alpha, \beta, \theta)$. Vectors are indicated with bold script for both random variables $\boldsymbol{x}$ and parameters $\boldsymbol{\theta}$. Other useful notations are indicated in the table below.

| **Basics** | |
| --- | --- |
| $\exp x$ | the base-2 exponent, $2^x$ |
| $\log x$ | the base-2 logarithm, $\log_2 x$ |
| $\{x_n\}_{n=1}^N$ | the set $\{x_1, x_2, \ldots, x_N\}$ |
| $x_i^j$ | $x_i$ raised to the power $j$ |
| $x_i^{(j)}$ | indexing by both $i$ and $j$ |

| **Linear algebra** | |
| --- | --- |
| $\boldsymbol{x}^{(i)}$ | a column vector of feature counts for instance $i$, often word counts |
| $\boldsymbol{x}_{j:k}$ | elements $j$ through $k$ (inclusive) of a vector $\boldsymbol{x}$ |
| $[\boldsymbol{x}; \boldsymbol{y}]$ | vertical concatenation of two column vectors |
| $[\boldsymbol{x}, \boldsymbol{y}]$ | horizontal concatenation of two column vectors |
| $\boldsymbol{e}_n$ | a "one-hot" vector with a value of $1$ at position $n$, and zero everywhere else |
| $\boldsymbol{\theta}^\top$ | the transpose of a column vector $\boldsymbol{\theta}$ |
| $\boldsymbol{\theta} \cdot \boldsymbol{x}^{(i)}$ | the dot product $\sum_{j=1}^N \theta_j \times x_j^{(i)}$ |
| $\mathbf{X}$ | a matrix |
| $x_{i,j}$ | row $i$, column $j$ of matrix $\mathbf{X}$ |
| $\text{Diag}(\boldsymbol{x})$ | a matrix with $\boldsymbol{x}$ on the diagonal, e.g., $\begin{pmatrix} x_1 & 0 & 0 \\ 0 & x_2 & 0 \\ 0 & 0 & x_3 \end{pmatrix}$ |
| $\mathbf{X}^{-1}$ | the inverse of matrix $\mathbf{X}$ |

**Text datasets**

| | |
|---|---|
| $w_m$ | word token at position $m$ |
| $N$ | number of training instances |
| $M$ | length of a sequence (of words or tags) |
| $V$ | number of words in vocabulary |
| $y^{(i)}$ | the true label for instance $i$ |
| $\hat{y}$ | a predicted label |
| $\mathcal{Y}$ | the set of all possible labels |
| $K$ | number of possible labels $K = |\mathcal{Y}|$ |
| $\square$ | the start token |
| $\blacksquare$ | the stop token |
| $\boldsymbol{y}^{(i)}$ | a structured label for instance $i$, such as a tag sequence |
| $\mathcal{Y}(\boldsymbol{w})$ | the set of possible labelings for the word sequence $\boldsymbol{w}$ |
| $\diamond$ | the start tag |
| $\blacklozenge$ | the stop tag |

**Probabilities**

| | |
|---|---|
| $\Pr(A)$ | probability of event $A$ |
| $\Pr(A \mid B)$ | probability of event $A$, conditioned on event $B$ |
| $\mathrm{p}_B(b)$ | the marginal probability of random variable $B$ taking value $b$; written $\mathrm{p}(b)$ when the choice of random variable is clear from context |
| $\mathrm{p}_{B|A}(b \mid a)$ | the probability of random variable $B$ taking value $b$, conditioned on $A$ taking value $a$; written $\mathrm{p}(b \mid a)$ when clear from context |
| $A \sim p$ | the random variable $A$ is distributed according to distribution $p$. For example, $X \sim \mathcal{N}(0, 1)$ states that the random variable $X$ is drawn from a normal distribution with zero mean and unit variance. |
| $A \mid B \sim p$ | conditioned on the random variable $B$, $A$ is distributed according to $p$.[2] |

**Machine learning**

| | |
|---|---|
| $\Psi(\boldsymbol{x}^{(i)}, y)$ | the score for assigning label $y$ to instance $i$ |
| $\boldsymbol{f}(\boldsymbol{x}^{(i)}, y)$ | the feature vector for instance $i$ with label $y$ |
| $\boldsymbol{\theta}$ | a (column) vector of weights |
| $\ell^{(i)}$ | loss on an individual instance $i$ |
| $L$ | objective function for an entire dataset |
| $\mathcal{L}$ | log-likelihood of a dataset |
| $\lambda$ | the amount of regularization |