

# まえがき

自然言語処理 (Natural Language Processing: NLP) は、人間の言語に対する自動的な計算処理全般を意味する集合的な用語である。ここには、人間が産出したテキストを入力とするアルゴリズムと、自然なテキストを出力とするアルゴリズムの両方が含まれる。そのようなアルゴリズムへの要求は変わることなく増え続けている。人々は毎年ますます多くの量のテキストを産出し続けているし、自分自身の言語で対話できるようなコンピュータ操作を期待している。そして、自然言語処理は非常に困難でもある。人間の言語は本質的に曖昧で、変化し続け、適切に定義されていないためである。

自然言語は生来的にシンボリックであるので、それを処理しようという最初の試みも、論理や規則やオントロジーに基づいたシンボリックなものであった。一方で、自然言語はとても曖昧で、ひどく定まらないものであったので、より統計的なアルゴリズムによるアプローチが必要とされた。実際、自然言語処理において昨今の支配的な手法は、全て統計的機械学習 (statistical machine learning) に基づいている。ここ 10 年以上にわたって、核となる自然言語処理手法は、線形モデルによる教師あり学習というアプローチに席卷されていた。それらの中心には、パーセプトロン、線形サポートベクトルマシン、ロジスティック回帰があり、非常に高い次元を持つスパース (疎ら) な素性ベクトルによる訓練が行われていた。

2014 年頃、この分野において、そのようなスパースな入力を扱う線形モデルから、密な入力を扱う非線形ニューラルネットワーク (nonlinear neural network models) への転換による成功例が散見されるようになった。いくつかのニューラルネットワーク技法は、線形モデルの単純な一般化で、ニューラルネットワークは線形分類器からの容易に取り替え可能な代替品として用いられた。一方で、より先進的で、考え方の変革を求め、新しいモデリングの可能性を提供するものもあった。特に、再帰的ニューラルネ

ットワーク (Recurrent Neural Networks: RNN) に基づく方法の一群は系列モデルで席捲していたマルコフ仮定への依存を軽減し、任意の長さの系列によって条件付けを行うことと、優れた素性抽出器を構成することとを可能にした。これらの進展は言語モデル、自動機械翻訳、そしてその他の様々な応用においてブレイクスルーをもたらした。

ニューラルネットワークの手法は、強力ではあるが、様々な理由から、それに取り組み始めるのに比較的高い障壁がある。本書を通じて、NLP に従事している専門家と、そして入門者に向けて、基本的な背景、専門用語、ツール、方法論を提供しようと試みている。それにより、言語を扱うニューラルネットワークモデルの原理を理解し、それを彼ら自身の研究や仕事に応用できるようになるだろう。あわせて、機械学習やニューラルネットワークの専門家たちにも、言語データを効果的に扱うことを可能にするような、背景、専門用語、ツール、そして考え方を提供できればと思う。

最後に、本書が、自然言語処理と機械学習のいずれについても専門知識を持たない人たちに対しても、その両方への親切な（いくぶん不完全かもしれないが）入門書となってくれればと望んでいる。

## 意図している読者層

本書は、自然言語処理のためのニューラルネットワーク技法について速習しようという、計算機科学やその関連分野の技術的背景のある読者を対象としている。主たる読者対象は、自然言語処理と機械学習を学ぶ大学院生であるが、すでに実績を持つような、自然言語処理あるいは機械学習の研究者にとっても（上級的话题をいくつか含めることで）、そして今まで機械学習やNLPに接したことのない人々にも（基礎を徹底的に論じることで）有用であるようにと努力している。もちろん、最後のグループは厳しい努力が必要である。

本書は自己完結してはいるが、数学、特に大学生レベルの確率、代数、微積分、そして、アルゴリズムとデータ構造についての基本的な知識を前提としている。事前に機械学習に接していることは大きな助けとはなるが、必須ではない。

本書は概説論文 [Goldberg, 2016] を発展させたもので、より徹底的な説明を与え、様々な理由から概説論文では扱えなかったいくつかのトピックを深く論じるために、大きく拡充するとともに、ある種の再構成を行っている。加えて、本書は、概説論文にはなかった、言語データに対するニューラルネットワークの応用についてのより具体的な例を含んでいる。本書がNLPや機械翻訳の背景を持たない人達にも有益となるように意図されているのに対し、概説論文はそれらの分野の知識を仮定したものとなっている。実際、およそ2006年から2014年に実践された、機械学習と線形モデルが重視された自然言語処理に詳しい読者は、概説論文の版の方が、早く読め、必要に応える優

れた構成となっていると感じるかもしれない。それでも、そのような読者も、単語埋め込みの章（第 10 章と第 11 章）、RNN を用いた条件付き生成の章（第 17 章）、そして、構造予測とマルチタスク学習 (Multi-Task Learning: MTL) の章（第 19 章と第 20 章）は読む価値があると思ってくれるかもしれない。

## 本書の焦点

本書は自己完結であるように意図されており、様々な手法が統一的な記法と枠組みを用いて紹介されている。しかし、本書の主たる目的はニューラルネットワーク（深層学習）の仕組みとその言語データへの応用を紹介することであって、機械学習理論や自然言語技術の基礎について深い説明を提供することではない。それらが必要な場合には、外部の情報を参照するように指示している。

同様に、本書はニューラルネットワークの仕組みについて、研究を進め、次なる進展を切り拓こうという人々にとっての網羅的な資源となることも意図していない（ただし、その良い入り口にはなるかもしれない）。むしろ、既存の有益な技術を取りあげ、それに関心のある言語処理の問題に、有益かつ創造的な方法で適用してみようというような興味を持つ読者に向けた著作である。

補足的文献 ニューラルネットワーク、その背後にある理論、最適化手法、その他の上級のトピックについての深い、一般的な議論については、他の既存の資源を参照して欲しい。特に Bengio et al. [2016] による書籍を強くお勧めする。

実用的な機械学習についての親しみやすいが厳格な入門としては、無償で入手できる Daumé III [2015] の書籍を強くお勧めする。機械学習をより理論的に扱っているものとしては、無償で入手できる Shalev-Shwartz and Ben-David [2014] の教科書と Mohri et al. [2012] の教科書を参照のこと。

NLP への圧倒的な入門としては、Jurafsky and Martin [2008] を参照のこと。Manning et al. [2008] による情報検索の書籍も言語データを扱うために必要な情報を含んでいる。

最後に、言語学的背景について速習しようとするのであれば、原著シリーズの Bender [2013] の書籍が簡潔でいて網羅的な内容であり、情報处理的な考え方になじんだ読者向けの記述となっている。Sag et al. [2003] の入門的な文法書の最初の章も読む価値がある。

本書執筆の間も、ニューラルネットワークと深層学習の研究の進歩は恐ろしく速く、最先端は動く標的であり、最新最良のものに遅れないでいることを望むべくもない。そ

のため、本書は、様々な機会であまく動作することがすでに証明されたような、より確立した頑健な手法を扱うことに焦点を置いている。さらに、完全に機能するかは明らかではないが、そのうち確立する、かつ/もしくは、含めるに値するほど期待できると筆者が考えて選んだ手法が含まれている。

2017年3月

Yoav Goldberg

# 謝 辞

本書は、同じトピックについて私が著した概説論文 [Goldberg, 2016] から生まれたものである。そして、その概説論文は、私が学ぼうとした際、そして学生や協力者に教えようとした際、深層学習と自然言語処理とが交わる領域についての上手に構成された明快な素材がなかったという私の不満から生まれたものである。そのため、私は、概説論文に（最初のドラフトから出版後のコメントまで、様々な形式で）コメントをしてくれた沢山の人達に、そして、書籍のドラフトのいろいろな段階にコメントしてくれた数多くの人達に恩義がある。あるコメントは、直接伺ったし、あるものは電子メール経由で、そして、あるものはツイッターの入り乱れた会話の中でなされた。本書それ自体にはコメントしていない（実際、何人かは読んでさえいない）けれども、関連したトピックについて議論してくれた人達の影響も受けている。何人かは、深層学習の専門家、何人かは NLP の専門家、何人かは両者に詳しい。そして、その他はこれらのトピックについて学ぼうとしている人達であった。何人かは（多くはないが）本当に詳細なコメントをして寄与してくれた。その他の人は小さな細部について議論し、多くはその中間である。しかし、彼らの全員が本書の最終版に影響を与えている。彼らは、アルファベット順で、Yoav Artzi, Yonatan Aumann, Jason Baldridge, Miguel Ballesteros, Mohit Bansal, Marco Baroni, Tal Baumel, Sam Bowman, Jordan Boyd-Graber, Chris Brockett, Ming-Wei Chang, David Chiang, Kyunghyun Cho, Grzegorz Chrupala, Alexander Clark, Raphael Cohen, Ryan Cotterell, Hal Daumé III, Nicholas Dronen, Chris Dyer, Jacob Eisenstein, Jason Eisner, Michael Elhadad, Yad Faeq, Manaal Faruqui, Amir Globerson, Frédéric Godin, Edward Grefenstette, Matthew Honnibal, Dirk Hovy, Moshe Koppel, Angeliki Lazaridou, Tal Linzen, Thang Luong, Chris Manning, Stephen Merity, Paul Michel, Margaret Mitchell, Piero Molino, Graham

Neubig, Joakim Nivre, Brendan O'Connor, Nikos Pappas, Fernando Pereira, Barbara Plank, Ana-Maria Popescu, Delip Rao, Tim Rocktäschel, Dan Roth, Alexander Rush, Naomi Saphra, Djamel Seddah, Erel Segal-Halevi, Avi Shmidman, Shlital Shmidman, Noah Smith, Anders Søgaard, Abe Stanway, Emma Strubell, Sandeep Subramanian, Liling Tan, Reut Tsarfaty, Peter Turney, Tim Vieira, Oriol Vinyals, Andreas Vlachos, Wenpeng Yin, そして, Torsten Zesch.

このリストにはもちろん, このトピックについてのその人たちの学術的著作を通じて私が会話した多くの研究者を含めていない.

本書は, また, Bar-Ilan University の the Natural Language Processing Group (そして, その緩やかな広がり) との交流から多くを得ており, それによって形作られた. 彼らは, Yossi Adi, Roei Aharoni, Oded Avraham, Ido Dagan, Jessica Fidler, Jacob Goldberger, Hila Gonen, Joseph Keshet, Eliyahu Kiperwasser, Ron Konigsberg, Omer Levy, Oren Melamud, Gabriel Stanovsky, Ori Shapira, Micah Shlain, Vered Shwartz, Hillel Taub-Tabib, そして, Rachel Wities. 多くの人はこれら二つのリストの両方に属するが, 重複を省いてリストを短くしている.

本書と概説論文の匿名の査読者の方達も, 名前が出ないにも関わらず (そして時には煩わしかったと思われるが), 信頼に足るコメント, 示唆, 訂正を行ってくれた. この最終作品が多くの面で劇的に改善されたのは間違いないと自信を持って言うことができる. 誰であろうと, ありがとう. そして, Graeme Hirst, Michael Morgan, Samantha Draper, C.L. Tondo の組織力にも感謝する.

いつも通り, 全ての誤りは私自身の責任である. それでも, もし見つけたら知らせしてほしい. もしあればだが, 次の版ではそれらを一覧にするようにしたい.

最後に, 妻 Noa に感謝する. 彼女は, 私が執筆の大騒ぎの中に紛れ込んでしまった時でも, 忍耐強く, 協力的であった. 両親 Esther と Avne, そして, 兄弟 Nadav にも感謝する. 彼らは, しばしば, 私が書籍を執筆するという発想に私以上に興奮していた. そして, The Streets Cafe (King George 支店) と Shne'or Cafe のスタッフにも感謝する. 彼らは, 私の気を散らすことは最小限に留めつつ, 執筆の期間を通じて私のお腹を良い具合に満たし, 飲み物を提供してくれた.

2017 年 3 月

Yoav Goldberg