

Jochen HIRSCHLE



Deep Natural Language Processing

Einstieg in Word Embedding,
Sequence-to-Sequence-Modelle
und Transformer mit Python



Online:
GitHub-Repository zum Buch

HANSER

Für Tuuli-Marja

HANSER

Jochen Hirsche

Deep Natural Language Processing

Einstieg in Word Embedding, Sequence-to-Sequence-Modelle
und Transformers mit Python

Der Autor:

Dr. Jochen Hirschle, Braunschweig

Alle in diesem Buch enthaltenen Informationen, Verfahren und Darstellungen wurden nach bestem Wissen zusammengestellt und mit Sorgfalt getestet. Dennoch sind Fehler nicht ganz auszuschließen. Aus diesem Grund sind die im vorliegenden Buch enthaltenen Informationen mit keiner Verpflichtung oder Garantie irgendeiner Art verbunden. Autor und Verlag übernehmen infolgedessen keine juristische Verantwortung und werden keine daraus folgende oder sonstige Haftung übernehmen, die auf irgendeine Art aus der Benutzung dieser Informationen – oder Teilen davon – entsteht.

Ebenso übernehmen Autor und Verlag keine Gewähr dafür, dass beschriebene Verfahren usw. frei von Schutzrechten Dritter sind. Die Wiedergabe von Gebrauchsnamen, Handelsnamen, Warenbezeichnungen usw. in diesem Buch berechtigt deshalb auch ohne besondere Kennzeichnung nicht zu der Annahme, dass solche Namen im Sinne der Warenzeichen- und Markenschutz-Gesetzgebung als frei zu betrachten wären und daher von jedermann benutzt werden dürften.

Bibliografische Information der Deutschen Nationalbibliothek:

Die Deutsche Nationalbibliothek verzeichnet diese Publikation in der Deutschen Nationalbibliografie; detaillierte bibliografische Daten sind im Internet über <http://dnb.d-nb.de> abrufbar.

Dieses Werk ist urheberrechtlich geschützt.

Alle Rechte, auch die der Übersetzung, des Nachdruckes und der Vervielfältigung des Buches, oder Teilen daraus, vorbehalten. Kein Teil des Werkes darf ohne schriftliche Genehmigung des Verlages in irgendeiner Form (Fotokopie, Mikrofilm oder ein anderes Verfahren) – auch nicht für Zwecke der Unterrichtsgestaltung – reproduziert oder unter Verwendung elektronischer Systeme verarbeitet, vervielfältigt oder verbreitet werden.

© 2022 Carl Hanser Verlag München, www.hanser-fachbuch.de

Lektorat: Sylvia Hasselbach

Copy editing: Sandra Gottmann, Wasserburg

Umschlagdesign: Marc Müller-Bremer, www.rebranding.de, München

Umschlagrealisation: Max Kostopoulos

Titelmotiv: © gettyimages.de/MR.Cole_Photographer

Layout: Manuela Treindl, Fürth

Ausstattung patentrechtlich geschützt. Kösel FD 351, Patent-Nr. 0748702

Print-ISBN: 978-3-446-47363-8

E-Book-ISBN: 978-3-446-47390-4

E-Pub-ISBN: 978-3-446-47409-3

Inhalt

Titelei

Impressum

Inhalt

1 Einleitung

2 Textdaten verarbeiten und vorverarbeiten

2.1 Grundlegende Techniken der Verarbeitung von Textdaten

2.2 Mit NumPy arbeiten

2.3 One-Hot-Encodierung und Bag-of-Words-Modell

3 Grundlagen maschinellen Lernens

3.1 Lineare Regression

3.1.1 Eine Gerade in eine Punktwolke legen

3.1.2 Die Lage der Geraden bestimmen

3.1.3 Die Qualität eines Modells bestimmen

3.1.4 Multivariate Regression

3.1.5 Praktische Umsetzung mit Python und Scikit-Learn

3.2 Logistische Regression

3.2.1 Verfahrensweise

3.2.2 Gütemaße

3.2.3 Praktische Umsetzung mit Scikit-Learn

3.3 Softmax-Regression

3.3.1 Verfahrensweise

3.3.2 Praktische Umsetzung mit Scikit-Learn

4 Einfache Verfahren zur Vektorisierung von Textdaten

4.1 One-Hot-Encodierung und Bag-of-Words-Ansatz

4.2 N-grams

4.3 TF-IDF-Vektorisierung

4.4 Umsetzung mit Scikit-Learn

4.4.1 Vektorisierung mit dem Count-Vectorizer

4.4.2 TF-IDF-Vektorisierung

4.4.3 Lemmatisierung

4.4.4 Einsatz eines N-gram-Modells

5 Deep Learning-Essentials

5.1 Neuronen und neuronale Netze

5.2 Wie neuronale Netze lernen

5.3 Architektur und Einstellungen eines neuronalen Netzes

5.3.1 Anzahl der Neuronen in der ersten aktiven Schicht

5.3.2 Anzahl der Neuronen in der Ausgabeschicht

5.3.3 Aktivierung der Neuronen der Ausgabeschicht

5.3.4 Auswahl einer passenden Verlustfunktion

5.3.5 Wahl des Optimierers

5.3.6 Aktivierung der Neuronen in der verdeckten Schicht

5.4 Ein neuronales Netz mit TensorFlow und Keras aufbauen und anlernen

5.4.1 Standardisierung der Features

5.4.2 Aufbau und Einstellungen eines neuronalen Netzes

5.4.3 Anlernen des Modells

5.4.4 Steuerung des Anlernprozesses (Early Stopping)

5.5 Generalisierung und Überanpassung

5.5.1 Regularisierung

5.5.2 Dropout

5.5.3 Praktische Umsetzung

6 Rekurrente Netze

6.1 Aufbau und Funktionsweise rekurrenter Netze

6.2 Long Short Term Memory (LSTM) und Gated Recurrent Units (GRU)

6.3 Praxis rekurrenter Netze: eine automatische Rechtschreibkorrektur

6.3.1 Umsetzung der Encodierung

6.3.2 Aufbau und Anlernen des rekurrenten Netzes

6.3.3 Mit einem bidirektionalen rekurrenten Layer arbeiten

6.4 Anlernen neuronaler Netze mit Generatoren

6.4.1 Generatoren und Generator-Funktionen in Python

6.4.2 Daten batchweise ziehen

6.4.3 Neuronale Netze mit Generatoren anlernen

6.4.4 Die Rechtschreibkorrektur mit einem Generator anlernen

7 Konvolutionale Netze

7.1 Funktionsweise konvolutionaler Netze

7.2 Sequenzdaten mit konvolutionalen Netzen verarbeiten

7.3 Praxis des Anlernens eines konvolutionalen Netzes mit Textdaten

8 Word Embedding

8.1 Funktionsweise

8.2 Aufgabenübergreifende semantische Räume: word2vec- und fastText-Verfahren

8.3 Mit Word Embedding-Verfahren in der Praxis arbeiten

8.3.1 Vorverarbeitung und Implementierung mit Keras

8.3.2 Der Heidegger-Algorithmus: ein generatives Modell zur Erzeugung von Texten

8.3.2.1 Aufbau eines generativen Modells

8.3.2.2 Vorbereitung der Daten

8.3.2.3 Aufbau und Anlernen des Netzes

8.3.2.4 Texte erzeugen

8.3.2.5 Synonyme Wörter identifizieren

8.4 Mit vortrainierten Worteinbettungen arbeiten (fastText)

8.4.1 fastText-Vektorräume aufbereiten

8.4.2 Austausch der Gewichte eines Embedding Layers

8.4.3 Den Vektorraum um unbekannte Wörter erweitern

9 Komplexe Lernarchitekturen umsetzen

9.1 Die funktionale API von TensorFlow

9.2 Ein Modell mit zwei Eingängen aufbauen und anlernen

9.2.1 Architektur des Modells

9.2.2 Anlernen des Modells

10 Sequence-to-Sequence-Modelle

10.1 Encoder-Decoder-Modelle mit Teacher Forcing

10.2 Attention-Mechanismus

10.3 Encoder-Decoder-Architekturen in der Praxis

10.3.1 Ein einfaches Encoder-Decoder-Modell

10.3.1.1 Vorbereitung der Daten

10.3.1.2 Aufbau des Encoder-Decoder-Modells

10.3.1.3 Das Inferenzmodell aufbauen und einsetzen

10.3.2 Encoder-Decoder-Modelle mit Attention-Mechanismus

10.3.2.1 Vorbereitung der Daten

10.3.2.2 Zusammenstellung des neuronalen Netzes

10.3.2.3 Anlernen des Modells

10.3.2.4 Aufbau des Inferenzmodells

10.3.2.5 Das Modell für Übersetzungen einsetzen

11 Transformers

11.1 Aufbau und Funktionsweise

11.1.1 Self-Attention

11.1.2 Die Transformer-Architektur

11.2 Subwort-Tokenisierung

11.3 Mit der Hugging Face-Bibliothek arbeiten

11.3.1 Hauptklassen der Transformers-Bibliothek

11.3.2 Mit der Hugging Face-Pipeline arbeiten

11.3.3 Mit der Tokenizer-Klasse arbeiten

11.3.4 Mit der Model-Klasse arbeiten

11.3.5 Fine Tuning vortrainierter Netze

11.3.5.1 Ein vortrainiertes Modell mit einem nichttrainierten Kopf laden

11.3.5.2 Eine Durchleitung organisieren

11.3.5.3 Teile des Netzes auf nichttrainierbar stellen

11.3.5.4 Das Modell anlernen

12 Diskussion und Ausblick

Literaturverzeichnis

1 Einleitung

Wer die Berichterstattung zum Thema künstliche Intelligenz in den Medien verfolgt, dem wird aufgefallen sein, dass die erste Euphorie inzwischen deutlich abgeklungen ist. Das liegt nicht nur daran, dass sich aufgrund der Vielzahl an Berichten ein Gewöhnungseffekt eingestellt hat. Auch haben in vielen Bereichen die Entwicklungen den hohen Erwartungen nicht standgehalten.

Anstatt es uns auf den Rücksitzen autonomer Fahrzeuge bequem zu machen, wie von vielen ExpertInnen noch vor wenigen Jahren vollmundig angekündigt¹, steuern wir unsere Autos im Jahr 2022 noch immer gestresst und von Hand durch die Staus der Innenstädte. Auch um die allgemeine künstliche Intelligenz (*General AI*) ist es inzwischen verdächtig still geworden. Dabei jagte zuvor ein aufsehenerregendes KI-Ereignis das nächste. Man denke nur an IBMs *Watson*, der medienwirksam mehrfach in der Quizshow *Jeopardy* gegen die Champions gewinnen konnte [Franklin2014]. Oder auch Googles *Alpha Go*, das im gleichnamigen Brettspiel gegen den weltbesten Go-Spieler antrat und ihn besiegte. Viele glaubten, dass die Entwicklungen in diesem atemberaubenden Tempo weitergehen würden und dass die Zukunftsvisionen aus Science-Fiction-Filmen und Serien, in

denen Humanoide von Menschen kaum noch zu unterscheiden sind, schon bald in greifbare Nähe rücken würden.

Tatsächlich ist der Alltag künstlicher Intelligenz heute wesentlich grauer geworden. Schlagzeilen machen KIs heute eher, weil sie von zweifelhaften Regimen zur umfassenden Überwachung von Menschen eingesetzt werden oder weil sie mit rassistischen oder sexistischen Verfehlungen auffallen [West2019]. Und die ersten Kontakte mit den kleinen Schwestern der Film-Supercomputer verlaufen auch nicht immer so, wie wir uns das vorgestellt haben. Wer hat nicht schon holprige Konversationen mit Alexa geführt, sich mit automatischen Antwortsystemen am Telefon herumgeschlagen, sinnlose Wortvervollständigungen auf Smartphones rückgängig gemacht oder ist mit dem Chatbot auf der Seite einer Bank oder Versicherung aneinandergeraten?

Womöglich hat Judea Pearl (ein Pionier der KI-Entwicklung) recht, wenn er etwas spöttisch zu Protokoll gibt, dass die eindrucksvollen Ergebnisse intelligenter Systeme im Wesentlichen auf *Curve Fitting* zurückgingen. Was er damit meint, ist, dass KI-Systeme heute zwar in der Lage sind, beliebig komplexe Regelmäßigkeiten in Daten aufzuspüren. Was sie dagegen nicht können, ist, zwischen kausalen Beziehungen und bloßen Korrelationen zu unterscheiden [Pearl2018]. Dadurch, so Pearl, seien die Möglichkeiten, weitere signifikante Fortschritte in diesem Bereich zu machen, eher begrenzt – die Revolution fällt also womöglich aus [Marcus2018].

Ob das stimmt oder nicht, sei dahingestellt. Gewiss ist allerdings, dass der Siegeszug künstlich intelligenter Systeme trotz aller Beschränktheit nicht aufzuhalten ist. Auch wenn die Anwendungen vielleicht nicht den Utopien oder Dystopien aus den Filmen nahekommen. Systeme, die Machine- und Deep-

Learning-Verfahren implementieren, haben sich in den letzten Jahren in rasantem Tempo ausgebreitet. Kaum ein Geschäftsfeld, in dem maschinelle Lernalgorithmen nicht zum Einsatz kommen. Sei es im Marketing, im Journalismus, im Finanzwesen oder in der verarbeitenden Industrie [James2021a].

Dabei handelt es sich meist um kleinteilige, problemzentrierte Anwendungen, um *Narrow AI* im Vergleich zu *General AI*. Solche Algorithmen erledigen spezifische Aufgaben für uns: sie übersetzen, produzieren Textzusammenfassungen, handeln an der Börse, warnen vor möglichen Ausfällen eines Aggregats, machen auf Hackerangriffe aufmerksam, generieren Antworten auf Fragen oder bestimmen, welche Videos und welche Werbung wir sehen und welche Posts in unseren Newsfeeds landen.

Wie umfassend diese Vereinnahmung bereits ist, lässt sich nur schwer sagen. Meist können wir von außen nicht sicher sein, ob eine Anwendung von einem maschinellen Lernsystem betrieben wird oder nicht. Sicher ist nur, dass die Technologie bereits heute Teil unseres Alltags ist und aus diesem Alltag nicht mehr verschwinden wird. Wir werden uns daran gewöhnen müssen, in Zukunft noch häufiger mit selbstlernenden Automaten zu interagieren und womöglich mit ihnen zu kooperieren. Viele Berufsbilder werden unter Anpassungsdruck geraten, sich wandeln, und Tätigkeiten, für die bis vor Kurzem noch menschliche Arbeitskraft und kognitive Fertigkeiten notwendig waren, werden ganz oder teilweise von Algorithmen übernommen [Acemoglu2018].

Dieses Buch handelt von dieser Art künstlicher Intelligenz. Von Anwendungen, die im Hintergrund operieren. Keine davon wird vermutlich bei Ihnen das Gefühl auslösen, dass Sie es mit einem künstlichen Bewusstsein à la *HAL* oder einem der Androiden aus

Westworld zu tun haben. Das heißt aber nicht, dass diese Anwendungen nicht auf ihre Weise intelligent und faszinierend sind.

Automaten, die menschliche Sprache imitieren, hinterlassen bei den meisten Menschen einen etwas ambivalenten Eindruck. Dass ein Algorithmus einen Satz produziert, der aus der Feder eines Philosophen oder Dichters stammen könnte, ist zumindest gewöhnungsbedürftig. Je mehr man aber hinter die Kulissen blickt und versteht, wie solche Systeme arbeiten, was sie lernen und was sie nicht lernen und wovon ihre Leistung abhängt, umso weniger mysteriös wirken sie.

Was Ihnen dieses Buch bietet

Dieses Buch bietet eine Einführung in den Bereich der Verarbeitung natürlicher Sprachdaten (Natural Language Processing) mit Deep Learning-Verfahren. Es richtet sich an alle LeserInnen, die sich für diese Schnittmenge aus linguistischer Analyse und maschinellem Lernen interessieren. An Personen, die genauer wissen möchten, wie sich Textdaten mit den neuesten statistischen Verfahren verarbeiten lassen und wie Übersetzungen, Textvervollständigungen, Textzusammenfassungen, Sentiment-Analysen und andere textbasierte Anwendungen in der Theorie funktionieren und wie sie sich in der Praxis implementieren lassen.

Um mit diesem Buch arbeiten zu können, sollten Sie ein wenig Programmiererfahrung mit Python mitbringen. Sie sollten wissen, wie Sie eine virtuelle Umgebung mit *Anaconda* aufsetzen, wie sich Pakete installieren lassen, und Sie sollten mit der Arbeit mit Klassen, Funktionen, Schleifen und Standardbehältern wie

Listen, Dictionaries, Tuples und Sets vertraut sein. Außerdem brauchen Sie ein paar grundlegende Konzepte aus der Statistik. So sollten Sie sich unter Begriffen wie *Variable*, *Verteilung*, *Mittelwert* und *Standardabweichung* etwas vorstellen können und generell keine größeren Phobien vor Zahlen haben. Viel mehr brauchen Sie nicht.

Dies ist kein Buch, in dem es in erster Linie um die mathematischen Grundlagen des Deep Learnings geht. Der Text kommt in weiten Teilen ohne Formeln aus, obwohl wir sie manchmal der Vollständigkeit halber abdrucken. Im Vordergrund stehen textbasierte Erklärungen der wichtigsten Konzepte und eine Umsetzung dieser Konzepte mithilfe gängiger Python-Frameworks wie *TensorFlow/Keras*, *Scikit-Learn* oder *Transformers*. Wenn Sie im Detail wissen möchten, wie die Algorithmen mathematisch implementiert sind, sollten Sie parallel ein anderes Buch lesen. Wenn Ihnen mathematische Formeln generell näherstehen als textbasierte Erklärungen, ist dieses Buch vermutlich eher nicht die richtige Wahl für Sie.

Wenn Sie sich aber dafür entscheiden, dieses Buch zu lesen, haben Sie am Ende eine ziemlich klare Vorstellung darüber, wie moderne Verfahren in der statistischen Sprachanalyse arbeiten. Sie lernen nicht nur die Funktionsweise, die Möglichkeiten und Grenzen dieser Verfahren kennen. Sie verstehen auch die Grundprinzipien rekurrenter und konvolutionaler Netze und wissen, wie moderne Architekturen wie Sequence-to-Sequence-Modelle oder Transformer arbeiten. Außerdem ist Ihnen klar, welche Verfahren für welche Aufgaben geeignet sind, und Sie können die meisten dieser Verfahren selbst implementieren. Sie können eigene neuronale Netze aufbauen und anlernen, um Texte zu klassifizieren, zu generieren oder in eine andere Sprache zu übersetzen.

Darüber hinaus bietet Ihnen dieses Buch einen Einstieg in die Verwendung neuronaler Netze, die von anderen vortrainiert und zur Verfügung gestellt werden. Sie lernen, wie Sie allgemeine Sprachmodelle, die über sehr große Datenmengen angelernt wurden, laden und einsetzen können, und Sie lernen, wie Sie diese Netze nachtrainieren können, um sie für spezifische Aufgaben nutzbar zu machen.

Welche Inhalte Sie erwarten

Dieses Buch bietet eine schrittweise und kompakte Einführung in den Themenkomplex Machine Learning/Natural Language Processing (NLP). Die Kapitel bauen inhaltlich aufeinander auf. Wenn Sie also eines der Kapitel lesen, sollten Sie die Inhalte aus den vorherigen Kapiteln kennen. Konzepte, Klassen oder Funktionen werden in der Regel nur an einer Stelle erklärt, und Sie können davon ausgehen, dass einmal eingeführtes Wissen später wieder gebraucht wird. Wenn Sie neu in diesem Gebiet sind, ist es daher ratsam, die Kapitel nacheinander zu lesen, damit das didaktische Konzept aufgeht. Wenn Sie schon über Kenntnisse im einen oder anderen Bereich verfügen, können Sie natürlich einzelne Kapitel oder auch mehrere Kapitel überspringen.

Wir starten unsere Reise in die Welt des Deep Learnings mit einer allgemeineren Einführung in den Bereich maschinellen Lernens. Bevor wir uns an komplexen Modellen abarbeiten, sollten wir eine konkrete Vorstellung davon haben, wie einfache Modelle funktionieren. Wir machen uns also mit den Grundlagen statistischer Lernverfahren und den Basics im Umgang mit Daten in Python vertraut.

Deshalb beginnen wir im Anschluss an die Einleitung im **zweiten Kapitel** mit einer Einführung in die Behälterklassen, die Python standardmäßig bietet, und in einige externe Pakete und Funktionen wie *NumPy*, die uns die Arbeit mit Textdaten und Matrizen erleichtern. Das **dritte Kapitel** widmet sich den allgemeinen Grundlagen maschinellen Lernens. Wir schauen uns an, wie einfache Lernzellen arbeiten, wie man Beziehungen zwischen x - und y -Variablen modelliert, was es mit der Optimierung von Gewichten und dem Gradientenabstiegsverfahren auf sich hat und wie sich mithilfe von Aktivierungsfunktionen beliebige Zielwerte (stetige und kategoriale Werte) produzieren lassen.

In **Kapitel 4** steigen wir dann in die Analyse von Textdaten ein. Dabei geht es zunächst um bewährte Vorverarbeitungsverfahren wie den Bag-of-Words-Ansatz, der es ermöglicht, Textdaten auf sehr basale Art mit maschinellen Standardverfahren zu verarbeiten. **Kapitel 5** bietet im Anschluss einen Einstieg in das eigentliche Zielthema des Buches: in die Arbeit mit neuronalen Netzen. Auch hier müssen wir uns zuerst einiger grundlegender Verfahren versichern, ehe wir uns um die speziellen, auf die Interpretation und Produktion von Texten zugeschnittenen Architekturen kümmern. In diesem Kapitel geht es daher um die Basics, die für die praktische Arbeit mit neuronalen Netzen erforderlich sind. Angefangen von der Funktionsweise einzelner Neuronen, über Schichten von Neuronen, das Durchschleusen von Daten durch diese Schichten bis hin zur Verlustfunktion und der Optimierung der Gewichte mithilfe des Backpropagation-Mechanismus. Danach wissen Sie, wie ein neuronales Netz lernt, wie Gewichte initialisiert und optimiert werden und welche Aufgaben Ihnen bei der Zusammenstellung eines solchen Netzes zukommen.

Mit diesem Wissen ausgerüstet, erörtern wir in **Kapitel 6** eine erste, für die Verarbeitung von Textdaten konzipierte Netzarchitektur: rekurrente Netze. Dabei handelt es sich um ein Verfahren, das im Speziellen die zeitliche Abfolge der Präsentation von Zeichen – zum Beispiel von Buchstaben in Wörtern oder von Wörtern in Sätzen – berücksichtigt und auf dieser Grundlage Interpretationen erzeugt. In **Kapitel 7** geht es nahtlos mit konvolutionalen Netzen weiter. Zwar hat sich dieses Verfahren vor allem bei der Verarbeitung von Bilddaten bewährt, es lässt sich aber ohne Weiteres auf Sequenzdaten übertragen. Die Vor- und Nachteile gegenüber rekurrenten Netzen sehen wir uns dabei anhand einer automatischen Rechtschreibkorrektur genauer an. Das folgende **Kapitel 8** thematisiert dann die Vorverarbeitung von Wörtern mit Worteinbettungsverfahren. Damit lassen sich Wörter bzw. Token als Vektoren in einem mehrdimensionalen Raum unter Berücksichtigung semantischer und grammatikalischer Beziehungen darstellen. Das Verfahren bildet heute die Grundlage für fast alle Deep Learning-Ansätze, die Textdaten als Rohmaterial verwenden. Wir beleuchten dabei nicht nur die Möglichkeiten, Worteinbettungen selbst anzulernen, sondern führen auch in die Verwendung vortrainierter Vektorräume wie *fastText* oder *Word2Vec* ein.

Beginnend mit **Kapitel 9** eruieren wir dann die Optionen, die komplexe neuronale Netze bei der Textinterpretation und Textgenerierung bieten. Insbesondere geht es um Encoder-Decoder-Modelle. Sie sind aus der modernen Textanalyse (sei es in Sequence-to-Sequence- oder in Transformer-Modellen) nicht mehr wegzudenken. Zunächst dreht sich aber alles um das Handwerkszeug, das wir brauchen, um solche Modelle praktisch umzusetzen: Es geht darum, wie wir Netze, die über zwei Eingänge verfügen oder in denen Datenströme bestimmte

Schichten überspringen, aufsetzen können. Diese Kenntnisse sind notwendig, wenn wir im nachfolgenden **Kapitel 10** Sequence-to-Sequence-Modelle genauer unter die Lupe nehmen. Mit dieser Architektur lassen sich nicht nur einzelne Wörter, sondern ganze Sätze oder Textpassagen erzeugen. Solche Architekturen werden für Übersetzungen oder Textzusammenfassungen eingesetzt. Darüber hinaus lernen wir in diesem Kontext den Attention-Mechanismus kennen, der in einer Variante, der Self-Attention, auch in Transformer-Modellen Karriere gemacht hat. Das ist das Thema von **Kapitel 11**. Darin zeigen wir nicht nur, wie Transformer-Modelle funktionieren, sondern wir sehen uns auch an, wie sich vortrainierte Modelle, die über den Transformer-Hub *Hugging Face* vertrieben werden, sich für verschiedene Aufgaben nachtrainieren lassen. Das Buch schließt im **zwölften Kapitel** mit einer Zusammenfassung einiger grundlegender Konzepte und einer Diskussion der Stärken und Schwächen neuronaler Netze.

Wie Sie mit den Codebeispielen arbeiten können

Die meisten Kapitel bestehen aus einem Mix aus Einführungen in Konzepte maschinellen Lernens und aus einem praxisorientierten Teil, in denen diese Konzepte mit Python umgesetzt werden. Wenn Sie die Beispiele auf Ihrem Rechner laufen lassen möchten, müssen Sie die Voraussetzungen dafür schaffen.

Wir verwenden Python in der *Version 3.7.6*, eingebettet in eine virtuelle Umgebung, die über die Data Science-Plattform *Anaconda* verwaltet wird. *Anaconda* bietet unter anderem den Vorteil, dass es bei der Installation externer Pakete sicherstellt, dass die Pakete untereinander harmonieren. Da wir einige Pakete

installieren müssen und da diese Pakete in ihrer Arbeit wiederum auf weitere Unterpakete angewiesen sind, ist das ziemlich hilfreich.

Um also arbeitsfähig zu sein, benötigen Sie die folgenden Bibliotheken in Ihrer virtuellen Umgebung:

```
matplotlib, Version=3.4.2
nltk, Version=3.6.2
numpy, Version=1.19.2
pandas, Version=1.0.3
scikit-learn, Version=0.22.1
tensorflow, Version=2.3.0
transformers, Version=4.10.2
```

Bis auf *Transformers* (das zum Zeitpunkt der Drucklegung mit pip installiert werden muss) sind alle Bibliotheken über Anaconda verwaltbar (installieren mit *conda install*). Sie können den Code vermutlich auch mit anderen, aktuelleren Versionen von Python und den genannten Bibliotheken ausführen. Allerdings ist es möglich, dass Sie dann an der einen oder anderen Stelle auf Warnungen, Fehlermeldungen oder auf andere Probleme stoßen, die wir nicht vorhersehen können.

Die Codebeispiele im Buch sind über *GitHub* als *Jupyter Notebooks* verfügbar. Sie können den Code entweder im Internet unter https://github.com/tplusone/hanser_deep_nlp abrufen und ansehen oder, wenn Sie die Software *Git* installiert haben, das gesamte Repository inklusive Codebeispielen und Beispieldaten über das Terminal mit dem folgenden Befehl auf Ihren Rechner ziehen:

```
git clone https://github.com/tplusone/hanser\_deep\_nlp.git
```

¹ Vgl. <https://www.wired.com/story/self-driving-cars-challenges/> [Abgerufen am 11.09.2021]

2 Textdaten verarbeiten und vorverarbeiten

Wenn es um die Verarbeitung numerischer und textbasierter Informationen geht, hat Python gegenüber anderen Programmiersprachen ein paar entscheidende Vorteile. Schon die Standardbibliothek bietet eine Vielzahl einfacher Funktionen und Klassen, um Textdaten in Form zu bringen, zu transformieren und zu strukturieren. So ist es ein Leichtes, einmal tokenisierte Texte in Arrays zu verwalten, Auszüge zu extrahieren oder mithilfe von Schleifen oder List Comprehensions Transformationen durchzuführen. Auch Casting-Operationen, die Arrays in Sets, Tuples oder Dictionarys verwandeln, funktionieren in den meisten Fällen vorhersehbar und problemlos. So lassen sich Texte für das Anlernen in Form bringen.

Um einen maschinellen Lernalgorithmus zu trainieren, brauchen wir aber andere Datenklassen. Statistische Verfahren operieren mit mathematischen Funktionen und mögen daher keine Überraschungen, wenn es um Datentypen und die Struktur der Datenbehälter geht. Da eines der Grundprinzipien in Python aber die Abkehr von Typsicherheit zugunsten von Duck-Typing ist, benötigen wir eine Alternative, die dieses Anforderungsprofil

erfüllt. Die Bibliothek *NumPy* (Numerical Python) ist dabei die erste Wahl. NumPy ist zum Glück auf die Datenklassen der Standardbibliothek abgestimmt, sodass Castings in beide Richtungen reibungslos funktionieren.

Im Folgenden sehen wir uns einige ausgewählte Techniken, Klassen und Bibliotheken an, die für die Vorverarbeitung von Texten zur Analyse mit Lernalgorithmen von Bedeutung sind. Dabei handelt es sich um keine auch nur annähernd vollständige Abhandlung der verschiedenen Möglichkeiten. Wir werfen aber einen Blick auf die Verfahren, die in den nachfolgenden Kapiteln immer wieder verwendet werden und die wir dort nicht noch einmal im Detail vorstellen. Neben der Tokenisierung und numerischen Encodierung von Wörtern geht es um die Repräsentation von Wörtern als One-Hot-Sets und natürlich um die Arbeit mit der Bibliothek NumPy.

2.1 Grundlegende Techniken der Verarbeitung von Textdaten

Texte liegen nach dem Einlesen in Python normalerweise als Strings vor. Eine Besonderheit der String-Klasse in Python ist, dass sie sich wie ein Iterable verhält. Die einzelnen Buchstaben werden als Characters auf Indexpositionen abgespeichert. Man kann deshalb über einen String sowohl iterieren als auch über den Aufruf einer Indexposition bzw. über Slicing einzelne oder mehrere Buchstaben herauslösen:

```
1. text = 'Die Sonne steht hoch am Himmel.'  
2. text[5], text[5:10]
```

Ausgabe:

```
('o', 'onne ')
```

Da in vielen statistischen Anwendungen Wörter oder Token als kleinste Analyseeinheiten fungieren, müssen wir Texte fast immer auf dieser Ebene zerlegen. Dieser Vorgang nennt sich *Tokenisierung*. Wir könnten uns mit einem Regex zwar einen eigenen Tokenizer zusammenbauen, einfacher geht es allerdings mit einem getesteten Produkt, das Wörter und Satzzeichen an verschiedenen Positionen erkennt und extrahiert. Die `word_tokenize`-Funktion aus dem NLTK-Modul erledigt genau diese Arbeit und gibt bei Übergabe eines Strings eine Liste der enthaltenen Wörter inklusive Satzzeichen zurück:

```
3. from nltk.tokenize import word_tokenize  
4.  
5. text = 'Die Sonne steht hoch am Himmel.'  
6. word_tokenize(text)
```

Ausgabe:

```
['Die', 'Sonne', 'steht', 'hoch', 'am', 'Himmel', '.']
```

In längeren Texten wiederholen sich die Wörter in der Regel mehrfach. Wenn wir von jedem dieser Wörter jeweils nur ein Exemplar behalten möchten, können wir die Liste in ein Set umwandeln. Eine angenehme Nebenwirkung dieses Vorgangs ist, dass automatisch alle Duplikate eliminiert werden:

```
7. words = ['die', 'Sonne', 'Sonne', 'scheint', 'schein  
t']  
8. set(words)
```

Ausgabe:

```
{'Sonne', 'die', 'scheint'}
```

Wie man Wörter in numerische Form bringt, um damit einen Machine Learning-Algorithmus zu füttern, sehen wir uns im nächsten Abschnitt genauer an. Manchmal brauchen wir allerdings keine spezielle Encodierung, sondern lediglich eine numerische Repräsentation der Wörter aus den Trainingsdaten (zum Beispiel, wenn wir eine Zielvariable vorbereiten). In diesem Fall bietet es sich an, zur Encodierung ein Dictionary zu verwenden. Es enthält für jedes Wort (*Key*) eine Ganzzahl (*Value*). Damit können wir die einzelnen Wörter aus einer Textsequenz nachschlagen und encodieren. In einem zweiten (zur Decodierung vorgesehenen) Dictionary sind die Keys und Values vertauscht: Bei Übergabe einer Ganzzahl erhalten wir das zugeordnete Wort zurück.

Die Ordnung der Wörter und Ganzzahlen in den Dictionarys können wir letztlich willkürlich festlegen, da es für den Lernalgorithmus unerheblich ist, welche Ganzzahlen für welche Wörter stehen. Wichtig ist nur, dass jedem Wort je eine eigene Ganzzahl zugeordnet ist. Allerdings bietet es sich natürlich an, die Reihenfolge der Zahlen mit der alphabetischen Ordnung der Wörter in Einklang zu bringen:

```
9. words = ['die', 'Sonne', 'Sonne', 'scheint', 'schein  
t']  
10. words = list(set(words))  
11. words = [ word.lower() for word in words ]  
12. words.sort()  
13. word_index = dict([ (word, idx) for idx, word in enumerate(words) ])  
14. index_word = dict([(idx, word) for word, idx in word_index.items()])  
15. word_index, index_word
```

Ausgabe:

```
({'die': 0, 'scheint': 1, 'sonne': 2},  
 {0: 'die', 1: 'scheint', 2: 'sonne'})
```