Praktikumsbericht

Machine Learning for Natural Language Processing

am Lehrstuhl für Informatik X

Konstantin Herud Thomas Schaffroth Maximilian Meißner

Abgabedatum: 16. August 2020 Betreuer: Daniel Schlör

Albin Zehe Konstantin Kobs Tobias Koopmann



Julius-Maximilians-Universität Würzburg Lehrstuhl für Informatik X Data Science

Zusammenfassung

Dieses Dokument soll Studenten an unserem Lehrstuhl bei der Erstellung ihrer Abschlussarbeit unterstützen. Wir zeigen eine beispielhafte Gliederung einer Arbeit und beschreiben die Inhalten der einzelnen Kapitel. Zusätzlich geben wir an vielen Stellen auch Hinweise zur Benutzung von LATEX für die Erstellung der Arbeit. Im Anhang ?? geben wir ein paar Hinweise zum Ablauf der Betreuung von Abschlussarbeiten an unserem Lehrstuhl.

Zur Handhabung dieses Pakets. In diesem Paket sind Vorlagen für verschiedene Dokumenttypen enthalten, die sie als Ausgangspunkt für ihre Arbeit verwenden können. Es gibt jeweils Vorlagen für deutsche und englische Arbeiten.

- template_thesis_de.tex, template_thesis_en.tex: Vorlage für Bachelorarbeit bzw. Masterarbeit
- template_seminar_de.tex, template_seminar_en.tex: Vorlage für Seminarausarbeitungen und Praktikumsberichte

Der Quelltext zu diesem Leitfaden ist ebenfalls im Paket enthalten. Diesen können Sie als praktisches Beispiel dafür verwenden, wie diese Dokumentenklasse angewandt wird.

Inhalt der Zusammenfassung. Schreiben Sie hier eine Zusammenfassung der Arbeit, vergleichbar mit dem Abstract auf wissenschaftlichen Papers. Sie dient dem Leser dazu, einen groben Überblick über die Inhalte zu gewinnen (Problemstellung, verwendeter Lösungsansatz, ggf. experimentelle Ergebnisse, gewonnene Erkentnisse). Der Umfang soll ca. eine halbe Seite betragen. Für Seminararbeiten ist diese Zusammenfassung nicht erforderlich.

Achtung: Bei Arbeiten auf Englisch fordern die Prüfungsordnungen, dass es eine deutsche Zusammenfassung gibt. Schreiben Sie in diesem Fall eine englische und eine deutsche Zusammenfassung (mit dem gleichen Inhalt). Die passenden IATEX-Befehle dafür finden Sie in den englischsprachigen Vorlagen.

WARNUNG: Die vorliegende Version des Leitfadens ist eine Vorabversion, die noch nicht vollständig ist. Sie bezieht sich größtenteils auf die Ausarbeitung von Bachelorund Masterarbeiten; Seminararbeiten unterscheiden sich davon etwas in Aufbau und Inhalt.

Inhaltsverzeichnis

1.	Einleitung	4
2.	Methodik	6
	2.1. Task 1: Dataset Preparation	6
	2.2. Task 2: Learning to Discriminate	9
	2.2.1. Stylometrische Klassifikation	9
Lit	teraturverzeichnis	12
Α.	Formelles/LaTeX	13

1. Einleitung

Thomas Schaffroth

Nachdem der in diesen Tagen gängige Begriff der Artificial Intelligence (AI) auf der Dartmouth Conference 1956 erstmals geprägt wurde, bezeichnet dieser bis heute das entsprechende Forschungs- und Entwicklungsfeld. Was damals mit dem computergestützten Lösen algebraischer Probleme oder der Demonstration geometrischer Theoreme begann, brachte im Laufe der Jahre einen umfangreichen, sich schnell entwickelndem Bereich der Informatik hervor. Nicht zuletzt stellten Innovationen, wie die Bayesschen Lernverfahren 1960, die Einführung des "Parallel Computings" 1980 oder die Erfindung des Backpropagation-Algorithmus für Neuronale Netze 1984, wesentliche Meilensteine in dieser Entwicklung dar.

Durch aufkommende Neuerungen zu Beginn der 2000er Jahre, wie der Herausforderung des Big Data sowie der Entfaltung des Internets und der mobilen Kommunikation wurden seitdem große Fortschritte in Feldern, wie der Computer Vision, Intelligenten Agenten und Mustererkennung erzielt. Es sind unter anderem ebendiese Weiterentwicklungen, die den Weg für moderne Herausforderungen, wie Spracherkennung, selbstfahrende Fahrzeuge und Natural Language Processing (NLP), geebnet haben [4]. Besonders das Themengebiet des NLP nimmt einen wichtigen Stellenwert in vielen Software-Anwendungen unseres täglichen Lebens ein. Einige der herausragendsten Beispiele sind E-Mail Plattformen, wie Microsoft Outlook (Spam-Klassifikation, Auto-Complete, etc.), sprachbasierte Assistenzsysteme, wie Apple Siri und Amazon Alexa oder Services für maschinelle Übersetzung, wie Google Translate [5].

Die zunehmende Bedeutung von NLP lässt sich neben seiner Rolle in unserem Alltag auch in der Forschung beobachten. Die ACL Anthology (AA) ist ein digitales Repository, dass zehntausende Veröffentlichungen von NLP-Papern aus der Familie der ACL- sowie anderer NLP-Konferenzen beinhaltet. Lag die Anzahl der Veröffentlichungen im Jahr 2000 noch bei 1050 wurden 2018 bereits 4173 publizierte Paper verzeichnet [3].

In diesem Praktikum wird im Folgenden das Natural Language Processing an sich mit wissenschaftlichen Arbeiten, die zu diesem Thema veröffentlicht wurden, zusammengeführt. Ziel dieser Arbeit ist die Bearbeitung der Fragestellung, ob die Erkennung der Herkunft von Autoren wissenschaftlicher Publikationen aufgrund textueller Eigenschaften ihrer Veröffentlichungen möglich ist, sowie welche Methoden bessere oder schlechtere Ergebnisse liefern. Der hierfür zu Verfügung stehende Datensatz beinhaltet 18000 Paper verschiedener NLP- und AI-Konferenzen in PDF- und Text-Form.

Das Praktikum gliedert sich in zwei Bestandteile. Zunächst müssen die Paper einer grundsätzlichen Datenbereinigung unterzogen werden. Informationen, von denen direkt oder indirekt auf die Herkunft der Autoren geschlossen werden kann, müssen dabei aus den Arbeiten entfernt werden. Beispiele hierfür sind Autoren-Namen, E-Mail-Adressen,

Ländernamen, Institutionen oder Referenzen im Text. Der Grund dafür ist, das für den zweiten Teil der Arbeit sichergestellt werden muss, dass die Klassifikation der Paper aufgrund allgemeiner Texteigenschaften, wie der Syntax und dem verwendeten Vokabular stattfindet.

Informationen, die sich auf die Herkunft des Autors beziehen, würden für eine Verfälschung des Klassifikationsergebnisses sorgen. Der folgende Schritt beschäftigt sich mit der Feature-Modellierung, dem Training und der Evaluierung im Kontext verschiedener Techniken der AI und NLP auf dem genannten Korpus in Bezug auf die Klassifikations-Aufgabe. Ziel soll sein, einen Klassifikator zu trainieren, der für eine möglichst große Zahl der Paper das Herkunftsland des jeweiligen Autors korrekt vorhersagt.

2. Methodik

Konstantin Herud

2.1. Task 1: Dataset Preparation

Da verschiedene Informationen das Herkunftsland von Autoren entweder direkt oder indirekt preisgeben, müssen diese aus der Datenbasis entfernt werden. Die Zielsetzung sieht hierbei vor Namen der Autoren, E-Mails, Institutionen und Firmen, Herkunftsländer, Förderungen und Danksagungen, Persönliche Daten sowie Referenzen zu entfernen. Referenzen sind beispielsweise indirekte Hinweise, weil Arbeitsgruppen der gleichen Herkunft und somit dem selben Fachgebiet die selbe Literatur referenzieren.

Da diese Aufgabe essentieller Bestandteil der Validität späterer Ergebnisse ist, wurde eine mehrstufige Daten-Pipeline entwickelt, um alle Informationen bestmöglich aus den Texten zu entfernen.

Textextraktion Im ersten Schritt muss der rohe Text der Dokumente eingelesen werden. Die Daten liegen sowohl in PDF- als auch in Textform vor, wobei letztere durch das Befehlszeilenprogramm pdftotext erstellt wurde und somit bereits Potenzial für Fehler bietet. Ursprünglich wurde deshalb ein Ansatz mit dem Werkzeug pdfminer verfolgt, um den Text der Dokumente selbstständig zu extrahieren und dabei direkt Strukturinformationen ausnutzen zu können (z. B. um Referenz-Blöcke zu erkennen). Bei einigen Dokumenten entstanden jedoch Fehler mit pdfminer, weshalb dieser Ansatz verworfen wurde. So werden in der Pipeline zunächst die Textdateien sowie zugehörige Metadaten eingelesen ¹. Da verschiedene spätere Schritte außerdem Informationen über die Zeichenposition des Abstract-Starts brauchen, wird dieser zudem bereits mit einem regulären Ausdruck, unterstützt durch verschiedene Fehlermaßnahmen, bestimmt. So wird beispielsweise überprüft, ob der Titel des Dokuments das Wort Abstract enthält, um zu frühe Treffer des regulären Ausdrucks zu überspringen. Falls der Start nicht gefunden werden kann oder eine bestimmte Obergrenze überschreitet, wird ein mittelwert-basierter Defaultwert verwendet.

Dokumentkopf Als Dokumentkopf definieren wir jeglichen Text (in Rohform), der vor dem Abstract erscheint. Die hohe Dichte unterschiedlicher zu entfernender Informationen, gepaart mit Artefakten und Strukturfehlern, machen diese Stufe zu einer der größten Herausforderungen der Aufgabe. So kommen beziehen sich vor allem Namen, E-Mails,

¹https://github.com/marekrei/ml_nlp_paper_data

Herkunftsländer, Institutionen und persönliche Daten auf diesen Schritt. Initial wurde dafür die Idee verfolgt, Named-Entity-Recognition mittels spaCy [2] einzusetzen, um all diese Informationen (außer E-Mails, welche leicht durch reguläre Ausdrücke zu erkennen sind) zu entfernen. SpaCy bietet dafür verschiedene vortrainierte Modelle, welche für diese Arbeit mit im Text durch reguläre Ausdrücke erkannten Informationen nachtrainiert wurden. Da dieser Ansatz Entitäten nicht nur zu unverlässlich erkannte (insbesondere keine Phrasen, die über mehrere Tokens reichen), sondern auch zu schlecht klassifzierte wurde ein weiterer Ansatz entwickelt. So wurde ein LSTM-basiertes [1], neuronales Netzwerk implementiert, um jede Zeile (im Textdokument, also keine ganzen Sätze) des Vor-Abstract-Texts in die Klassen Autor, E-Mail, Private Informationen, Institution oder Anderes zu unterteilen. Länder werden in einem späteren Schritt entfernt. Zum Training wurden 250 Dokumente zeilenweise mit einer Kommandozeilenanwendung annotiert. Dabei wurden zufällig Dokumente aus dem Korpus gezogen, allerdings wurde sichergestellt, dass jede Konferenz vorkommt. Zudem sind verschiedene Fehlerfälle enthalten, bei denen der Start des Abstracts nicht korrekt erkannt wurde und somit späterer Text enthalten ist. Das Modell erreicht auf weiteren 50 Dokumenten einen Makro-F1-Wert von etwa 96%. Die Architektur verfolgt einen ähnlichen Ansatz wie Gleichung 2.4–2.7, nutzt jedoch nach X_{zd} ein zweites LSTM zur zeilenweisen Klassifikation.

Referenzen Referenzen unterteilen sich in zwei Typen: Verweise innerhalb des Texts und ausführlich Referenzblöcke. Erstere sind aufgrund ihrer einheitlichen Struktur leicht mittels regulären Ausdrücken zu ermitteln. Dafür wurden zahlreiche zufällige Dokumente betrachtet und drei Strukturtypen identifiziert, für die eigene reguläre Ausdrücke entwickelt wurden: Klammer-Ausdrücke, die Jahreszahlen im 20. und 21. Jahrhundert enthalten, Phrasen mit "et. al." und Aufzählungen (z.B. verbunden durch "und", "," oder "&" etc.), gefolgt von Jahreszahlen. Referenzblöcke zu identifizieren stellt sich jedoch als zweitschwerste Aufgabe heraus. Ein initialer Ansatz sah vor, die zuvor im Text gefundenen Referenzen zusammen mit verschiedenen Schlüsselwörtern dazu zu nutzen, Zeilen-Blöcke als Referenzen zu erkennen. Da durch die Konvertierung von zweispaltigen, mehrseitigen Dokumenten in eindimensionalen Text jedoch häufig Strukturfehler auftreten, sind Teile der Referenzblöcke oftmals stark über die Dokumente verteilt, weshalb dieser initiale Ansatz zu keinen befriedigenden Ergebnissen führte. Deshalb wurde ebenfalls ein neuronales Netz, mit einer Architektur äquivalent zur vorherigen, trainiert, um zeilenweise zu annotieren, ob es sich um Referenzen handelt. Dafür wurden etwa 50 Dokumente manuell annotiert, wobei neben zufällig gewählten Werken insbesondere darauf geachtet wurde, Fehlerfälle zu verwenden. Das Netz verarbeitet den Text iterativ in Stücken von 80 Zeilen (um Padding und somit Grafikspeicher-Anforderungen möglichst gering zu halten, wurde kein kompletter End-to-End Ansatz gewählt). Um dennoch von Informationen über die Position dieser 80 Zeilen im Dokument zu profitieren, wird auf die d-dimensionalen (1, ..., i, ..., d), internen Zustände des Netzwerks ein Positionssignal PE addiert (vgl. [6], Gleichung 2.1 & 2.2).

$$PE(zeile, 2i) = \sin\left(\frac{zeile}{10000^{2i/d}}\right)$$
 (2.1)

$$PE(zeile, 2i+1) = \cos\left(\frac{zeile}{10000^{2i/d}}\right)$$
 (2.2)

Das Modell erreicht auf 20% separaten, zufälligen Daten einen Makro-F1-Wert von beinahe 99,7%. Sowohl das Dokumentkopf- als auch das Referenzblocknetz sind mit ihren Gewichten aufgrund der hoch-rekurrenten Struktur unter einem Megabyte groß.

Gutachter Generell wurden wenige Gutachter innerhalb der Dokumente festgestellt. Diese geringe Menge folgt jedoch einem festen Schema, welches mittels eines regulären Ausdrucks ermittelt wird.

Danksagungen Danksagungen sind meist wenige Zeilen in geschlossenen Blöcken und somit leichter zu ermitteln als Referenzblöcke. Vorerst stand im Raum, ebenfalls ein neuronales Netz einzusetzen, allerdings genügen reguläre Ausdrücke. So wird nach dem letzten Vorkommen der gängigen Überschrift "Acknowledgment" gesucht (wobei unterschiedliche Schreibweisen zu beachten sind, i.e. Acknowledge?ments?). Das Ende der Danksagung wird durch eine Liste mit Schlüsselwörtern, kombiniert mit den zuvor gefunden Referenzen sowie einem regulären Ausdruck für Referenzen (da Danksagungen diesen häufig vorausgehen), ermittelt. Kann dennoch kein Ende gefunden werden, wird der Text bis zum Ende des Paragraphen gewählt.

E-Mails Der Großteil aller E-Mails wurde in diesem Schritt bereits durch die Kopfzielen-Stufe entfernt (insbesondere schwer identifizierbare Fälle, deren gezielt irreführende Schreibweise z. B. Spam verhindern soll). Der Rest aller E-Mails wird in diesem Schritt mit zwei weiteren regulären Ausdrücken entfernt.

Herkunftsländer Länder und Nationalitäten werden mit dem Python-Modul *geotext* ² ermittelt. Das Modul arbeitet mit regulären Ausdrücken und der geographischen Datenbank *GeoNames*, die jegliche Länder und über elf Millionen Ortsnamen enthält.

Fußnoten Fußnoten sind unscharf definiert und aufgrund der fehlenden visuellen Struktur sehr schwer im rohen Text zu identifizieren. Um diese Informationen dennoch zu entfernen, wird der Text zunächst mittels spaCy in Sätze unterteilt. Anschließend werden alle Sätze mit einer Liste von Schlüsselwörtern, wie "sponsored" oder "internship", sowie bereits ermittelten Autor- und Organisationsnamen untersucht. Wird ein Schlüsselwort gefunden, wird der entsprechende Satz entfernt.

²https://github.com/elyase/geotext

Sprache Zuletzt wird nicht-englische Sprache entfernt. Diese Aufgabe stellt ebenfalls eine Herausforderung dar, da viele Rohtext-Fragmente aufgrund kryptischer Formeln oder anderen Artefakten von gängigen Sprachidentifikationssystemen fälschlicherweise als nicht-englisch erkannt werden. Dennoch wird die spaCy-Erweiterung *CLD* eingesetzt, welche die *Compact Language Detection 2* Bibliothek ³ anbindet. Diese ermöglicht die Klassifikation von rohem Text in über 80 verschiedene Sprachen mittels Zeichen-N-Grammen und einem Naiven-Bayes-Modell. Dieser Ansatz führt zu einer akzeptablen Sensitivität, allerdings zu einer schlechten Genauigkeit. Um letztere zu verbessern, wird lediglich eine geringe Teilmenge relevanter Sprachen der 80+ möglichen berücksichtigt.

Die einzelnen Stufen der Pipeline lassen sich nicht nur jeweils über verschiedene Prozesse verteilen und vervielfältigen, die Pipeline lässt sich auch als Ganzes parallelisieren. Zusätzlich ermöglicht die Anwendung Hardware-Beschleunigung für die neuronalen Netzwerke. Das Wissen über zu entfernende Informationen wird akkumuliert und letztlich in eine CSV-Datei geschrieben, außerdem wird der inkrementell bereinigte sowie der ursprüngliche Text zwischen Stufen weitergegeben und abschließend in eine Text-Datei geschrieben.

2.2. Task 2: Learning to Discriminate

Um die bereinigten Dokumente hinsichtlich des Herkunftslands ihrer Autoren zu klassifizieren, wurden zwei Ansätze verfolgt: Zum einen Klassifikation mittels verschiedenen stylometrischen Merkmalen, zum anderen ein *End-to-End*-Deep-Learning-Modell.

2.2.1. Stylometrische Klassifikation

- z: Zeilen im Dokument
- t: Tokens pro Zeile (Padding kürzerer Zeilen)
- d: Dimension des Modells (z. B. 150)
- e: Embedding-Dimension
- c: Anzahl Klassen

³https://github.com/CLD2Owners/cld2

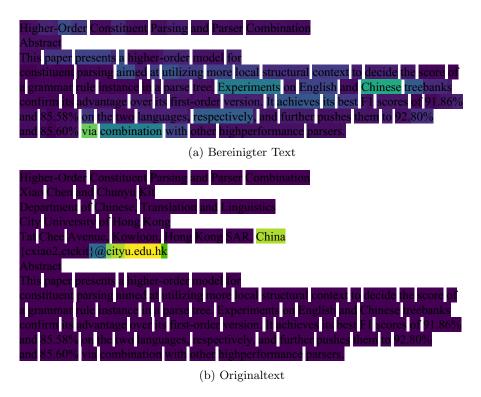


Abb. 2.1.: Visualisierung der gelernten Gewichtung einzelner Subtokens zur Klassifikation des Herkunftslandes mit bereinigtem und ursprünglichem Text.

$$X_{zte} = \text{Embedding}(X_{zt})$$
 (2.3)

$$X_{ztd} = LSTM_t(X_{zte})$$
 (2.4)

$$A_{zt} = \operatorname{softmax}_t \left(\sum_d X_{ztd} W_d^1 \right) \tag{2.5}$$

$$X_{zd} = \sum_{t} X_{ztd} A_{zt} \tag{2.6}$$

$$A_z = \operatorname{softmax}_z \left(\sum_d X_{zd} W_d^2 \right) \tag{2.7}$$

$$X_d = \sum_z X_{zd} A_z \tag{2.8}$$

$$Y_c = Y_d W_{dc}^3 \tag{2.9}$$

$$W_{L,C_i} = \left(\frac{|C_{\text{max}}|}{|C_i|}\right)^{\alpha}, \quad \alpha = 0.6$$
(2.10)

		AU	CA	CN	FR	DE	IN	IL	JP	SG	СН	UK	USA
Wahrheit	AU	6		8		1						7	6
	CA	1	8	4					2		1	7	32
	CN	1	1	188		1			2	1		1	17
	FR			1	36	3			1			2	4
	DE		2	2	2	47			1			12	11
	IN			1			6					2	11
	IL		1		1			14	2				10
	$_{ m JP}$			2	1				45		1	1	6
	SG		1	6						9		2	8
	CH			3		7					4	4	8
	UK	3	5	1	6	1	2		1		2	66	18
	USA	3	13	38	17	6	1	2	8	2	3	32	746

Tab. 2.1.: Konfusionsmatrix der Klassen Australien (AU), Canada (CA), China (CN), Frankreich(FR), Deutschland (DE), Indien (IN), Israel (IL), Japan (JP), Singapur (SG), Schweitz (CH), Großbritannien (UK), Amerika (USA). Leere Zellen implizieren Nullen.

Literaturverzeichnis

- [1] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural computation*, 9:1735–80, 12 1997.
- [2] M. Honnibal and I. Montani. spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing. To appear, 2017.
- [3] S. M. Mohammad. The state of nlp literature: A diachronic analysis of the acl anthology, 2019.
- [4] J. A. Perez, F. Deligianni, D. Ravi, and G.-Z. Yang. Artificial intelligence and robotics, 2018.
- [5] S. Vajjala, B. Majumder, A. Gupta, and H. Surana. *Practical Natural Language Processing*. O'Reilly Media Inc., 2020.
- [6] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. u. Kaiser, and I. Polosukhin. Attention is all you need. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems* 30, pages 5998–6008. Curran Associates, Inc., 2017.

A. Formelles/LaTeX