

LUDWIG-MAXIMILIANS-UNIVERSITÄT MÜNCHEN  
PROGRAMMING LANGUAGES AND ARTIFICIAL INTELLIGENCE



## Titel der Arbeit

Ruben Triwari

Bachelorarbeit  
im Studiengang 'Informatik plus Mathematik'

Betreuer: Prof. Dr. Johannes Kinder

Mentor: Moritz Dannehl, M.Sc.

Ablieferungstermin: 28. August 2024

# Inhaltsverzeichnis

<b>1</b>	<b>Einführung</b>	<b>1</b>
<b>2</b>	<b>Grundlagen und Termini</b>	<b>2</b>
2.1	Maschinelles Lernen . . . . .	2
2.1.1	Definition . . . . .	2
2.2	Sentence Transformer . . . . .	2
2.3	Large Language Models . . . . .	2
2.4	Code2Vec . . . . .	2
<b>3</b>	<b>Methodik</b>	<b>2</b>
3.1	Datensatz . . . . .	2
3.2	Datenpipeline . . . . .	3
3.3	Stabilität von SentenceTransformer . . . . .	3
<b>4</b>	<b>Funktionskommentare</b>	<b>5</b>
4.1	Motivation . . . . .	5
4.2	Methodik . . . . .	5
<b>5</b>	<b>Code2Vec</b>	<b>5</b>
5.1	Motivation . . . . .	5
5.2	Adaption auf C . . . . .	5
5.3	Training . . . . .	5
<b>6</b>	<b>Funktionsnamen</b>	<b>5</b>
6.1	Motivation . . . . .	5
6.2	Methodik . . . . .	5
<b>7</b>	<b>Coddelama-Erklärungen</b>	<b>5</b>
7.1	Motivation . . . . .	5
7.2	Codellama . . . . .	5
7.3	Prompt Engineering und Temperature . . . . .	5
<b>8</b>	<b>Ergebnisse</b>	<b>5</b>
8.1	Evaluierung durch Experten . . . . .	5
8.1.1	Methodik . . . . .	5
8.1.2	Auswertung und Ergebnisse . . . . .	5
8.2	Qualitative Evaluierung . . . . .	5
8.3	Quantitative Evaluierung . . . . .	5
<b>9</b>	<b>Limitation</b>	<b>5</b>
<b>10</b>	<b>Diskussion</b>	<b>5</b>
<b>11</b>	<b>Fazit</b>	<b>5</b>

<b>12 Results: Comparing natural language supervised methods for creating Rich Binary Labels</b>	<b>5</b>
<b>13 Conclusion</b>	<b>6</b>
<b>14 Notes on form</b>	<b>6</b>
14.1 Formatting . . . . .	6
14.2 Citation . . . . .	7
<b>15 General Addenda</b>	<b>9</b>
15.1 Detailed Addition . . . . .	9
<b>16 Figures</b>	<b>9</b>
16.1 Example 1 . . . . .	9
16.2 Example 2 . . . . .	9
<b>Literatur</b>	<b>12</b>

# Abstract

## 1 Einführung

In den letzten Jahren gab es große Fortschritte in der natürlichen Sprachverarbeitung, besonders hervorzuheben sind Large Language Models die sich mittlerweile in vielen Bereichen der Informatik in die Lösungsansätze für Problemen in jeweiligen Bereichen eingeschlichen haben. Diese Arbeit untersucht nun, ob diese Fortschritte in der natürlichen Sprachverarbeitung eine Hilfestellung leisten können um Source Code Funktionen semantisch sinnvoll in einen Vektor mit reelwertigen Zahlen zu codieren. Diese Vektoren können dann später als Label verwendet werden um ein Modell zu trainieren was Binary Code als Input nimmt und diesen ebenfalls in einen semantischen Vektor mit reelwertigen Zahlen codiert. Das resultierende Modell kann hinterher verwendet werden um Reverse Engeneering zu erleichtern. Ein einfaches Beispiel ist folgendes: Man stelle sich vor, dass man eine Funktion die in Binary Code vorliegt, mühselig manuell verstanden was für eine Aufgabe die Funktion in der Code Base hat. Nun kann man diese Funktion codieren und über die Gesamte Code Base einen Nearest Neighbor Search durchführen und all ähnlichen Funktionen ausgeben lassen. Das spart zeit, denn nun hat man eine Idee was diese anderen Funktionen für eine Aufgabe in der Code Base erfüllen könnten.

Das oben beschriebene Problem Source Code Vektoren in sinnvoll semantische reelwertige Vektoren zu codieren ist sehr ähnlich zu einen Problem in der natürlichen Sprachverarbeitung und dort bereits gelöst. Die rede ist von dem Problem einen gegebenen Satz in einen semantisch sinnvollen Vektor abzubilden. Es ist nageliegend zu versuchen dieses Ergebnis der natürlichen Sprachverarbeitung zu benutzen um eine Lösung für unser Problem zu konstruieren. Die intuitivste Idee ist es einfach die Funktionsnamen, die in natürlicher Sprache verfasst sind als beschreibung der Funktion zu verwenden. Diese Beschreibunf können wir nun mühelos codieren, da sie in natürlicher Sprache vorliegt. Eine zweite Idee ist, die Kommentare der Funktionen, die in natürlicher Sprache verfasst sind, als Beschreibung der Funktion zu verwenden. Am viel versprechsten ist es die Funktionen von einen Large Language Modell in natürlicher Sprache beschreiben zu lassen. Als letztes habe ich noch ein bestehendes Modell Code2Vec verwendet und es für dieses Problem angepasst.

## 2 Grundlagen und Termini

### 2.1 Maschinelles Lernen

In diesen Abschnitt wird zunächst maschinelles lernen definiert und dann darauf aufbauend grundlegende Trainingsarten vorgestellt. Heutzutage ist maschinelles Lernen weitverbreitet und wird nahezu in jeden Bereich der Informatik verwendet. Maschinelles Lernen wird überall eingesetzt wo eine analytische Lösung eines Problems zu aufwendig oder gar überhaupt nicht existiert.

#### 2.1.1 Definition

Diese Lösung durch maschinelles Lernen versucht aus den Daten ein Muster abzuleiten. Bei einer endlichen Menge an Daten ist klar, dass das resultierende Modell nur eine Approximation der gesuchten Lösung ist. Trotz des Bekanntheitsgrades, gibt es den Irrglauben, dass maschinelles lernen nur was mit Neuronale Netze zu tun hat, diese Annahme ist im allgemeinen falsch. Im folgenden wird eine allgemeine Definition von maschinellen Lernen vorgestellt:

**Definition 1** Sei  $X$  eine beliebige Input Menge,  $Y$  eine beliebige Output Menge,  $f \in \{X \rightarrow Y\}$  die gesuchte Lösung des Problems,  $\mathbb{D}$  eine beliebige Menge aus gegebenen Datenpunkten,  $H_1 \subset \{X \rightarrow Y\}$  ein Hypothesenraum, und  $A_1 : \mathcal{P}(\{X \rightarrow Y\}) \times \mathcal{P}(\mathbb{D}) \rightarrow \{X \rightarrow Y\}$  ein Lernalgorithmus. Dann ist das Ziel, bei gegebenen Daten, den Hypothesenraum  $H_1$  und den Lernalgorithmus  $A_1$  so zu wählen, sodass

$$A_1(H_1, \mathbb{D}) \approx f.$$

Maschinelles lernen ist also die Suche nach einem Lernalgorithmus und Hypothesenraum, die dann in Kombination mit gegebenen Daten, die optimale Lösung approximieren. Dabei ist hervorzuheben, dass der Datensatz das Herzstück jeder Problemstellung im Bereich des maschinellen Lernens ist. Ist der Datensatz zu klein oder überhaupt nicht repräsentativ für das gegebene Problem, wird der Lernalgorithmus die falschen Muster erkennen und dadurch eine fehlerhafte Approximation produzieren.

### 2.2 Sentence Transformer

### 2.3 Large Language Models

### 2.4 Code2Vec

## 3 Methodik

### 3.1 Datensatz

Im Maschinellen lernen hat der Datensatz bzw. die Trainingsdaten den größten Einfluss auf die Güte des Modells.

## **3.2 Datenpipeline**

## **3.3 Stabilität von SentenceTransformer**



## **4 Funktionskommentare**

### **4.1 Motivation**

### **4.2 Methodik**

## **5 Code2Vec**

### **5.1 Motivation**

### **5.2 Adaption auf C**

### **5.3 Training**

## **6 Funktionsnamen**

### **6.1 Motivation**

### **6.2 Methodik**

## **7 Coddelama-Erklärungen**

### **7.1 Motivation**

### **7.2 Codellama**

### **7.3 Prompt Engeneering und Temperature**

## **8 Ergebnisse**

### **8.1 Evaluierung durch Experten**

#### **8.1.1 Methodik**

#### **8.1.2 Auswertung und Ergebnisse**

### **8.2 Qualitative Evaluierung**

### **8.3 Quantitative Evaluierung**

## **9 Limitation**

## **10 Diskussion**

## **11 Fazit**

## **12 Results: Comparing natural language supervised methods for creating Rich Binary Labels**

- Stabilität von Sentence Transformer



- Kommentare von Funktionen um Embeddings zu generieren
- Funktionsnamen von Funktionen um Embeddings zu generieren
- Code2Vec um Embeddings zu generieren
- CodeLlama Erklärungen von Funktionen um Embeddings zu generieren
- Evaluierung durch tSNE-Plots
- Evaluierung durch Experten
- Evaluierung durch Formel

$$I_k : \mathbf{N} \times \mathbf{N} \times \mathbf{N}^k \rightarrow [0, 1]$$

$$I_k(x, i, v) = \begin{cases} 1 & , \exists j \in \mathbf{N} : x = v_j \wedge i = j \\ \frac{1}{2} & , \exists j \in \mathbf{N} : x = v_j \wedge i \neq j \\ 0 & , \text{otherwise} \end{cases}$$

$$E_k : \mathbf{N}^k \times \mathbf{N}^k \rightarrow [0, 1]$$

$$E_k(u, v) = \frac{1}{G_k} \sum_{i=1}^k \frac{I_k(u_i, i, v_i)}{\log_2(i+1)}$$

$$\text{wo } G_k := \sum_{i=1}^k \frac{1}{\log_2(i+1)}.$$

$$CMP_k : \mathbf{R}^{N \times l} \times \mathbf{R}^{N \times l} \times \{\mathbf{R}^l \times \mathbf{R}^{N \times l} \rightarrow \mathbf{N}^k\} \times \{\mathcal{P}([0, 1]) \rightarrow [0, 1]\} \rightarrow [0, 1]$$

$$CMP_k(X, Y, f_k, agg) = agg(\{E_k(f_k(X_{i,j}, X), f_k(Y_{i,j}, Y)) | j \in \{1, 2, 3, \dots, N\}\})$$

## 13 Conclusion

## 14 Notes on form

### 14.1 Formatting

This LaTeX template uses the following formatting:

- font: Linux Libertine O (alternatively: Times New Roman)
- font size: 12 pt
- left and right margin: 3.5 cm, top and bottom margin: 3 cm
- align: left
- line spacing: one and a half (alternative: 15 pt line spacing with 12 pt font size)

When implementing the specifications in Word, it is essential to define style sheets.

## 14.2 Citation

The citation method follows the author-year system. Place reference is in the text, footnotes should only be used for explanations and comments. The following notes are taken from the *language* bibliography template from `ron.artstein.org`:

The *Language* style sheet makes a distinction between two kinds of in-text citations: citing a work and citing an author.

- Citing a work:
  - Two authors are joined by an ampersand (&).
  - More than two authors are abbreviated with *et al.*
  - No parentheses are placed around the year (though parentheses may contain the whole citation).
- Citing an author:
  - Two authors are joined by *and*.
  - More than two authors are abbreviated with *and colleagues*.
  - The year is surrounded by parentheses (with page numbers, if present).

To provide for both kinds of citations, `language.bst` capitalizes on the fact that `natbib` citation commands come in two flavors. In a typical style compatible with `natbib`, ordinary commands such as `\citet` and `\citep` produce short citations abbreviated with *et al.*, whereas starred commands such as `\citet*` and `\citep*` produce a citation with a full author list. Since *Language* does not require citations with full authors, the style `language.bst` repurposes the starred commands to be used for citing the author. The following table shows how the `natbib` citation commands work with `language.bst`.

Command	Two authors	More than two authors
<code>\citet</code>	Hale & White Eagle (1980)	Sprouse et al. (2011)
<code>\citet*</code>	Hale und White Eagle (1980)	Sprouse and colleagues (2011)
<code>\citep</code>	(Hale & White Eagle 1980)	(Sprouse et al. 2011)
<code>\citep*</code>	(Hale und White Eagle 1980)	(Sprouse and colleagues 2011)
<code>\citealt</code>	Hale & White Eagle 1980	Sprouse et al. 2011
<code>\citealt*</code>	Hale und White Eagle 1980	Sprouse and colleagues 2011
<code>\citealp</code>	Hale & White Eagle 1980	Sprouse et al. 2011
<code>\citealp*</code>	Hale und White Eagle 1980	Sprouse and colleagues 2011
<code>\citeauthor</code>	Hale & White Eagle	Sprouse et al.
<code>\citeauthor*</code>	Hale und White Eagle	Sprouse and colleagues
<code>\citefullauthor</code>	Hale und White Eagle	Sprouse and colleagues

Authors of *Language* articles would typically use `\citet*`, `\citep`, `\citealt` and `\citeauthor*`, though they could use any of the above commands. There is no command for giving a full list of authors.

## Bibliography

The bibliography of this template includes the references of the *language* stylesheet as a sample bibliography.

## **15 General Addenda**

If there are several additions you want to add, but they do not fit into the thesis itself, they belong here.

### **15.1 Detailed Addition**

Even sections are possible, but usually only used for several elements in, e.g. tables, images, etc.

## **16 Figures**

### **16.1 Example 1**

### **16.2 Example 2**

## **Abbildungsverzeichnis**

## **Tabellenverzeichnis**



## Literatur

- BUTT, MIRIAM, und WILHELM GEUDER (Hg.) 1998. *The projection of arguments: Lexical and compositional factors*. Stanford, CA: CSLI Publications.
- CROFT, WILLIAM. 1998. Event structure in argument linking. In Butt & Geuder, 21–63.
- DONOHUE, MARK. 2009. Geography is more robust than linguistics. Science e-letter, 13 August 2009. URL <http://www.sciencemag.org/cgi/eletters/324/5926/464-c>.
- DORIAN, NANCY C. (Hg.) 1989. *Investigating obsolescence*. Cambridge: Cambridge University Press.
- GROPEN, JESS; STEVEN PINKER; MICHELLE HOLLANDER; RICHARD GOLDBERG; und RONALD WILSON. 1989. The learnability and acquisition of the dative alternation in English. *Language* 65.203–57.
- HALE, KENNETH, und JOSIE WHITE EAGLE. 1980. A preliminary metrical account of Winnebago accent. *International Journal of American Linguistics* 46.117–32.
- HASPELMATH, MARTIN. 1993. *A grammar of lezgian*. Walter de Gruyter.
- HYMES, DELL H. 1974a. *Foundations in sociolinguistics: An ethnographic approach*. Philadelphia: University of Pennsylvania Press.
- HYMES, DELL H. (Hg.) 1974b. *Studies in the history of linguistics: Traditions and paradigms*. Bloomington: Indiana University Press.
- HYMES, DELL H. 1980. *Language in education: Ethnolinguistic essays*. Washington, DC: Center for Applied Linguistics.
- MINER, KENNETH. 1990. Winnebago accent: The rest of the data. Lawrence: University of Kansas, MS.
- MORGAN, TJH; NT UOMINI; LE RENDELL; L CHOUINARD-THULY; SE STREET; HM LEWIS; CP CROSS; C EVANS; R KEARNEY; I DE LA TORRE; ET AL. 2015. Experimental evidence for the co-evolution of hominin tool-making teaching and language. *Nature communications* 6.
- PERLMUTTER, DAVID M. 1978. Impersonal passives and the unaccusative hypothesis. *Berkeley Linguistics Society* 4.157–89.
- POSER, WILLIAM. 1984. *The phonetics and phonology of tone and intonation in Japanese*. Cambridge, MA: MIT Dissertation.
- PRINCE, ELLEN. 1991. Relative clauses, resumptive pronouns, and kind-sentences. Paper presented at the annual meeting of the Linguistic Society of America, Chicago.
- RICE, KEREN. 1989. *A grammar of Slave*. Berlin: Mouton de Gruyter.

- SALTZMAN, ELLIOT; HOSUNG NAM; JELENA KRIVOKAPIC; und LOUIS GOLDSTEIN. 2008. A task-dynamic toolkit for modeling the effects of prosodic structure on articulation. *Proceedings of the 4th International Conference on Speech Prosody (Speech Prosody 2008)*, Campinas, 175–84. URL <http://aune.lpl.univ-aix.fr/~sprosig/sp2008/papers/3inv.pdf>.
- VAN DER SANDT, ROB A. 1992. Presupposition projection as anaphora resolution. *Journal of Semantics* 9.333–77.
- SINGLER, JOHN VICTOR. 1992. Review of Melanesian English and the Oceanic substrate, by Roger M. Keesing. *Language* 68.176–82.
- SPROUSE, JON; MATT WAGERS; und COLIN PHILLIPS. 2011. A test of the relation between working memory capacity and syntactic island effects. *Language*, to appear.
- STOCKWELL, ROBERT P. 1993. Obituary of Dwight L. Bolinger. *Language* 69.99–112.
- SUNDELL, TIMOTHY R. 2009. Metalinguistic disagreement. Ann Arbor: University of Michigan, MS. URL <http://faculty.wcas.northwestern.edu/~trs341/papers.html>.
- TIERSMA, PETER M. 1993. Linguistic issues in the law. *Language* 69.113–37.
- WILSON, DEIRDRE. 1975. *Presuppositions and non-truth-conditional semantics*. London: Academic Press.
- YIP, MOIRA. 1991. Coronals, consonant clusters, and the coda condition. *The special status of coronals: Internal and external evidence*, hrsg. von Carole Paradis und Jean-François Prunet, 61–78. San Diego, CA: Academic Press.