



TECHNISCHE HOCHSCHULE NÜRNBERG
GEORG SIMON OHM

Faculty of Computer Science

IT-Based textgeneration using NLP methods

State of the art and design of a prototype

Bachelor Thesis in
Business Informatics and Management

from

Tim Löhr

Student ID 3060802

First advisor: Prof. Dr. Alfred Holl

Second advisor: Prof. Dr. Florian Gallwitz

© 2020

This work and all its parts are (protected by copyright). Any use outside the narrow limits of copyright law without the author's consent is prohibited and liable to prosecution. This applies in particular to duplications, translations, microfilming as well as storage and processing in electronic systems.

Angaben des bzw. der Studierenden:

Name: _____ Vorname: _____ Matrikel-Nr.: _____

Fakultät: Studiengang:

Semester:

Titel der Abschlussarbeit:

Ich versichere, dass ich die Arbeit selbständig verfasst, nicht anderweitig für Prüfungszwecke vorgelegt, alle benutzten Quellen und Hilfsmittel angegeben sowie wörtliche und sinngemäße Zitate als solche gekennzeichnet habe.

Ort, Datum, Unterschrift Studierende/Studierender

Erklärung zur Veröffentlichung der vorstehend bezeichneten Abschlussarbeit

Die Entscheidung über die vollständige oder auszugsweise Veröffentlichung der Abschlussarbeit liegt grundsätzlich erst einmal allein in der Zuständigkeit der/des studentischen Verfasserin/Verfassers. Nach dem Urheberrechtsgesetz (UrhG) erwirbt die Verfasserin/der Verfasser einer Abschlussarbeit mit Anfertigung ihrer/seiner Arbeit das alleinige Urheberrecht und grundsätzlich auch die hieraus resultierenden Nutzungsrechte wie z.B. Erstveröffentlichung (§ 12 UrhG), Verbreitung (§ 17 UrhG), Vervielfältigung (§ 16 UrhG), Online-Nutzung usw., also alle Rechte, die die nicht-kommerzielle oder kommerzielle Verwertung betreffen.

Die Hochschule und deren Beschäftigte werden Abschlussarbeiten oder Teile davon nicht ohne Zustimmung der/des studentischen Verfasserin/Verfassers veröffentlichen, insbesondere nicht öffentlich zugänglich in die Bibliothek der Hochschule einstellen.

Hiermit ☐ genehmige ich, wenn und soweit keine entgegenstehenden Vereinbarungen mit Dritten getroffen worden sind,

☐ genehmige ich nicht,

dass die oben genannte Abschlussarbeit durch die Technische Hochschule Nürnberg Georg Simon Ohm, ggf. nach Ablauf einer mittels eines auf der Abschlussarbeit aufgebrachten Sperrvermerks kenntlich gemachten Sperrfrist

von Jahren (0 - 5 Jahren ab Datum der Abgabe der Arbeit),

der Öffentlichkeit zugänglich gemacht wird. Im Falle der Genehmigung erfolgt diese unwiderruflich; hierzu wird der Abschlussarbeit ein Exemplar im digitalisierten PDF-Format auf einem Datenträger beigelegt. Bestimmungen der jeweils geltenden Studien- und Prüfungsordnung über Art und Umfang der im Rahmen der Arbeit abzugebenden Exemplare und Materialien werden hierdurch nicht berührt.

Ort, Datum, Unterschrift Studierende/Studierender

Preface I

The following thesis was created during my 7th and last semester at the University of Applied Science - Georg Simon OHM.

Within my last three semesters I realized, that my major interest among all IT related topics is artificial intelligence.

Together with my professor *Prof. Dr. Alfred Holl* I worked out a method-matrix for the entire structure of this paper. Without his cooperative support overseas while I was studying abroad at the City University of Hong Kong, this thesis would have not been possible for me.

Even though Natural Language Processing is just a subfield of machine learning, the current state-of-the-art research is far beyond what I can research within a bachelor thesis. In this way, I decided to write my thesis about the subfield *textgeneration* within NLP. My state-of-the-art research includes all *hot topics* within NLP and my prototyp focuses only on the textgeneration part, to dive deeper into what NLP and especially textgeneration is able to accomplish in the year 2020.

Preface II

My interest started basically with my IT project, in which my team and I programmed an autonomously driving remote control car with a deep neural network together with a Raspberry Pi 3. From this first project on, I selected all my elective courses to be related with machine learning or data science in any possible way. I wanted to further increase my knowledge, so I searched for a website which provides courses related to AI. I found *www.udacity.com*, which offers courses in cooperation with top IT companys, such as Google, Airbnb or Microsoft. Out of curiosity, I bought the course *Natural Language Processing*. After successfully finishing it, I was encouraged to write my bachelor thesis in a subfield of *Natural Language Processing*.

For my research I encountered a lot of recently published and old papers from <https://arxiv.org/>. To read through the papers requires a lot of prior knowledge, especially in mathematics, which I learned in my abroad semester in Hong Kong.

Machine Learning and more specifically NLP is not an intuitive study. I provided the common terminologys from top researchers and tried to make the entry into this field as smooth as possible if the reader has no prior knowledge about this topic.

I still recommend some basic linear algebra and calculus knowledge to understand the formulas more easily.

Thank you very much for reading.

Abstract

– At the end , finally finished :) –

Contents

1. Intro	1
1.1. Case study of a current NLP system	2
1.1.1. Case study	2
1.1.2. Useful application areas of NLP systems	2
1.1.3. Useful application areas of NTG systems	2
2. State of the Art	3
2.1. Relevante Aspekte der Mathematik	3
2.2. Geschichte des NLP	3
2.3. Aktuelle Trends der Technologie	3
2.3.1. Einsatzgebiete von NLP-Systemen	3
2.3.2. Einsatzgebiete von NTG-Systemen	3
3. Prototyp	5
3.1. Zielsetzung	5
3.2. Fachkonzept	5
3.2.1. Struktur	6
3.2.2. Neuronales Netz	6
3.2.3. Prozessmodellierung	6
3.2.4. Datenflussmodellierung	6
3.3. Implementierung	6
3.4. Evaluation	7
4. Generierung von Übertragbarem Wissen	9
A. Supplemental Information	11
List of Figures	13
List of Tables	15
List of Listings	17

Chapter 1.

Intro

In recent months and years, neural networks have produced many *state-of-the-art* results in almost all possible disciplines of machine learning. One of the disciplines is Natural Language Processing, *NLP* for short. This term covers many other hot research disciplines, including

- Sentiment Analysis
- Machine Translation
- Voice Recognition
- Text Generation (Neural Text Generation *NTG*)

Another term for text generation is denoted by *Language Modelling*, because it uses the words and grammar as input for the model. In the last 5 years, there were mainly two approaches for modelling NLP, namely the rule-based system and the template-based system (Figure 1.1). Today the *State-of-the-Art* are the neural end-to-end systems, which lead to a far more advanced output [?]. These new systems offer more flexibility and scale with a proportionately better results with less required data, because the complexity and thus the required computing power has increased. However, this fact leads to a complexity problem, because it becomes very difficult to understand the decisions of the neural network. The neural network is basically still a *black box* to a large extent, although it gives surprisingly good results, especially in NLP. Nevertheless, neural network models for text processing are difficult to understand, so nowadays compromises between rule-based systems still have to be made and hybrid systems have to be used.

The neural text generation, also called *NTG*, has many other interesting application fields, including

- Speech recording and conversion to text
- Conversation systems e.g. chatbots
- Text summary

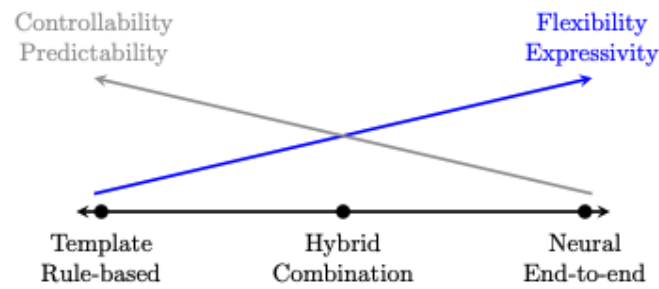


Figure 1.1.: Rule-Based vs. Neural-Text-Generations System [?]

In order to train language models, they must be taught the probability of occurring words in relation to the preceding words. There are several approaches to achieve this goal. Language models can be trained on the level of words, whole sentences or even whole paragraphs. The granularity in which the training takes place is called *n-grams*, where *n* represents the number of preceding words.

1.1. Case study of a current NLP system

1.1.1. Case study

Image-to-Text | Captionbot Microsoft

1.1.2. Useful application areas of NLP systems

IoT, Grammarly, etc

1.1.3. Useful application areas of NTG systems

Chapter 2.

State of the Art

2.1. Relevante Aspekte der Mathematik

Notationen etc.

2.2. Geschichte des NLP

Zeitabfolge der geschichtlichen Hintergründe

2.3. Aktuelle Trends der Technologie

Neuronales Ende-zu-Ende

2.3.1. Einsatzgebiete von NLP-Systemen

Speech Recognition, Machine Translation

2.3.2. Einsatzgebiete von NTG-Systemen

Image-to-Text, Weatherforecast

Chapter 3.

Prototyp

In this chapter, we're actually using some code!

```
1 x = 1
2 if x == 1:
3     # indented four spaces
4     print("x is 1.")
```

Listing 3.1: This is an example of inline listing

You can also include listings from a file directly:

```
1 x = 1
2 if x == 1:
3     # indented four spaces
4     print("x is 1.")
```

Listing 3.2: This is an example of included listing

3.1. Zielsetzung

Image Captioning

3.2. Fachkonzept

Fachkonzept

3.2.1. Struktur

The different steps of Text Generation

- Importing Dependencies
- Loading the Data
- Creating Character/Word mappings
- Data Preprocessing
- Modelling
- Generating text

3.2.2. Neuronales Netz

LSTM <https://www.analyticsvidhya.com/blog/2017/12/fundamentals-of-deep-learning-introduction-to-lstm/>

Experimenting with different models

- A more trained model
- A deeper model
- A wider model
- A gigantic model

3.2.3. Prozessmodellierung

Funktionen etc.

3.2.4. Datenflussmodellierung

Diagramm

3.3. Implementierung

Code

3.4. Evaluation

Print Ergebnisse

Bild

Image Caption

Chapter 4.

Generierung von Übertragbarem Wissen

Modulare Erweiterbarkeit meines Projekts. Einordnung in Gesellschaftlichen Kontext

Appendix A.

Supplemental Information

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

This is the second paragraph. Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

And after the second paragraph follows the third paragraph. Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

After this fourth paragraph, we start a new paragraph sequence. Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are

written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

List of Figures

1.1. Rule-Based vs. Neural-Text-Generations System [?]	2
--	---

List of Tables

List of Listings

3.1. This is an example of inline listing	5
3.2. This is an example of included listing	5

