



project work

HS16 Studiengang Informatik

Spaghetti Bolognese mit Parmesan

Autoren

Raphael Emberger, Kal-El,
Musashi Miyamoto

Datum

20. November 2017

Bitte füllen Sie das Titelblatt aus und berücksichtigen Sie Folgendes:

-> Bitte auf keinen Fall Schriftart und Schriftgrösse ändern. Text soll lediglich überschrieben werden!

-> Bitte pro Tabellenzeile max. 4 Textzeilen!

- Titel: Fügen Sie Ihren Studiengang direkt nach dem Wort „Project Work“ ein (max. 2 Zeilen).
- Titel der Arbeit: Überschreiben Sie den Lauftext mit dem Titel Ihrer Arbeit (max. 4 Zeilen).
- Autoren: Tragen Sie Ihre Vor- und Nachnamen ein (alphabetisch nach Name).
- Betreuer: Tragen Sie Ihren Betreuer / Ihre Betreuer ein (alphabetisch nach Name).
- Ohne Nebenbetreuung, Industriepartner oder externe Betreuung, ganze Tabellenzeile löschen.
- Am Schluss löschen Sie den ganzen Beschrieb (grau) und speichern das Dokument als pdf. ab.

Zusammenfassung

- Zusammenfassung

Vorwort

- Stellt den persönlichen Bezug zur Arbeit dar und spricht Dank aus.

DECLARATION OF ORIGINALITY

Project Work at the School of Engineering

By submitting this project work, the undersigned student confirms that this thesis is his/her own work and was written without the help of a third party. (Group works: the performance of the other group members are not considered as third party).

The student declares that all sources in the text (including Internet pages) and appendices have been correctly disclosed. This means that there has been no plagiarism, i.e. no sections of the Bachelor thesis have been partially or wholly taken from other texts and represented as the student's own work or included without being correctly referenced.

Any misconduct will be dealt with according to paragraphs 39 and 40 of the General Academic Regulations for Bachelor's and Master's Degree courses at the Zurich University of Applied Sciences (Rahmenprüfungsordnung ZHAW (RPO)) and subject to the provisions for disciplinary action stipulated in the University regulations.

City, Date:

Signature:

.....

.....

.....

.....

The original signed and dated document (no copies) must be included after the title sheet in the ZHAW version of all Bachelor thesis submitted.

Inhaltsverzeichnis

1. Introduction	8
1.1. Initial Position	9
1.2. Task	10
2. Theoretical Principles	11
2.1. Definitions	11
2.2. Recurrent Neural Networks	11
2.3. Seq2Seq	11
2.4. Attention	11
2.5. Performance Evaluation	11
3. Experiments	12
4. Results	13
5. Discussion and Prospects	14
6. Index	15
6.1. Literaturverzeichnis	15
6.2. Glossary	16
6.3. Abbildungsverzeichnis	17
6.4. Tabellenverzeichnis	18
6.5. Symbol Glossary	19
6.6. Acronym Glossary	20
A. Appendix	22
A.1. Projektmanagement	22
A.2. Final Words	22

Abstract

- Summary

This is just some normal text that goes here

Preface

- Stellt den persönlichen Bezug zur Arbeit dar und spricht Dank aus.

thank-yous go here

1. Introduction

The rapid pace at which the human race has overcome barriers of communication in the past 100 years is astounding. Starting at the introduction of traditional carrier mail [1] to the invention of the internet [citation needed], subjectively, one of our only common goals is the improvement of communication. [citation needed]

However, the internet is unique in its ability to connect us nearly instantaneously with people from all over the globe at the touch of a button. This brings into focus one of the last remaining barriers of communication we have to overcome. Languages.

In 2009 we spoke 6909 unique languages [?]. Therefore, finding a way to efficiently translate between any of these languages is key to further enable interaction. The three main challenges which present themselves are as follows:

1. Ambiguity - Words may contain multiple meanings and depending on their context, only a subset apply.
2. Non-Standard Terminology - This refers to the use of language constructs which do not adhere to the official language documentation. A popular example of this is the abbreviations and emojis used in tweets.
3. Named Entities - To a machine, a name appears like any other word. It's a set of characters. However, they are typically not translatable and thus the machine needs to be in a position to accurately identify named entities.

The rapid advancements in machine learning as well as increase in computational power has enabled computer scientists develop machine translators. Past machine assisted translation software as well as current iterations can be associated with one of the following five groups.[Citation Needed].

1. Rule-Based
2. Statistical
3. Example-Based
4. Hybrid Machine Translation
5. Neural Machine Translation

Although different implementations of the Neural Machine Translation (hence forth NMT) approach varies widely in terms of accuracy, speed and reliability, it has, in general, shown the most promise of succeeding.

This has led to an increase in attention from the scientific community resulting in a large number of slight variations, each claiming to produce better results.

This report focuses on two main aspects:

1. Build a functioning NMT and evaluate it with the field's common methods
2. Rebuild and reproduce the results of a number of recently published papers in order to
 - a) Compare our model against the current state of the art.
 - b) Reproduce, verify and compare the results of other models with our own.

- c) Determine and why our model succeeded or fell short of expectations when compared to others.

This paper is structured as follows: Chapter two explains the fundamental concepts of machine learning, homing in on terminology and theory commonly applied in NMT and common evaluation methods. Chapter three is dedicated to show casing our selected approach while chapter four briefly summarizes the selected approaches. Chapter five and six details the selected data and methodology. The experiments, results and the comparison is described in chapter seven while chapter eight and nine summarizes our learnings and possible next steps.

- humans are the most communicative mammal on this planet (source needed)
- our species has relentlessly pursued the improvement of communication methods. First carrier mail, then the telegram followed by the radio, telephone, tv and most recently, the internet.
- The last innovation enabled us truly to connect with [insert percentage of people on the planet who have access to internet] of earth's population.
- One major challenge remains. Language barriers.
- Rise of machine learning has begun to help us overcome this barrier. [Examples of Microsoft (Skype), Google (translate/nmt) and DeepI]
- machines becoming better at processing human language (accuracy)
- translations are still not always 100% correct due to factors such as
 - double meaning
 - context
 - sentiment (sarcasm vs. criticism)
- this project is aimed at building our own neural machine translation agent, rebuilding current nmts and verifying their published performance as well as measure our nmt against the published performance and our measured performance
- attention to correct words
- database structures
- multiple ways to ask for identical information
- multiple solutions proposed
- KBQA: Learning Question Answering over QA Corpora and Knowledge Bases
- Eric, Manning - 2017 - Key-Value Retrieval Networks for Task-Oriented Dialogue - With Highlights
- Asking your Assistant (Google, Siri or S-Voice) weather you have an appointment tomorrow and ask follow-up questions about this appointment is currently not possible (due to above challenges but could be if these papers prove implementable)

1.1. Initial Position

- No Response from KBQA for Code
- Refusal to share code from Manning
- Ultimate new goal: Implement Manning's solution without his code

- Nennt bestehende Arbeiten/Literatur zum Thema -> Literaturrecherche
- Stand der Technik: Bisherige Lösungen des Problems und deren Grenzen
- (Nennt kurz den Industriepartner und/oder weitere Kooperationspartner und dessen/deren Interesse am Thema Fragestellung)

1.2. Task

- Small Steps
 - implement seq2seq network for translation
 - * implement char-based
 - * implement word-based
 - * try multiple different implementations (reversed-input, multiple LSTMs) and compare against each other
 - * get decent results on both and move on
 - implement seq2seq with attention
 - * attempt various attention mechanism
 - One Large Step
 - map best working models and tools to KBQA and get better results than Stanford
 - Rub better results in Eric's face.
 - Profit.
-
- Formuliert das Ziel der Arbeit
 - Verweist auf die offizielle Aufgabenstellung des/der Dozierenden im Anhang
 - (Pflichtenheft, Spezifikation)
 - (Spezifiziert die Anforderungen an das Resultat der Arbeit)
 - (Übersicht über die Arbeit: stellt die folgenden Teile der Arbeit kurz vor)
 - (Angaben zum Zielpublikum: nennt das für die Arbeit vorausgesetzte Wissen)
 - (Terminologie: Definiert die in der Arbeit verwendeten Begriffe)

2. Theoretical Principles

The first part of this chapter is dedicated to explaining fundamental terminology and theoretical concepts in machine learning, followed by concepts specific to NMT.

It is important to note that the following definitions and explanations are restricted to information specific to NMT. As such concepts related to unsupervised learning are excluded from this report.

2.1. Definitions

Reference A *reference* refers to a the correct translation against which the translation produced by the NMT can be measured. It is either an entire paragraph, sentence or word. Traditionally a *reference* refers to one complete sentence.

Hypothesis In the context of machine translation a *hypothesis* refers to the output produced by the NMT, given an input sentence.

2.2. Recurrent Neural Networks

- Standard Neural Networks
- Recurrent Neural Networks
 - Problems
 - Solutions

2.3. Seq2Seq

- encoder
- decoder

2.4. Attention

- Mechanisms

2.5. Performance Evaluation

- Translation
 - Bleu
 - others
- KB-Retrieval
 - Bleu
 - sent2vec

3. Experiments

- Folgende waren schlechter als in deren Literatur beschrieben
 - Keras Tutorial Char-Based
 - MagicMagic Keras Char = Keras Word
 - Our attempt at Word Based
- Hidden State(?)
- Used Keras Tutorial from TF-Talk
- Google NMT
- Find out why ours didn't perform as well as the above two

- (Beschreibt die Grundüberlegungen der realisierten Lösung (Konstruktion/Entwurf) und die Realisierung als Simulation, als Prototyp oder als Software-Komponente)
- (Definiert Messgrößen, beschreibt Mess- oder Versuchsaufbau, beschreibt und dokumentiert Durchführung der Messungen/Versuche)
- (Experimente)
- (Lösungsweg)
- (Modell)
- (Tests und Validierung)
- (Theoretische Herleitung der Lösung)

4. Results

- (Zusammenfassung der Resultate)

5. Discussion and Prospects

Wie in XXX nachzulesen, gibt es sogenannte Gleichungen. Ω

- Bespricht die erzielten Ergebnisse bezüglich ihrer Erwartbarkeit, Aussagekraft und Relevanz
- Interpretation und Validierung der Resultate
- Rückblick auf Aufgabenstellung, erreicht bzw. nicht erreicht
- Legt dar, wie an die Resultate (konkret vom Industriepartner oder weiteren Forschungsarbeiten; allgemein) angeschlossen werden kann; legt dar, welche Chancen die Resultate bieten

6. Index

6.1. Literaturverzeichnis

Union, U. P. (2016), 'History'.

URL: <http://www.upu.int/en/the-upu/history/about-history.html> 8

6.2. Glossary

élite select group or class 16

élitism advocacy of dominance by an élite 14

6.3. Abbildungsverzeichnis

6.4. Tabellenverzeichnis

6.5. Symbol Glossary

Ω unit of electrical resistance 14

6.6. Acronym Glossary

HRZ Hochschulrechenzentrum 14

[title=Index]

A. Appendix

A.1. Projektmanagement

- Offizielle Aufgabenstellung, Projektauftrag
- (Zeitplan)
- (Besprechungsprotokolle oder Journals)

A.2. Final Words

- CD mit dem vollständigen Bericht als pdf-File inklusive Film- und Fotomaterial
- (Schaltpläne und Ablaufschemata)
- (Spezifikationen u. Datenblätter der verwendeten Messgeräte und/oder Komponenten)
- (Berechnungen, Messwerte, Simulationsresultate)
- (Stoffdaten)
- (Fehlerrechnungen mit Messunsicherheiten)
- (Grafische Darstellungen, Fotos)
- (Datenträger mit weiteren Daten(z. B. Software-Komponenten) inkl. Verzeichnis der auf diesem Datenträger abgelegten Dateien)
- (Softwarecode)