# Refining the taxonomy of Wikidata with neural networks

## Bachelorarbeit

zur Erlangung des Grades einer Bachelor of Science (B.Sc.)
im Studiengang Informatik

vorgelegt von

## Alex Baier

Erstgutachter:    Prof. Dr. Steffen Staab
Institute for Web Science and Technologies

Zweitgutachter:    Max Mustermann
Institute for Web Science and Technologies

Koblenz, im  Dezember 2016

# Erklärung

Ich versichere, dass ich die vorliegende Arbeit selbständig verfasst und keine anderen als die angegebenen Quellen und Hilfsmittel benutzt habe.

|  | Ja | Nein |
|---|---|---|
| Mit der Einstellung dieser Arbeit in die Bibliothek bin ich einverstanden. | ☐ | ☐ |
| Der Veröffentlichung dieser Arbeit im Internet stimme ich zu. | ☐ | ☐ |
| Der Text dieser Arbeit ist unter einer Creative Commons Lizenz verfügbar. | ☐ | ☐ |
| Der Quellcode ist unter einer Creative Commons Lizenz verfügbar. | ☐ | ☐ |
| Die erhobenen Daten sind unter einer Creative Commons Lizenz verfügbar. | ☐ | ☐ |

....................................................................................................

(Ort, Datum)                                                          (Unterschrift)

# Contents

# List of Figures

# 1 Introduction

# 2 Foundations

## 2.1 Definitions

**Definition 1** (Directed Graph). A graph G is an ordered pair $G = (V, E)$, where $V$ is a set of vertices, and $E = \{(v_1, v_2) \mid v_1, v_2 \in V\}$ is a set of ordered pairs called directed edges, connecting the the vertices.

**Definition 2** (Predecessor, Successor). Let $G = (V, E)$ be a directed graph.
$v_1 \in V$ is a predecessor of $v_2 \in V$, if there exists an edge so that $(v_1, v_2) \in E$.
Let $v \in V$ be a vertice of G, then $pred(v) = \{w \mid (w, v) \in E\}$ is the set of predecessors of $v$.
$v_1 \in V$ is a successor of $v_2 \in V$, if there exists an edge so that $(v_2, v_1) \in E$.
Let $v \in V$ be a vertice of G, then $succ(v) = \{w \mid (v, w) \in E\}$ is the set of successors of $v$.

**Definition 3** (Walk). Let $G = (V, E)$ be a directed graph.
A walk $W$ of length $n \in \mathbb{N}$ is a sequence of vertices $W = (v_1, \ldots, v_n)$ with $v_1, \ldots, v_n \in V$, so that $(v_i, v_{i+1}) \in E \; \forall i = 1, \ldots, n-1$.

**Definition 4** (Cycle). A walk $W = (v_1, \ldots, v_n)$ of length $n$ is called a cycle, if $v_1 = v_n$.

**Definition 5** (Path). A walk $P = (v_1, \ldots, v_n)$ is a path from $v_1$ to $v_n$, if $v_i \neq v_j$ for all $i, j = 1, \ldots, n$ with $i \neq j$.

**Definition 6** (Acyclic Graph). A directed graph $G$ is called acyclic graph, if there are no cycles in $G$.

**Definition 7** (Statement). **TODO 1: Define statement.**

**Definition 8** (Class). A class is a tuple $(id, label, Statements, Instances, wiki)$:

- $id \in \mathbb{N}$, which is a numerical Wikidata item ID;

- $label$, which is the, to $id$ corresponding, English label in Wikidata;

- $Statements$ is a set of statements about the class;

- $Instances \in \mathcal{P}(\mathbb{N})$ is the set of numerical Wikidata item IDs, which are instances of the class;

- $wiki$ is the, to the class corresponding, English Wikipedia article text.

**Definition 9** (Taxonomy). A taxonomy $T = (C, S)$ is a acyclic graph, where $C$ is a set of classes, and $S$ is a set of subclass-of relations between these classes.

1

**Definition 10** (Subclass Relation). Let $T = (C, S)$ be a taxonomy.
The transitive, ordered relation $\lhd_{subclass}$ is defined.
Let $c_1, c_2 \in C$. $c_1 \lhd_{subclass} c_2$, if there is a path $P = (c_1, \ldots, c_2)$ from $c_1$ to $c_2$ in $T$.

**Definition 11** (Root class). Let $T = (C, S)$ be a taxonomy.
$r \in C$ is called root class of $T$, if $|succ(r)| = 0$.
$root(T) = \{r \in C \mid |succ(r)| = 0\}$ is the set of all root classes in $T$.

Finally we can define our problem as the following task:

**Problem.** Let $W_1$ be the taxonomy of Wikidata, where only labeled root classes are considered. On 7th November 2016 the following state applies $|root(W_1)| = 5332$.
$W_1 = (C, S)$ is the input for the described problem.
Let $W_2$ be the refined output taxonomy.
A refinement method is needed to significantly reduce the number of root classes in the Wikidata taxonomy. After the refinement method is applied on $W_1$, which outputs $W_2$, the following should be true: $|root(W_2)| \ll |root(W_1)|$.

The refinement process can be reduced to the following smaller task:
Let $r \in root(W_1)$.
Find a $c \in C$ with $\neg(c \lhd_{subclass} r)$, so that $c$ is the most similar super class of $r$.
Connecting $r$ to $c$ with an edge produces the output taxonomy $W_2 = (C, S \cup \{(r, c)\})$.
Accordingly $|root(W_2)| = |root(W_1)| - 1$ applies.
Repeating this smaller task will eventually yield $|W_2| \ll |W_1|$.

The problem can therefore be defined as developing a method, which finds, given a taxonomy $W = (C, S)$ and a root class $r = root(W)$, the most similar superclass of $r$.

# 3  Analysis of Wikidata's taxonomy

# 4  Analysis and comparison of neural networks

# 5  Architecture of neural network

# 6  Implementation

## 6.1  Collection of training data

## 6.2  Implementation of network

## 6.3  Training

# 7  Evaluation

## 7.1  Execution on all root classes

## 7.2  Survey with Wikidata community

# References

[1] Wikidata game. `https://tools.wmflabs.org/wikidata-game/`.

[2] Wikidata statistics. `https://tools.wmflabs.org/wikidata-todo/stats.php`.

[3] Ebru Arisoy, Tara N. Sainath, Brian Kingsbury, and Bhuvana Ramabhadran. Deep neural network language models. In Proceedings of the NAACL-HLT 2012 Workshop: Will We Ever Really Replace the N-gram Model? On the Future of Language Modeling for HLT, WLM '12, pages 20–28, Stroudsburg, PA, USA, 2012. Association for Computational Linguistics.

[4] Shaosheng Cao, Wei Lu, and Qiongkai Xu. Deep neural networks for learning graph representations. In Dale Schuurmans and Michael P. Wellman, editors, AAAI, pages 1145–1152. AAAI Press, 2016.

[5] Claudia d'Amato, Steffen Staab, Andrea G. B. Tettamanzi, Tran Duc Minh, and Fabien L. Gandon. Ontology enrichment by discovering multi-relational association rules from ontological knowledge bases. In Sascha Ossowski, editor, SAC, pages 333–338. ACM, 2016.

[6] Karol Gregor, Ivo Danihelka, Alex Graves, and Daan Wierstra. DRAW: A recurrent neural network for image generation. CoRR, abs/1502.04623, 2015.

[7] Nal Kalchbrenner, Edward Grefenstette, and Phil Blunsom. A convolutional neural network for modelling sentences. CoRR, abs/1404.2188, 2014.

[8] Alexander Maedche, Viktor Pekar, and Steffen Staab. Ontology Learning Part One — on Discovering Taxonomic Relations from the Web, pages 301–319. Springer Berlin Heidelberg, Berlin, Heidelberg, 2003.

[9] Alexander Maedche and Steffen Staab. Ontology learning for the semantic web. IEEE Intelligent Systems, 16(2):72–79, March 2001.

[10] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. CoRR, abs/1301.3781, 2013.

[11] Tomas Mikolov, Martin Karafiát, Lukás Burget, Jan Cernocký, and Sanjeev Khudanpur. Recurrent neural network based language model. In INTERSPEECH 2010, 11th Annual Conference of the International Speech Communication Association, Makuhari, Chiba, Japan, September 26-30, 2010, pages 1045–1048, 2010.

[12] Viktor Pekar and Steffen Staab. Taxonomy learning - factoring the structure of a taxonomy into a semantic classification decision. In Proceedings of the 19th Conference on Computational Linguistics, COLING-2002, August 24 - September 1, 2002 , Taipei, Taiwan, 2002, 2002.

[13] Giulio Petrucci, Chiara Ghidini, and Marco Rospocher. Using recurrent neural network for learning expressive ontologies. CoRR, abs/1607.04110, 2016.

[14] Alessandro Sperduti and Antonina Starita. Supervised neural networks for the classification of structures. IEEE Transactions on Neural Networks, 8(3):714–735, may 1997.

[15] Serge Stratan. Wikidata taxonomy browser (beta). `http://sergestratan.bitbucket.org/`.

[16] Wilson Wong, Wei Liu, and Mohammed Bennamoun. Ontology learning from text: A look back and into the future. ACM Comput. Surv., 44(4):20:1–20:36, September 2012.

[17] Guoqiang Peter Zhang. Neural networks for classification: a survey. In and Cybernetics - Part C: Applications and Reviews, 2000.