

Masters Thesis in Computer Science

Extending Semantic Hypergraphs by Neural Embedding-based Semantic Similarity for Pattern Matching

Max Reinhard

Matrikelnummer: 359417

June 22, 2023

Supervised by Prof. Dr. Manfred Hauswirth
Additional guidance by Prof. Dr. Camille Roth* and Dr. Thilo Ernst[†]

*Centre Marc Bloch (An-Institut der Humboldt-Universität zu Berlin)

[†]Fraunhofer-Institut für offene Kommunikationssysteme

Abstract Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

Contents

1	Introduction	4
2	Fundamentals and Related Work	5
2.1	Semantic Hypergraph	5
2.1.1	Structure	5
2.1.2	Syntax	5
2.2	Semantic Similarity	5
2.2.1	Different Similarity Measures	5
2.2.2	Types of Semantic Similarity	5
2.3	Embedding-based Similarity	5
2.3.1	Embedding Types	5
2.3.2	Distance Measures	5
3	Problem Statement	6
3.1	Research Questions	6
4	Solution Approach	7
4.1	Integration into the Pattern Matching Process	7
4.1.1	semsim Functional Pattern	7
4.1.2	Sub-pattern Similarity Thresholds	7
4.2	Fixed Word Embedding-based Matching	7
4.3	Contextual Embedding-based Matching	7
4.3.1	Context References	7
4.3.2	Token Mapping	7
4.4	Similarity Threshold Control	7
4.4.1	Breakpoint Discovery	7
5	Implementation	8
5.1	Relevant external Software Libraries used	8
5.2	Modules newly added to the SH Framework	8
5.3	Modifications of the SH Pattern Matching	8
5.4	Modifications to the Hypergraph database	8
6	Results and Evaluation	9
6.1	Case Study: Conflicts	9
6.1.1	Quantitative Results	10
6.1.2	Qualitative Results	10
7	Conclusion	12
8	Future Work	13
8.1	Implementation Improvements	13
8.2	Further Evaluations	13

1 Introduction

- Context: The big problem
- Problem statement: The small problem
- Methodology / Strategy
- Structure

Notes:

- Huge amounts of text, which can provide insight about stuff
- Automatic tools can provide assistance for humans to process all the text
- This generally means filtering the original text corpus or otherwise reducing amount of information the information that has to be processed by humans
- Filtering introduces a bias
- Especially for scientific purposes it is relevant to mitigate bias or at least understand what bias has been introduced (to make it transparent)
- Semantic Hypergraphs can be a valuable tool for that because...

Human life in times of widespread use of the internet and smartphones is most certainly more than ever interspersed with text-based communication...

A Semantic Hypergraphd (**menezes_semantic_2021**) is a form of representation for Natural Language (NL) and therefore knowledge. *NL* sentences can be modelled as a recursive hypergraph which can be represented in a formal language. The framework allows to specify semantic patterns in this formal language which can be matched against an existing *SH*.

The aim of the SH framework is to provide a *open* and *adaptive* framework to analyse text corpora, especially in the domain of computational social science (CSS) (**lazer2009computational**). The framework is *open* in the sense that it's representation formalism is inspectable and intelligible by humans and that the pattern matching follows explicit rules. The framework is adaptive in the sense that the parsing is based on adaptive subsystems (ML-based) and therefore allows for an error-tolerant parsing from *NL* to *SH* in regards to grammatical and syntactical correctness (???).

2 Fundamentals and Related Work

2.1 Semantic Hypergraph

2.1.1 Structure

2.1.2 Syntax

Square bracket notation

2.2 Semantic Similarity

2.2.1 Different Similarity Measures

String Similarity

Lievenstein distance, etc..

Lexical Similarity

tf-idf, etc.?

2.2.2 Types of Semantic Similarity

Lexical Databases

WordNet and alike (not the scope of this work)

2.3 Embedding-based Similarity

2.3.1 Embedding Types

Fixed Word Embeddings

Contextual Ebeddings

2.3.2 Distance Measures

Mean reference vector vs. pairwise distance

3 Problem Statement

3.1 Research Questions

4 Solution Approach

Combining Semantic Hypergraphs with neural embeddings

4.1 Integration into the Pattern Matching Process

4.1.1 `semsim` Functional Pattern

pattern works only for atoms

4.1.2 Sub-pattern Similarity Thresholds

4.2 Fixed Word Embedding-based Matching

word2vec via gensim

discussion about using transformer models for single word embeddings?

single-word and multi-word reference

Square bracket notation

4.3 Contextual Embedding-based Matching

4.3.1 Context References

4.3.2 Token Mapping

4.4 Similarity Threshold Control

4.4.1 Breakpoint Discovery

detect change points in number of matches

see <https://github.com/deepcharles/ruptures>

half-max point and quarter/three-quarter points (percentiles, not quantiles) fit function and search for inflection as well as maximum derivative points, problematic in cases with less continuous change in number of matches.

how to approach this for practical applications?

5 Implementation

5.1 Relevant external Software Libraries used

Here list libs and models to be referenced later.

Word2Vec Gensim SentenceTransformers Transformers SpaCy

5.2 Modules newly added to the SH Framework

5.3 Modifications of the SH Pattern Matching

5.4 Modifications to the Hypergraph database

6 Results and Evaluation

In this chapter...

6.1 Case Study: Conflicts

This case study follows the approach presented in (menezes_semantic_2021) where expressions of conflict are extracted from a SH constructed from a corpus of news titles that were shared on the social media platform *Reddit*. Specifically all titles shared between January 1st, 2013 and August 1st, 2017 on *r/worldnews*¹, which is described as: “A place for major news from around the world, excluding US-internal news.” (Number of headers: 479384)

Pattern 1 is used to extract conflicts between two parties, where the SOURCE shows some form of aggression against the TARGET, potentially regarding some TOPIC:

$$(\text{ PRED/P.so,x SOURCE/C TARGET/C [against,for,of,over]/T TOPIC/[RS] }) \wedge \\ (\text{ lemma/J >PRED/P [accuse,arrest,clash,condemn,kill,slam,warn]/P })$$

Pattern 1: Conflict pattern

To investigate whether it is possible to capture the abstract concept of a country using the multi-word `semsim` pattern introduced in 4.2, a list of the worlds 20 most populous countries (`wiki_list_of_countries`) is used (listed in descending order by population size):

India, China, USA, Indonesia, Pakistan, Nigeria, Brazil, Bangladesh, Russia, Mexico, Japan, Philippines, Ethiopia, Egypt, Vietnam, Congo, Iran, Turkey, Germany, France

To avoid the repetition of that list in the following pattern, we introduce a variable:

COUNTRIES = [india,china,usa,indonesia,pakistan,nigeria,brazil,bangladesh,russia,mexico,ajapan,philippines,ethiopia,egypt,vietnam,congo,iran,turkey,germany,france]

Pattern 2: Countries variable

Pattern 3 shows the resulting pattern for conflicts between countries:

$$(\text{ PRED/P.so,x SOURCE/C TARGET/C semsim [against,for,of,over]/T TOPIC/[RS] }) \wedge \\ (\text{ semsim/J >/PRED/P [accuse,arrest,clash,condemn,kill,slam,warn]/P }) \wedge \\ (\text{ semsim/J >SOURCE/C COUNTRIES/C }) \wedge (\text{ semsim/J >TARGET/C COUNTRIES/C })$$

Pattern 3: Country conflict pattern

¹<http://reddit.com/r/worldnews>

In pattern 4 a sub-pattern specific threshold $t_{sim}^{countries}$ for the countries **semsim** sub-pattern is introduced.

$$\begin{aligned} & (\text{PRED/P.so,x SOURCE/C TARGET/C semsim [against,for,of,over]/T TOPIC/[RS] }) \wedge \\ & \quad (\text{semsim/J >/PRED/P [accuse,arrest,clash,condemn,kill,slam,warn]/P }) \wedge \\ & \quad \quad (\text{semsim/J >SOURCE/C COUNTRIES/C } t_{sim}^{countries}) \wedge \\ & \quad \quad (\text{semsim/J >TARGET/C COUNTRIES/C } t_{sim}^{countries}) \end{aligned}$$

Pattern 4: Country conflict pattern

6.1.1 Quantitative Results

Pattern 3 is matched against the described *Reddit r/worldnews* hypergraph. The similarity threshold t_{sim} for the **semsim** function (see ??) is varied. t_{sim} is either varied for the entire pattern or for a specific **semsim** sub-pattern.

6.1.2 Qualitative Results

Table 6.1: hyper dyper table

Scenario Name	Pattern	Samples	Variable Threshold	Reference Edges	Ref. Edges Source
1_original-pattern	Pattern X	Erdogan slams ridicule of 'Muslims discovered Americas' claim Iran forces 'kill Kurdish rebels on Iraq border Ukraine Accuses Russia of Invasion	-/-	-/-	-/-
2-1_sensim-fix_preds	Pattern X	Pakistani police kill feared militant leader in mysterious pre-dawn shootout Al-Shabaab militants claim responsibility for deadly attack on Garissa University College in Kenya Casualties as Congo troops, UN forces fight rebels	'preds': 0.19 Percentile: 50	-/-	-/-
2-2_sensim-fix_preds	Pattern X	Iranian police have arrested merchants for selling clothing that featured the flags of the United States and Britain, two longtime foes of the Islamic republic Syrian Air Force Strikes kill 38 ISIS fighters Seven Libyan soldiers killed fighting off Islamists near Benghazi: source	'preds': 0.54 Percentile: 50	-/-	-/-

7 Conclusion

8 Future Work

8.1 Implementation Improvements

implemnt multiprocessing, i.e. server process for both hypergraph and semsim matchers.
other option would be to leverage python shared memory capabilities but is likely to be less stable and has less scaling potential

8.2 Further Evaluations