

1 Rule-based geographical  
2 relationship extraction in Dutch  
3 news articles

4 Antoon W. Meijer  
5 11016914

6 Bachelor thesis  
7 Credits: 18 EC

8 Bachelor Opleiding Kunstmatige Intelligentie

9 University of Amsterdam  
10 Faculty of Science  
11 Science Park 904  
12 1098 XH Amsterdam

13 *Supervisor*  
14 dr. M.J. Marx

15 Innovation Center for Artificial Intelligence(?)  
16 Faculty of Science  
17 University of Amsterdam  
18 Science Park 904  
19 1098 XH Amsterdam

20 January 31st, 2019

<sup>21</sup> **Abstract**

## Contents

23	<b>1 Introduction</b>	<b>3</b>
24	1.1 Research question . . . . .	3
25	1.2 Overview . . . . .	3
26	<b>2 Related work</b>	<b>4</b>
27	<b>3 Method</b>	<b>5</b>
28	3.1 Dataset . . . . .	5
29	3.2 Toponym detection . . . . .	5
30	3.3 Data selection . . . . .	5
31	3.4 Spatial identifier selection . . . . .	6
32	3.5 Parsing spatial identifiers to predicates . . . . .	7
33	3.6 Converting toponym-predicate pairs into coordinates . . . . .	8
34	3.6.1 Coordinate source & search method . . . . .	8
35	3.6.2 ON(1) . . . . .	8
36	3.6.3 IN(1) . . . . .	9
37	3.6.4 AT(1) . . . . .	9
38	3.6.5 BETWEEN(2) . . . . .	9
39	3.6.6 INTERSECT(1/2) . . . . .	9
40	3.7 Predicate result merging . . . . .	10
41	3.8 Resulting output . . . . .	10
42	<b>4 Results</b>	<b>11</b>
43	4.1 Evaluation method . . . . .	11
44	<b>5 Discussion</b>	<b>12</b>
45	<b>6 Conclusion</b>	<b>13</b>

# 1 Introduction

As the amount of published articles and other Dutch text on the Internet continues to increase, it would be helpful for researchers and others interested in location-specific articles to be able to search for articles based on locations mentioned in the article. It would enable them to spend less time sifting through irrelevant articles, and to spend more time doing their own research. A geolocating model could also be beneficial for other purposes, such as automatic traffic accident location statistics.

Olieman et al<sup>[1]</sup> have attempted to accommodate this need by creating an application, LocLinkVis, that plots Dutch input text on a map based on geographical entities in the text. By doing so, users can more quickly see if an article is relevant. The approach taken by Olieman et al<sup>[1]</sup> identifies toponyms in text by using a gazetteer populated with OpenStreetMap data. After toponym disambiguation, the most likely candidate is shown on a map. However, LocLinkVis does not take into account that some locations might consist of more than one toponym. In such a case, the application simply shows the location of all involved toponyms, instead of understanding how these toponyms are linked and need to be interpreted as one location.

## 1.1 Research question

This research is aimed at improving the accuracy of geolocating complex locations mentioned in Dutch articles, through the use of toponyms and surrounding spatial identifier words. Accurately geolocating complex locations is especially important, as non-complex locations are already handled with sufficient accuracy by commonly-used geographical search engines such as OpenStreetMap and Google Maps. In this research, complex locations are defined as locations that involve multiple toponyms to accurately describe. Below is an example of a complex location:

*Een 25-jarige man uit Arnhem is dinsdagmiddag om het leven gekomen bij een ongeval op de N59 ter hoogte van Oude-Tonge.*

In the example above, the location is not the entire length of the N59 road, but only the section near Oude-Tonge. Such locations are usually not accurately geolocated by the common geographical search engines. This research thus aims to answer the following question:

*In order to accurately geolocate a complex location, how can the relationship between several toponyms extracted from Dutch text best be modeled?*

*Using spatial identifier words to model the relationship between toponyms, how accurately can complex locations be geolocated?*

## 1.2 Overview

[NOG SCHRIJVEN]

## <sup>90</sup> 2 Related work

## 91 3 Method

92 This section outlines ... [INTRO SCHRIJVEN, SAMENVATTING VAN METHOD  
93 SECTIE]

### 94 3.1 Dataset

95 Even though the scope of the research is wider, the dataset that was used solely  
96 consists of traffic accident news articles. This choice was made because of the  
97 fact that traffic accident articles almost always contain toponyms to describe  
98 where the accident took place, with a relatively high degree of complex locations  
99 as well, thus making it interesting for this research.

100  
101 The dataset used originates from the website *flitsservice.nl* and has been edited  
102 by Hendriks<sup>[2]</sup> to separate each article’s body from its title. The titles of the  
103 articles will not be used in this research, because these never contain complex  
104 locations due to their short length, and title toponyms are usually repeated  
105 in the article. The dataset contains 7784 articles. For testing purposes, the  
106 dataset has been split in two parts. The first 4998 articles are used as training  
107 set, the last 2786 are used as testing set. The publishing dates of the articles  
108 vary between 2003 and 2019.

### 109 3.2 Toponym detection

110 The first step to geolocating an article is to find the toponyms in the text. This  
111 was done using the model constructed by Overdijk<sup>[3]</sup>, which is based on Lo-  
112 cLinkVis by Olieman et al.<sup>[1]</sup>. Using the named entity recognition features of  
113 spaCy was also considered, but because of superior accuracy results reported  
114 by Overdijk<sup>[3]</sup>, his model was chosen.

115  
116 The model first tokenizes the input text using spaCy, a natural language pro-  
117 cessing toolkit available for Dutch language processing. Then, it detects munic-  
118 ipalities mentioned in the tokenized text and searches for streets that are within  
119 the scope of the detected municipalities. There are some drawbacks to this  
120 method, such as when there is no municipality detected, other toponyms will  
121 also not be detected. This poses a problem for articles that contain locations  
122 from outside of the Netherlands. Also, the model does not recognize landmarks  
123 and other high-detail geographical features. However, the model does detect  
124 highways and secondary roads (N-roads) that are not tied to a specific munic-  
125 ipality. The model also classifies the toponyms into several categories, but these  
126 labels are not used in this research.

### 127 3.3 Data selection

128 After toponym detection, all sentences that contain no toponyms are deleted.  
129 In the sentences that do contain a toponym, all words further than  $n$  words  
130 away from the closest toponym are also deleted. The value of  $n$  will be decided  
131 later in section Model Tuning. A toponym with its surrounding words will from  
132 now on be referred to as a *span*. In a span, there can be more than one toponym

133 if there was more than one toponym in its source sentence. [WINDOW SIZE  
134 NOEMEN]

### 135 3.4 Spatial identifier selection

136 When dealing with complex locations, the key to accurate geolocation lies in  
137 understanding the relation between the toponyms specifying the location. This  
138 relation is expressed in spatial identifier words surrounding the toponyms. To  
139 obtain a better understanding of what words are important for spatial identi-  
140 fication, a part of the training set was manually classified as either containing  
141 complex locations or containing non-complex locations. Subsequently, an anal-  
142 ysis was done to visualize the difference in word frequencies between all articles  
143 and complex-marked articles.

144  
145 The complex/non-complex annotated dataset contains the toponym spans of  
146 the first two hundred articles of the training set. Eighty-two articles are marked  
147 as containing complex locations, the rest are marked non-complex. Because  
148 spatial identifier words are never behind a toponym, the words located after the  
149 last toponym in a span are not counted in this analysis. The word count in both  
categories is normalized by the amount of words in each respective category.

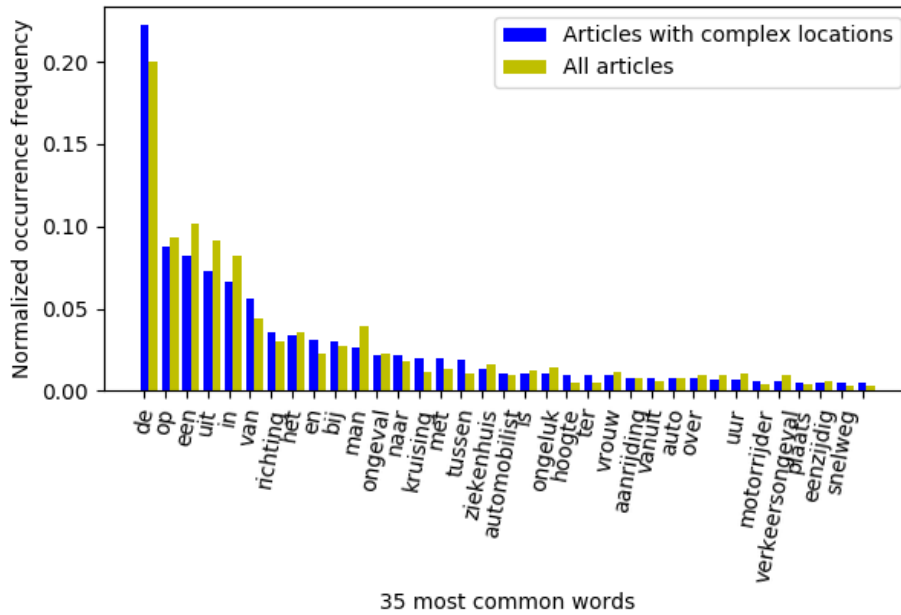


Figure 1: Frequency of thirty-five most common words in complex location articles, together with frequency of said word in all of first two hundred articles.

150  
151 To get a sense of how accurately this subset of two hundred articles represents  
152 the entire training set, Figure 2 was plotted. As one can see, the occurrence  
153 frequencies of the most common thirty-five words in both datasets correspond  
154 quite well.

155  
156 Figure 1 shows what words are most frequently occurring in articles containing

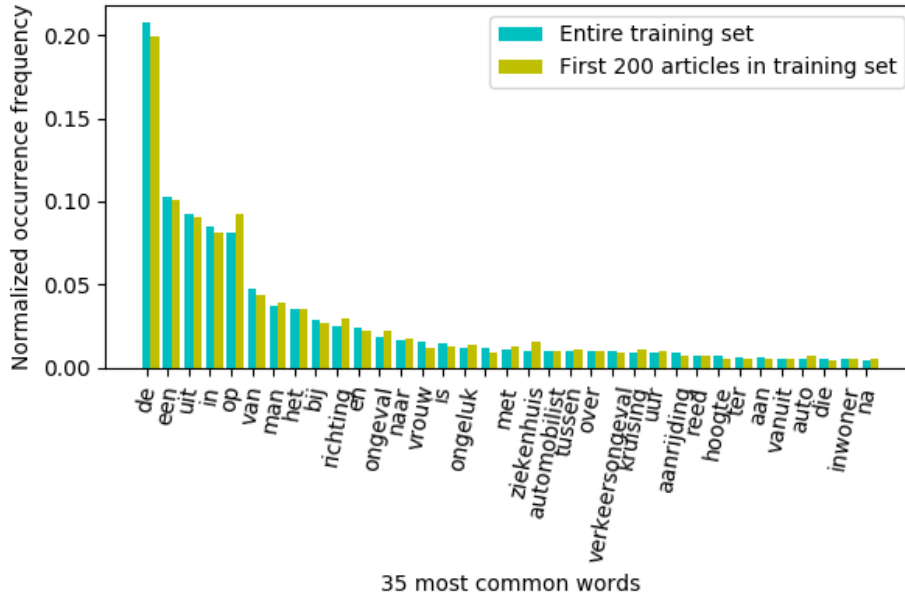


Figure 2: Frequency of thirty-five most common words in the whole training set, together with frequency of said word in first two hundred articles.

complex locations, but it also indicates how often these words are actually occurring in the general training set.

Drawing from Figure 1, word occurrence frequencies are mostly the same between complex locations and average articles, with some minor exceptions. When only looking at words that contain geographical information, there are some words with bigger frequency differences, such as “kruising”, “tussen” and “hoogte”. Even though these words do not appear in the training set very often, they appear to be important to complex locations. Of course, there are also spatial identifiers occurring with similar frequency in both categories that are still very important to understand, such as “op”, “uit”, “in”, “richting”, “bij” and “naar”.

### 3.5 Parsing spatial identifiers to predicates

Now that the most common spatial identifiers are selected, they can be parsed into specific rules to make geocoding possible. In the table below, the resulting predicates can be found on the left-hand side, with the selected input words on the right-hand side. The number behind the predicates determines how many locations the predicate takes. Some input words, such as “kruispunt” and “te”, are not among the most common thirty-five words, but can easily be added to an existing predicate since they have the same meaning. Sometimes, there occurs more than one of the selected prepositions within the window size before the toponym. In this case, the preposition closest to the toponym is selected.



Predicate	Spatial identifier
ON(1)	op over
IN(1)	in te
AT(1)	bij hoogte
BETWEEN(2)	tussen
INTERSECT(1/2)	kruising kruispunt splitsing

### 181 3.6 Converting toponym-predicate pairs into coordinates

182 After converting spatial identifier words to predicates, each toponym is con-  
 183 verted to coordinates in a way that depends on the predicate. This step gives  
 184 two types of results: a list of coordinate points from road toponyms, processed  
 185 by “ON(1)” or “INTERSECT(1/2)”, or a simple two- or four-coordinate bound-  
 186 ing box from other toponyms, processed by the other predicates. Both resulting  
 187 types have their own result list where they are added to. These lists are later  
 188 combined, this process is described in section 3.7.

#### 189 3.6.1 Coordinate source & search method

190 To retrieve the coordinates associated with a toponym, data is extracted from  
 191 OpenStreetMap. This is done using Nominatim and Overpass API. Nominatim  
 192 is used to search for cities and villages with a clear center and surrounding  
 193 bounding box, whereas the Overpass API is used to search for streets. Nominatim  
 194 is the preferred search method due to its speed and ease-of-use, but it shows  
 195 poor recall of highway segments and other long road segments. Both search  
 196 methods are programmed to primarily search for coordinates in the Nether-  
 197 lands, with Overpass excluding any results abroad, and Nominatim preferring  
 198 results in the Netherlands.

#### 199 3.6.2 ON(1)

200 Through observation, the toponym following the “ON” predicate has been deter-  
 201 mined to mostly be a street of some kind. Therefore, the toponym coordinates  
 202 are searched for through Overpass. No bounding box is given other than the  
 203 Netherlands, so search results usually include multiple streets throughout the  
 204 country if the street name happens to be a commonly used one. The toponym  
 205 and result coordinates are added to the list of road search results.

206 **3.6.3 IN(1)**

207 Toponyms following “IN” are largely cities, villages, landmarks and public build-  
208 ings. Toponyms are thus searched for using Nominatim. Search results of the  
209 class “boundary” or “place” are preferred, because both results contain accurate  
210 bounding box of a city or village. If these classes are not found, the first search  
211 result is returned. The resulting coordinates are added to the list containing all  
212 found bounding boxes in that specific article.

213 **3.6.4 AT(1)**

214 **3.6.5 BETWEEN(2)**

215 Toponyms occurring after “BETWEEN” can be of any kind, though most are  
216 cities and villages. Both toponyms are looked up through Nominatim, with the  
217 result of class “place” given priority. In case there is no result of class “place”,  
218 the first result is returned.

219  
220 The bounding box between these two toponyms takes the shape of a square,  
221 where both toponyms are at opposite ends of a diagonal in the square. The  
222 other two points forming the bounding box are located at an angle of forty-five  
223 degrees from both sides of the diagonal, according to the properties of a square.  
224 To determine the distance between either point and either toponym, the length  
225 of the diagonal is divided by the square root of two.

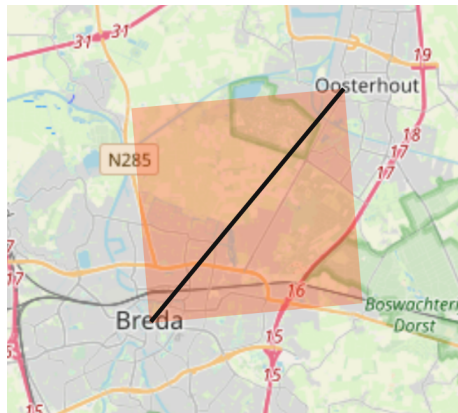


Figure 3: Example of bounding box created by “BETWEEN” predicate. Red zone is bounding box, black line is diagonal between toponyms.

226 **3.6.6 INTERSECT(1/2)**

227 Toponyms occurring in this predicate are assumed to be streets. This predicate  
228 sometimes occurs with only one toponym. In this case, the other toponym that  
229 makes up the intersection is usually located somewhere else in the sentence or  
230 in a neighbouring sentence. Therefore, “INTERSECT(1)” is not processed, but  
231 immediately added to the coordinate point result list. This predicate will be  
232 dealt with later on, when the two resulting lists are merged.

233  
234 In the case of “INTERSECT(2)”, both toponyms are looked up through the  
235 Overpass API. The Overpass API searches for nodes where both toponyms overlap  
236 and it returns the coordinates of all nodes where this occurs. The resulting  
237 coordinates are added to the list of road search results.

### 238 3.7 Predicate result merging

239 Because street toponym predicates retrieve coordinates for all streets in the  
240 Netherlands with the toponym name, these need to be filtered in order to be  
241 of any use. That is why “INTERSECT(2)” and “ON(1)” predicate results are  
242 matched against any bounding box in the bounding box result list. If any of  
243 the result point coordinates fall within a bounding box, all coordinates that lie  
244 outside of the bounding box are removed from the coordinate point list, and the  
245 bounding box itself is also removed. This is how street coordinate results are  
246 narrowed down.

247  
248 The still unprocessed “INTERSECT(1)” predicate-toponym pair is matched  
249 using “INTERSECT(2)”, together with another toponym coming from other  
250 “INTERSECT(1)” or “ON(1)” predicates that are in the street result list. If  
251 any matching nodes are returned, both input toponyms and their predicates are  
252 deleted from the street result list.

253  
254 Any entries in both result lists that do not match with any other entry remain  
255 in the list. The merging method described above is used after processing  
256 all predicate-toponym pairs in a sentence. Afterwards, all remaining entries in  
257 both lists are added to two similar lists that contain locations from the entire  
258 article. When all sentences in the article are processed, the merging process is  
259 run again these two lists. The reason for performing both sentence-based and  
260 article-based merging instead of only article-based merging, is that toponyms  
261 in the same sentence are more likely to be related than toponyms in the article.

### 262 3.8 Resulting output

263 The model outputs two lists per article, with one list containing bounding boxes,  
264 and the other list containing point coordinates to signify streets and other points.  
265 both

## 4 Results

### 4.1 Evaluation method

To evaluate the model described in section 3, the first eighty articles of the test set were annotated. Every toponym in the article was manually tagged, including the spatial identifier in front of each toponym. Also, every location formed by the toponyms was annotated. Each location was marked by its coordinates, the toponym that describes the location most specifically, and a precision radius in meters. All locations consisting of toponyms with spatial identifiers outside the set selected for this research were ignored, since the model does not use these toponyms. The eighty articles contain a total of eighty-seven annotated locations. To isolate the model from imperfect recall during toponym detection, the model will also be evaluated using the annotated toponyms with spatial identifiers as input.

To calculate the accuracy the model, the number of correct geolocations is divided by the number of total geolocations. To evaluate a performed geolocation by the model, the average coordinate is calculated for each bounding box or road coordinate list. If this average coordinate is within radius distance of the annotated coordinate, the geolocation is marked correct. If it is outside radius distance, the geolocation is marked incorrect.

	Accuracy
Full model	0.33
Annotated toponym input	0.53

## 287 **5 Discussion**

288 - Can this study be generalized to all Dutch text, although the dataset contains  
289 only traffic accident articles? - Some roads and other toponyms may have  
290 changed since the publication date of the article. - Center of a bounding box of  
291 a village is sometimes far away of actual center of village. Possibly outside of  
292 radius range.

293 **6 Conclusion**

## 294 References

- 295 [1] A. Olieman, J. Kamps, and R. M. Claros, “Loclinkvis: A geographic in-  
296 formation retrieval-based system for large-scale exploratory search,” *arXiv*  
297 *preprint arXiv:1509.02010*, 2015.
- 298 [2] B. Hendriks, “Retrieving, cleaning and analysing dutch news articles about  
299 traffic accidents,” 2019.
- 300 [3] C. Overdijk, “Location recognition and geocoding in dutch traffic news ar-  
301 ticles,” 2020.