

1

2

3

4

5

6

7

8

9

10

YOUR TITLE

SUBMITTED IN PARTIAL FULFILLMENT FOR THE DEGREE OF

MASTER OF SCIENCE

AXEL HIRSCHHEL

10656146

MASTER INFORMATION STUDIES

DATA SCIENCE

FACULTY OF SCIENCE

UNIVERSITY OF AMSTERDAM

YOUR DATE OF DEFENCE IN THE FORMAT YYYY-MM-DD

11

	Internal Supervisor	External Supervisor
Title, Name	Dr Maarten Marx	Tom Kunzler
Affiliation	UvA, FNWI, IvI	Open State Foudation
Email	maartenmarx@uva.nl	tom@openstate.eu

14 **Todo list**

15 **Contents**

16	Todo list	1
17	1 Introduction	3
18	2 Related Work	5
19	2.1 RQ1	5
20	2.2 RQ2	5
21	3 Methodology	6
22	3.1 Description of the data	6
23	3.2 Wat plotjes en tabelletjes	7
24	3.3 Methods	9
25	3.3.1 RQ1	9
26	3.3.2 RQ2	9
27	4 Evaluation	10
28	5 Conclusions	11
29	5.1 Acknowledgements	11
30	References	12
31	A Slides	12

32

Abstract

33

34 Thesis requirements

- 35 • Your thesis is written in ACM style with two columns (`documentclass[sigconf]acmart`).
- 36 • It is maximally 10 pages long, excluding the title page and the appendix,
37 but including references, figures, etc

1 Introduction

Text classification is considered as one of the most important challenges within natural language processing. Classifying documents is vital, as it enables users to easily query and retrieve useful information. Moreover, it allows the automatization of many processes such as spam classification and sentiment analysis. Given the wide variety of application, many algorithms have been developed to tackle the problem (Aggarwal & Zhai, 2012).

Bayesian Classifiers are considered a class of classification algorithms, and they classify document based on word occurrences within documents. This word presence is used to calculate the probability that certain documents are part of a topic. The two prominent versions of bayesian classifiers are multi-variate Bernoulli models and multinomial models. Another widely used class of text classifiers are support vector machines (SVM). Within SVM the algorithm creates linear hyperplanes which split the data into classes based on a bag-of-words representation of texts. Using kernel tricks hyperplanes can be constructed which can find more complex relations than linear (Aggarwal & Zhai, 2012).

Recently, deep neural networks have been employed on classification problems as well. Most notably convolutional neural nets (CNN) have outperformed other methods on baseline classification problems. CNN have been generally been used on image data, but research in word and document embedding spaces such as Word2Vec (Mikolov, Chen, Corrado, & Dean, 2013) and Paragraph2Vec (Le & Mikolov, 2014) allow the use of CNN on text as well. Transforming words within documents into a multi-dimensional vectors allows the use of convolutional filters, which shift over the documents and detect patterns within documents (Kim, 2014).

In contrast to many of the baseline challenges within text classification, real-world application of classification often involves other challenges as well. Within this research documents of Dutch municipalities are classified, which is a difficult task due to three properties. Firstly, no labelled training data is available, which means that training needs to be done on data from the central document. Secondly, many of the documents are multi-topic and it is interesting to discuss how well algorithms deal with this. Thirdly, within classes a large variety of documents exist, as the documents differ in length and style. The research question and subquestions are thus:

- How well does CNN perform on classifying Dutch municipality documents compared to SVM and Bayesian Classifiers?
 - How does the detail of topics influence the performances of all algorithms?
 - What thresholds are optimal in detecting topics of multi-labeled documents.
 - How well do the performances of the algorithms on the dataset of the central government generalize to the dataset of the municipalities?
 - How can all algorithms be optimized to deal with the large in-class variety of the topics?

Within the next chapter, the literature review, current approaches to the mentioned challenges and the general idea behind the algorithms is further ex-

85 plained. Then the specific set-up for this research is discussed within the
86 methodology, which also provides information on how the research questions
87 are answered. The results of this research is described next, and provides a
88 detailed overview of the performance of all algorithms with various evaluation
89 metrics. Lastly, the answers to the research questions are formulated and the
90 conclusions from this article are discussed.

91 **2 Related Work**

92 Deze sectie bestaat uit een aantal "blokken", waarin je per blok de relevante
93 literatuur beschrijft.

94 Neem alleen literatuur op die van belang is voor jouw onderzoeksvraag en
95 deelvragen.

96 Typisch heb je 1 blok voor je hoofdvraag en per deelvraag **RQi** een blok.

97 **2.1 RQ1**

98 **2.2 RQ2**

99 **3 Methodology**

100 **3.1 Description of the data**

101 Data verzameling en beschrijving van de data

102 Hoe is de data verzameld, en hoe heb jij die data verkregen?

103 Wat staat er in de data? Niet alleen maar een technisch verhaal, maar ook
104 inhoudelijk. DE lezer moet een goed idee krijgen over de technische inhoud en
105 wat het betekent.

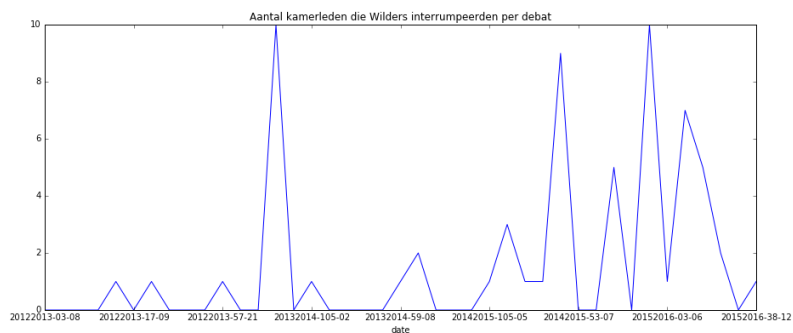


Figure 1: Aantal interrupties van Wilders in de Tweede Kamer door de tijd (periode 2012-2016).

3.2 Wat plotjes en tabelletjes

Zie het IPython Notebook `PandasAndLatex.ipynb` voor de code om vanuit pandas een poltje op te slaan en een dataframe als tabel op te slaan. Het werkt ideaal!

De interrupties van Wilders staan beschreven in Figure 1 en Tabel 1.

date	indegree	interruptie_volgorde
20122013-03-08	0.0	
20122013-07-16	0.0	
20122013-100-03	0.0	
20122013-100-06	0.0	
20122013-17-06	1.0	Pechtold-3
20122013-17-09	0.0	
20122013-21-04	1.0	Pechtold-3
20122013-22-08	0.0	
20122013-32-06	0.0	
20122013-48-23	0.0	
20122013-57-21	1.0	Pechtold-6
20122013-76-03	0.0	
20122013-76-06	0.0	
20132014-05-02	10.0	Roemer-4 Van Haersma Buma-4 Pechtold-4 Slob-5 ...
20132014-06-04	0.0	
20132014-105-02	1.0	Pechtold-10
20132014-105-06	0.0	
20132014-14-03	0.0	
20132014-14-06	0.0	
20132014-52-18	0.0	
20132014-59-08	1.0	Klaver-3
20142015-02-08	2.0	Pechtold-6 Slob-4
20142015-03-06	0.0	
20142015-09-09	0.0	
20142015-100-05	0.0	
20142015-105-05	1.0	Pechtold-2
20142015-111-04	3.0	Pechtold-6 Kuzu-8 Klaver-3
20142015-111-07	1.0	Pechtold-2
20142015-39-71	1.0	Pechtold-2
20142015-41-07	9.0	Samsom-2 Pechtold-3 Kuzu-6 Zijlstra-5 Van Ojik...
20142015-53-07	0.0	
20142015-61-23	0.0	
20142015-79-07	5.0	Klaver-10 Gesthuizen-3 Voordewind-2 Pechtold-6...
20142015-95-06	0.0	
20152016-02-07	10.0	Pechtold-5 Slob-7 Klaver-11 Kuzu-24 Öztürk-1 S...
20152016-03-06	1.0	Pechtold-5
20152016-14-02	7.0	Klaver-9 Roemer-4 Samsom-2 Van Haersma Buma-5 ...
20152016-14-05	5.0	Van Haersma Buma-13 Pechtold-4 Zijlstra-1 Klav...
20152016-27-03	2.0	Segers-4 Kuzu-10
20152016-38-10	0.0	
20152016-38-12	1.0	Klein-2

Table 1: Door wie werd Wilders onderbroken en hoe vaak per debat.

111 **3.3 Methods**

112 Hoe je je vraag gaat beantwoorden.

113 Dit is de langste sectie van je scriptie.

114 Als iets erg technisch wordt kan je een deel naar de Appendix verplaatsen.

115 Probeer er een lopend verhaal van te maken.

116 Het is heel handig dit ook weer op te delen nav je deelvragen:

117 **3.3.1 RQ1**

118 **3.3.2 RQ2**

119 4 Evaluation

120 Met een subsectie voor elke deelvraag.

121 In hoeverre is je vraag beantwoord?

122 Een mooie graphic/visualisatie is hier heel gewenst.

123 Hou het kort maar krachtig.

124 **5 Conclusions**

125 Hierin beantwoord je jouw hoofdvraag op basis van het eerder vergaarde bewijs.

126 **5.1 Acknowledgements**

127 Hier kan je bedanken wie je maar wilt.

128 **References**

- 129 Aggarwal, C. C., & Zhai, C. (2012). A survey of text classification algorithms.
130 In *Mining text data* (pp. 163–222). Springer.
- 131 Kim, Y. (2014). Convolutional neural networks for sentence classification. *arXiv*
132 *preprint arXiv:1408.5882*.
- 133 Le, Q., & Mikolov, T. (2014). Distributed representations of sentences and
134 documents. In *International conference on machine learning* (pp. 1188–
135 1196).
- 136 Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013). Efficient estimation of
137 word representations in vector space. *arXiv preprint arXiv:1301.3781*.

138 **A Slides**