

1                   METHODEN VOOR HET VOORSPELLEN VAN  
2                   PARTIJ-AFFILIATIE IN DE TWEEDE KAMER  
  
3                   SUBMITTED IN PARTIAL FULFILLMENT FOR THE DEGREE OF  
4                   BACHELOR OF SCIENCE  
  
5                   JASPER VAN DER HEIDE  
6                   10732721  
  
7                   BACHELOR INFORMATIEKUNDE  
8                   FACULTY OF SCIENCE  
9                   UNIVERSITY OF AMSTERDAM  
10                  YOUR DATE OF DEFENCE IN THE FORMAT YYYY-MM-DD

11

	<b>Internal Supervisor</b>	<b>Second Supervisor</b>
<b>Title, Name</b>	Dr Maarten Marx	
<b>Affiliation</b>	UvA, FNWI, IvI	
<b>Email</b>	maartenmarx@uva.nl .	



13	<b>Inhoudsopgave</b>	
14	<b>1 Introductie</b>	<b>3</b>
15	<b>2 Gerelateerd werk</b>	<b>3</b>
16	2.1 Classificatiemethoden in gerelateerde werken . . . . .	4
17	2.2 Invloed van oppositie of regering . . . . .	4
18	<b>3 Methodologie</b>	<b>5</b>
19	3.1 De data . . . . .	5
20	3.2 Methoden . . . . .	6
21	3.2.1 Deelvraag 1 . . . . .	6
22	3.2.2 Deelvraag 2 . . . . .	6
23	3.2.3 Deelvraag 3 . . . . .	7
24	<b>4 Evaluatie</b>	<b>7</b>
25	<b>5 Conclusies</b>	<b>8</b>
26	<b>A Slides</b>	<b>9</b>

27

## Samenvatting

28

## 1 Introductie

Teksten van politieke partijen kunnen bruikbaar zijn voor het bepalen van ideologische positie van andere teksten, aangezien zij zowel een tekst leveren als ook een bekende ideologie. Deze informatie kan vervolgens toegepast worden bij andere teksten die wellicht ideologisch van aard zijn. Bijvoorbeeld, aan de hand van deze informatie kan men teksten uit kranten classificeren op basis van ideologie[1, 2].

In diverse landen zijn al verschillende onderzoeken gedaan naar het classificeren van partij-affiliatie op basis van teksten van politici[3, 1]. Mede omdat elk land een andere politiek stelsel en cultuur heeft, verschillen de resultaten. Daarnaast gebruikt elk onderzoek ook een andere methode voor het classificeren.

Een onderzoek gericht op het Nederlandse parlement ontbreekt hierbij nog.

Dit onderzoek richt zich daarom op een breder scala aan mogelijke methoden en daarnaast ook specifiek op de Nederlandse politiek. De onderzoeksvraag luidt daarom dus ook: "In hoeverre is classificatie op basis van partij-affiliatie aan de hand van spreekbeurten in de Tweede Kamer het gevolg van ideologie?"

Deze vraag wordt beantwoord door de antwoorden te vinden op de volgende deelvragen:

1. Wat is het beste classificatiemodel voor classificatie van partij-affiliatie in de Tweede Kamer en wat is het resultaat van dit model?
2. In hoeverre is deze classificatie afhankelijk van partijnamen en namen van Kamerleden?
3. In hoeverre wordt deze classificatie bepaald door partij-status (d.w.z. oppositie of regering)?
4. In hoeverre is deze classificatie afhankelijk van of een partij rechts of links is?

**Overzicht van scriptie** In sectie 2 zal gerelateerd werk besproken worden, met name vergelijkbare onderzoeken uit andere landen. In sectie 3 zal de methodologie van de verschillende deelvragen behandeld worden. In sectie 4 zullen vervolgens de resultaten weergegeven worden. In sectie 5 zal een evaluatie plaatsvinden van zowel de resultaten als de gehanteerde methodologie. In sectie 6 wordt ten slotte de onderzoeksvraag beantwoord.

## 2 Gerelateerd werk

Diermeier et al. deden onderzoek naar het classificeren op basis van ideologische positie in de Amerikaanse Senaat[4]. Ze trainden hun classificatie op de speeches van de 25 meest liberale en de 25 meest conservatieve senatoren van het 101e tot en met het 107e congres en testten op de 25 meest liberale en de 25 meest conservatieve senatoren van het 108e congres. Later in het onderzoek vergeleken ze ook de 25 gematigd conservatieve en de 25 gematigd liberale senatoren. Dit onderzoek wist de ideologie van de senatoren te voorspellen met een 94 procent nauwkeurigheid voor de classificatie van de extremen, maar slechts een 52 procent nauwkeurigheid voor de classificatie van de gematigde senatoren.

Als een vervolg op dit onderzoek deden Graeme Hirst et al. een vergelijkbaar onderzoek naar zowel het Canadese Parlement [2]. Hierbij werd zowel gekeken naar de Engelse als Franse teksten. Afhankelijk van taal en dataset, vinden zij in dit onderzoek nauwkeurigheden van 83.2 procent en hoger. In tegenstelling tot het onderzoek van Diermeier et al., vinden zij minder dat de woorden van de sprekers een uiting zijn van ideologie.

Ferreira probeerde interventies van politici te classificeren op basis van geslacht, leeftijdsgroep, partij-affiliatie en oriëntatie in het Portugese parlement [3]. In alle classificaties kon men een F1-score van 0.87 of hoger bereiken.

## 2.1 Classificatiemethoden in gerelateerde werken

In het onderzoek van Diermeier et al. werd gebruik gemaakt van support vector machines. Verder maakten ze gebruik van *tf-idf* met een minimale woordfrequentie van 50 en een documentfrequentie van 10, *Part-Of-Speech tagging* en werden alle eennamen verwijderd.

In het onderzoek van Graeme Hirst et al. maakten ze gebruik van support vector [2]. Ze experimenteerden met verschillende vormen van pre-processing, inclusief stemmen en het verwijderen van woorden op basis van te hoge of te lage frequentie. Deze variaties maakten in hun onderzoek geen grote verschillen en uiteindelijk is gekozen voor het niet stemmen, het weglaten van woorden die in minder dan vijf documenten voorkomen en resultaten van zowel met als zonder de top 500 meest frequente woorden. Daarnaast werd geëxperimenteerd met vier wegen voor woorden: *boolean*, *tf*, *tf-norm* en *tf-idf*, waarvan *tf-idf* het beste resultaat opleverde.

In het onderzoek van Ferreira werd gebruik gemaakt van twee classificatiemethoden: Logistische regressie en MIRA[3]. Logistische regressie werd aangevuld met *group Lasso* regularisatie. Voor wegen van woorden werd geëxperimenteerd met *tf*, *tf-idf*,  $\Delta$ -*tf-idf* en  $\Delta$ -*BM-25*. Daarnaast wordt er gebruik gemaakt van woordclustering, *Concise Semantic Analysis* en stylometrische eigenschappen. Op *Part-Of-Speech tagging* na hadden stylometrische eigenschappen een duidelijke negatieve invloed op de classificatie.

## 2.2 Invloed van oppositie of regering

Graeme Hirst et al. vonden in hun onderzoek dat de classificatie van spreker in het Canadese parlement op basis van partij-affiliatie meer zegt over de status van de partij (regering of oppositie).[2] Zo vergeleken zij de top tien karakteristieke woorden van de liberalen en conservatieven in het 36e parlement (liberalen in de regering) en het 39e parlement (conservatieven in de regering). Hier vonden zij dat vier van de tien woorden van de liberalen (regering) in het 36e parlement bij het 39e parlement bij de conservatieven (regering) te vinden waren. Andersom gebeurde hetzelfde met één van de tien woorden van de conservatieven (oppositie) in het 36e parlement naar liberalen (oppositie) in het 39e parlement.

In hetzelfde onderzoek trainden ze ook hun classifiers op het ene parlement en testten deze op het andere parlement. Hierbij vonden zij in beide gevallen een nauwkeurigheid ver onder de baseline. Daarnaast deden ze ook nog een classificatie op de sprekers die in beide parlementen zaten en een andere classificatie op sprekers die niet in beide parlementen zaten. Bij de eerste

116 classificatie vonden ze nauwkeurigheden rond de baseline, terwijl in de tweede  
117 situatie nauwkeurigheden gevonden werden ver boven de baseline.

118 Deze resultaten leidden de onderzoekers tot de conclusie dat de classificatie  
119 voornamelijk het gevolg is van de status van de partij en minder van ideologie.

## 120 3 Methodologie

### 121 3.1 De data

122 De data die gebruikt worden, zijn de Handelingen van de Tweede Kamer gedu-  
123 rende het missionaire kabinet-Rutte II (5 november 2012 tot 22 maart 2017).  
124 Deze data is in xml-formaat van de website officiële bekendmakingen.nl gehaald,  
125 samen met corresponderende metadata xml-bestanden. De bestanden van de  
126 Handelingen bevatten voornamelijk informatie over spreekbeurten tijdens een  
127 debat, waaronder naam van een spreker, partij-affiliatie, inhoud van de spreek-  
128 beurt en het soort spreekbeurt. Deze gegevens zijn samengevoegd tot een tabel  
129 en opgeslagen als csv-bestand.

130 UITLEGGEN WELKE SPREEKBEURTEN EXACT GEKOZEN ZIJN  
131 en WAT EEN DOCUMENT IS

132 Deze dataset bevat naast de verkozen partijen van de 2012 Tweede Kamerver-  
133 kiezingen, ook afsplitsingen van die partijen (tien in totaal) en bezoeken van  
134 vertegenwoordigingen van die partijen uit de Eerste Kamer (tien in totaal).  
135 Omdat van beide categoriën relatief weinig data is en er overlap zit met hun  
136 oorspronkelijke partij, zijn deze er uit gehaald.

Tabel 1: Aantal spreekbeurten per partij gedurende het missionaire kabinet-Rutte II.

50PLUS	413
CDA	2216
ChristenUnie	1223
D66	2211
GroenLinks	1193
PVV	1880
PvdA	2269
PvdD	480
SGP	770
SP	2573
VVD	2157

## 137 3.2 Methoden

### 138 3.2.1 Deelvraag 1

139 Om deze deelvraag te beantwoorden zullen een aantal classificatiemethoden ver-  
140 geleken worden. Aangezien het onmogelijk om alle classificatiemethoden te ver-  
141 geleken, beperkt dit onderzoek zich tot classificatiemethoden die goede resul-  
142 taten hebben opgeleverd in andere onderzoeken, genoemd in 2.1. Hieronder  
143 worden de verschillende onderdelen besproken.

144 **Pre-processing** Voor pre-processing is gebruik gemaakt van tokenisation, lo-  
145 wercasing en stemming. Voor tokenisation is de reguliere expressie  
146  $w+$  gebruikt, die daarmee alleen de letters van het alfabet overhoudt. Deze  
147 woorden zijn vervolgens allemaal omgezet in kleine letters. Vervolgens is er  
148 gevarieerd tussen wel of geen gebruik maken van stemming. In het geval van  
149 stemming is gebruik gemaakt van de Snowball Stemmer via de Python NLTK  
150 module.

151 **Bag-of-words model** Bag-of-words model is de meest gebruikte representa-  
152 tie van data in vergelijkbare onderzoeken. Bij het bag-of-words model wordt  
153 elk document gerepresenteerd door een vector, waarbij elke kolom een woord  
154 voorstelt met een bijbehorende waarde. Voornaamste beperking van dit model  
155 is dat het geen rekening houdt met de volgorde van woorden, wat een groot  
156 effect kan hebben op de betekenis van een document.

157 Voor dit onderzoek zijn de volgende wegen voor woorden getest: *boolean*  
158 (wel of niet aanwezig), *tf* (woordfrequentie), *tf-norm* (woordfrequentie genor-  
159 maliseerd door documentlengte) en *tf-idf* [2] Daarnaast wordt in dit onderzoek  
160 geëxperimenteerd met een minimale of maximale woord- of documentfrequentie.

161 **Support-Vector Machines** Machine Learning Algorithms

- 162 1. Support-Vector Machines
- 163 2. Logistische Regressie
- 164 3. Naive Bayes

165 Voor de classificatiemethoden wordt waar mogelijk gebruik gemaakt van  
166 functies van de Python module scikit-learn[5], aangevuld met zelf geschreven  
167 code als dit niet reeds beschikbaar is. Bij al deze classificatiemethoden wordt  
168 gevarieerd met meerdere parameters door middel van een gridsearch. Hierbij  
169 wordt gebruikt gemaakt van 5-fold cross-validation. De uitslagen worden be-  
170 oordeeld op basis van gewogen f1-scores.

### 171 3.2.2 Deelvraag 2

172 Voor deze deelvraag wordt wederom een classificatie gedaan met de classifica-  
173 tiemethode die resulteerde uit deelvraag 1. In deze classificatie worden alle  
174 partijnamen vervangen door de tag PARTIJNAAM en alle namen van Kamerle-  
175 den vervangen door de KAMERLIDNAAM. Deze resultaten worden vervolgens  
176 vergeleken met de resultaten uit deelvraag 1.

### 177 3.2.3 Deelvraag 3

178 Om deze deelvraag te beantwoorden zullen de drie experimenten die Graeme  
179 Hirst et al. uitvoerden voor dezelfde vraag gereproduceerd worden op de dataset  
180 van de Tweede Kamer. Bij deze deelvraag zal de beste classifier uit deelvraag  
181 1 gebruikt worden.

182 Als vergelijkingsmateriaal is voor deze experiment een tweede dataset no-  
183 dig uit een ander kabinet. Hiervoor is het wenselijk dat dit kabinet bestaat  
184 uit andere partijen dan kabinet-Rutte II. Er moet voor het derde experiment  
185 variatie zijn in de Kamerleden tussen de twee kabinetten, maar ook voldoende  
186 Kamerleden die in beide perioden in de kamer zaten. Daarnaast is het ook  
187 wenselijk als het niet te ver terug is, zodat onderwerpen en taalgebruik enig-  
188 zins overeenkomstig zijn. Omdat kabinet-Rutte I een minderheidskabinet was  
189 met een bijzondere partij-status voor de PVV, is ervoor gekozen om de Tweede  
190 Kamer tijdens het missionaire kabinet-Balkenende IV (22 februari 2007 tot 20  
191 februari 2010) te gebruiken.

192 In het eerste experiment zullen de tien meest karakteristieke woorden per  
193 partij van het ene parlement vergeleken worden met de tien meest karakteristieke  
194 woorden per partij van het andere parlement. Als de classificatie op basis van  
195 ideologie is in plaats van partij-status, is het te verwachten dat de woorden bij  
196 een partij blijven en niet gekoppeld zijn aan in oppositie of regering zitten.

197 In het tweede experiment worden classifiers getraind op het ene parlement  
198 en getest op het andere parlement. Als de classificatie op basis van ideologie  
199 is in plaats van partij-status, is de verwachting dat er nog steeds aanzienlijke  
200 voorspellingen gedaan worden, aangezien de ideologie naar verwachting redelijk  
201 stabiel is binnen tien jaar (hoewel woordgebruik varieert). Als de scores aan-  
202 zienlijk lager zijn, kan dit het gevolg zijn van het veranderen van partij-status  
203 van partijen.

204 In het derde experiment zullen twee classificaties vergeleken worden. De  
205 eerste op Kamerleden die in beide parlementen zaten en een classificatie op  
206 Kamerleden die maar in één van de twee parlementen hebben gezeten.

## 207 4 Evaluatie

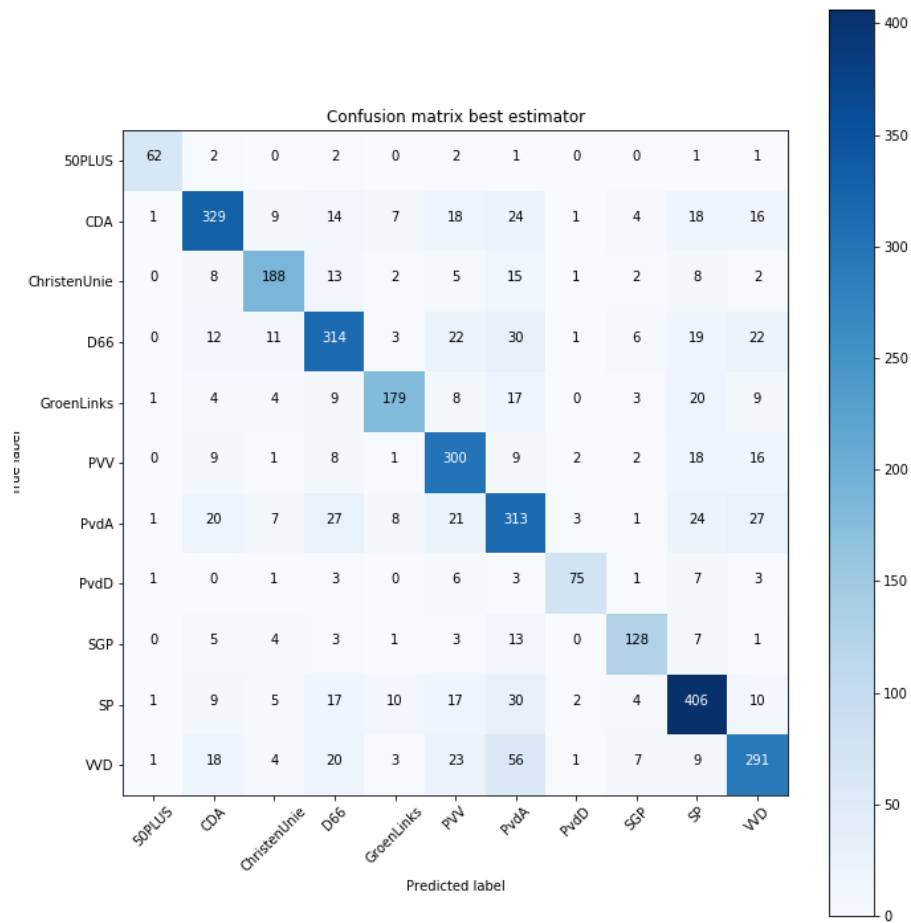
208 Met een subsectie voor elke deelvraag.

209 In hoeverre is je vraag beantwoord?

210 Een mooie graphic/visualisatie is hier heel gewenst.

211 Hou het kort maar krachtig.





212

213

## 214 5 Conclusies

215 Hierin beantwoord je jouw hoofdvraag op basis van het eerder vergaarde be-  
 216 wijs.

## 217 Referenties

- 218 [1] Felix Bießmann. Automating political bias prediction. *CoRR*,  
 219 abs/1608.02195, 2016.
- 220 [2] Graeme Hirst, Yaroslav Riabinin, Jory Graham, and Magali Boizot-roche.  
 221 Text to ideology or text to party status? \*.
- 222 [3] Vasco Ferreira. Using textual transcripts of parliamentary interventions for  
 223 profiling portuguese politicians. 2016.

- 224 [4] Daniel Diermeier, Jean-François Godbout, Bei Yu, and Stefan Kaufmann.  
225 Language and ideology in congress. *British Journal of Political Science*,  
226 42(1):31–55, 2012.
- 227 [5] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel,  
228 M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Pas-  
229 sos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-  
230 learn: Machine learning in Python. *Journal of Machine Learning Research*,  
231 12:2825–2830, 2011.
- 232 [6] Bei Yu, Stefan Kaufmann, and Daniel Diermeier. Classifying party affiliation  
233 from political speech. *Journal of Information Technology & Politics*, 5(1):33–  
234 48, 2008.

## 235 A Slides