

1 IDEOLOGIE EN CLASSIFICATIE IN DE HANDELINGEN
2 VAN DE TWEEDE KAMER
3
4 SUBMITTED IN PARTIAL FULFILLMENT FOR THE DEGREE OF
5 BACHELOR OF SCIENCE
6
7 JASPER VAN DER HEIDE
8 10732721
9
10 BACHELOR INFORMATIEKUNDE
11 FACULTEIT DER NATUURWETENSCHAPPEN, WISKUNDE EN
INFORMATICA
UNIVERSITEIT VAN AMSTERDAM
2018-06-28

	Begeleider	Tweede lezer
12 Titel, Naam	Dr Maarten Marx	
Affiliatie	UvA, FNWI, IvI	
Email	maartenmarx@uva.nl .	



14	Inhoudsopgave	
15	1 Introductie	3
16	2 Gerelateerd werk	3
17	2.1 Classificatiemethoden	4
18	2.2 Invloed van oppositie of regering	5
19	3 Methodologie	5
20	3.1 De data	5
21	3.2 Methoden	6
22	3.2.1 Deelvraag 1	6
23	3.2.2 Deelvraag 2	8
24	3.2.3 Deelvraag 3	8
25	3.2.4 Deelvraag 4	9
26	4 Evaluatie	10
27	4.1 Resultaten	10
28	4.1.1 Deelvraag 1	10
29	4.1.2 Deelvraag 2	12
30	4.2 Discussie	13
31	4.2.1 Deelvraag 1	13
32	4.2.2 Deelvraag 4	13
33	5 Conclusies	13
34	A Slides	14

35

Samenvatting

36

37 1 Introductie

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

44 In diverse landen zijn al verschillende onderzoeken gedaan naar het clas-
45 sificeren van partij-affiliatie op basis van teksten van politici[3, 1]. Mede omdat
46 elk land een ander stelsel en cultuur heeft, verschillen de resultaten. Elk onder-
47 zoek gebruikt ook een andere methode voor het classificeren. Daarnaast vinden
48 sommige onderzoeken dat deze classificatie minder het gevolg is van ideologie
49 maar meer van bijvoorbeeld regering tegenover oppositie.[2]

50 Een onderzoek gericht op het Nederlandse parlement ontbreekt hierbij nog.
51 Ook beperken veel onderzoeken zich vaak tot één classificatiemethode.

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

56 Deze vraag wordt beantwoord door de antwoorden te vinden op de vol-
57 gende deelvragen:

- 58 1. Wat is het beste classificatiemodel voor classificatie van partij-affiliatie in
59 de Tweede Kamer en wat is het resultaat van dit model?
- 60 2. In hoeverre is deze classificatie afhankelijk van partijnamen en namen van
61 Kamerleden?
- 62 3. In hoeverre wordt deze classificatie bepaald door partij-status (d.w.z. op-
63 positie of regering)?
- 64 4. In hoeverre wordt deze classificatie bepaald door links/rechts verdeling?

65 Daarom zal eerst gekeken worden naar classificatiemethoden en resultaten in
66 vergelijkbare onderzoeken. Van deze classificatiemethoden zullen een aantal
67 toegepast worden op teksten van de Tweede Kamer. Vervolgens zal door middel
68 van de overige deelvragen bepaald worden in hoeverre dit een reflectie is van
69 ideologie.

70 **Overzicht van scriptie** In sectie 2 zal gerelateerd werk besproken worden,
71 met name vergelijkbare onderzoeken in andere landen. In sectie 3 zal de me-
72 thodologie van de verschillende deelvragen behandeld worden. In sectie 4 zul-
73 len vervolgens de resultaten weergegeven worden. In sectie 5 zal een evaluatie
74 plaatsvinden van zowel de resultaten als de gehanteerde methodologie. In sectie
75 6 wordt ten slotte de onderzoeksvraag beantwoord.

76 2 Gerelateerd werk

77 Diermeier et al. deden onderzoek naar het classificeren op basis van ideologi-
78 sche 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. Een document werd hier gezien als de samenvoeging van alle spreekbeurten van een spreker. 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.

Het onderzoek van Bhand et al. richtte zich op het classificeren van leden van het Amerikaanse congres in 2005, op basis van affiliatie (Republikeins of Democratisch)[5]. Zij vonden hiervoor uiteindelijk een F_1 score van 0.684647.

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 F_1 score van 0.87 of hoger bereiken.

In het onderzoek van Høyland et al. werd een classificatiemodel voor partij-affiliatie op basis van teksten getraind op het vijfde Europese Parlement (1999-2004) en getest op het zesde Europese Parlement[6]. Hier verkregen zij een *macro average* F_1 score van 0.464.

2.1 Classificatiemethoden

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 eigennamen verwijderd.

In het onderzoek van Graeme Hirst et al. maakten ze gebruik van support vector machines[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 Bhand et al. gebruikten ze verschillende n-grams, inclusief verschillende manieren van *smoothing*[5]. Zij gebruikte als weging altijd de aanwezigheid van een woord. Als classificatiemodellen experimenteerden ze support vector machines en naive bayes classificatie. Voor het selecteren van *features* experimenteerden ze met een simpele minimale frequentie en het gebruik van een top aantal woorden op basis van mutual information. Uiteindelijk was het beste model bij hen een met support vector machine, met uni- en bigrams, gekozen op basis van mutual information.

In het onderzoek van Ferreira werd gebruik gemaakt van twee classificatiemethoden: Logistische regressie en MIRA[3]. Logistische regressie werd

127 aangevuld met *group Lasso* regularisatie. Voor wegenen van woorden werd
128 geëxperimenteerd met *tf*, *tf-idf*, Δ -*tf-idf* en Δ -*BM-25*. Daarnaast wordt er ge-
129 bruik gemaakt van woordclustering, *Concise Semantic Analysis* en stylome-
130 trische eigenschappen. Op *Part-Of-Speech tagging* na hadden stylometrische
131 eigenschappen een duidelijke negatieve invloed op de classificatie.

132 In het onderzoek van Høyland et al. werd gebruik gemaakt van een multi
133 class support vector machine[6]. Als beste waarde voor de regularisatieterm,
134 de C-parameter, vonden zij 0.8. Daarnaast gebruikten zij *dependency disambi-*
135 *guated stems* wat bij hen een F_1 score van twee procent hoger opleverden dan
136 normale stemming.

137 2.2 Invloed van oppositie of regering

138 Graeme Hirst et al. vonden in hun onderzoek dat de classificatie van spreker in
139 het Canadese parlement op basis van partij-affiliatie meer zegt over de status
140 van de partij (regering of oppositie).[2] Zo vergeleken zij de top tien karakteris-
141 tieke woorden van de liberalen en conservatieven in het 36e parlement (liberalen
142 in de regering) en het 39e parlement (conservatieven in de regering. Hier vonden
143 zij dat vier van de tien woorden van de liberalen (regering) in het 36e parlement
144 bij het 39e parlement bij de conservatieven (regering) te vinden waren. An-
145 dersom gebeurde hetzelfde met één van de tien woorden van de conservatieven
146 (oppositie) in het 36e parlement naar liberalen (oppositie) in het 39e parlement.

147 In hetzelfde onderzoek trainden ze ook hun classifiers op het ene par-
148 lement en testten deze op het andere parlement. Hierbij vonden zij in beide
149 gevallen een nauwkeurigheid ver onder de baseline. Daarnaast deden ze ook
150 nog een classificatie op de sprekers die in beide parlementen zaten en een an-
151 dere classificatie op sprekers die niet in beide parlementen zaten. Bij de eerste
152 classificatie vonden ze nauwkeurigheden rond de baseline, terwijl in de tweede
153 situatie nauwkeurigheden gevonden werden ver boven de baseline.

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

156 3 Methodologie

157 3.1 De data

158 De data die gebruikt worden, zijn de Handelingen van de Tweede Kamer gedu-
159 rende het missionaire kabinet-Rutte II (5 november 2012 tot 22 maart 2017).
160 Er is gekozen voor dit kabinet, omdat de data hiervoor makkelijk verkrijgbaar
161 was, het kabinet lang zat, waardoor er veel data is, en het recent is waardoor
162 het makkelijker te interpreteren is. Deze data zijn in xml-formaat van de web-
163 site officiële bekendmakingen.nl gehaald, samen met corresponderende metadata
164 xml-bestanden. De bestanden van de Handelingen bevatten voornamelijk infor-
165 matie over spreekbeurten tijdens een debat, waaronder naam van een spreker,
166 partij-affiliatie, inhoud van de spreekbeurt en het soort spreekbeurt. Deze ge-
167 gevens zijn samengevoegd tot een tabel en opgeslagen als csv-bestand.

168 Deze dataset bestaat uit een aantal soorten spreekbeurten; debat bijdra-
169 gen, interrupties en antwoorden. Debat bijdrage is de eerste onafgebroken
170 spreekbeurt die een spreker geeft achter het spreekgestoelte, aangeduid in de

171 xml-file met het attribuut *nieuw="ja"*. Interrupties zijn de vragen die andere
 172 politici stellen vanachter de interruptiemicrofoon aan de spreker. De antwoorden
 173 zijn vervolgens de reactie van een spreker achter het spreekgestoelte op een in-
 174 terruptie. Aangezien een debat bijdrage geïnterrupteerd kan worden, kan deze
 175 inhoudelijk doorlopen in een antwoord van een spreker. Er is in dit onderzoek
 176 ervoor gekozen om gebruik te maken van een debat bijdrage samengevoegd tot
 177 één document met alle bijbehorende antwoorden van die spreker. Daarnaast zijn
 178 er verschillende soorten sprekers; de voorzitter, Tweede Kamerleden, leden van
 179 het kabinet en gastsprekers. Daarnaast is alleen gekozen voor sprekers waarvan
 180 er een partij-affiliatie vermeld staat, dit is niet het geval voor leden van het
 181 kabinet, de voorzitter en gastsprekers (met uitzondering van Nederlandse leden
 182 van het Europees Parlement).

183 Deze dataset bevat vervolgens naast de verkozen partijen van de 2012
 184 Tweede Kamerverkiezingen, ook afsplitsingen van die partijen (tien in totaal) en
 185 bezoeken van vertegenwoordigingen van Nederlandse partijen uit het Europees
 186 Parlement (tien in totaal). Omdat van beide categorieën relatief weinig data is
 187 en er overlap zit met hun oorspronkelijke partij, zijn deze er uit gehaald. Op
 188 basis van de aantallen is er voor classificatie een baseline nauwkeurigheid van
 189 0.15 (door altijd grootste partij te kiezen) en baseline F_1 score van 0.11 (door
 190 willekeurig te voorspellen gewogen bij aantal spreekbeurten in klasse).

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

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

191 3.2 Methoden

192 3.2.1 Deelvraag 1

193 Om deze deelvraag te beantwoorden zullen een aantal classificatiemethoden ver-
 194 geleken worden. Aangezien het onmogelijk is om alle classificatiemethoden te
 195 vergelijken, beperkt dit onderzoek zich tot classificatiemethoden die gebruikt
 196 zijn in vergelijkbare onderzoeken, zoals besproken in 2.1. Er is ervoor geko-
 197 zen om alleen gebruik te maken van methoden waarvan reeds implementaties
 198 beschikbaar waren in Python. Hieronder worden de verschillende onderdelen
 199 besproken.

200 **Pre-processing** Voor pre-processing is gebruik gemaakt van tokenisation en
 201 lowercasing. Voor tokenisation is de reguliere expressie

202 $w+$ gebruikt, die daarmee alleen de letters en cijfers overhoudt. Deze woorden
203 zijn vervolgens allemaal omgezet in kleine letters. Vervolgens is er gevarieerd
204 tussen wel of geen gebruik maken van stemming. In het geval van stemming is
205 gebruik gemaakt van de Snowball Stemmer via de Python NLTK module.

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

212 Voor dit onderzoek zijn de volgende wegen voor woorden getest: *boolean*
213 (wel of niet aanwezig), *tf* (woordfrequentie), *tf-norm* (woordfrequentie genor-
214 maliseerd door documentlengte) en *tf-idf*. Daarnaast wordt in dit onderzoek
215 geëxperimenteerd met een minimale of maximale woord- of documentfrien-
216 tie. Ook is gekeken naar het effect van combinaties van unigrams, bigrams en
217 trigrams.

218 **Support Vector Machines en Logistische Regressie** De meest voorko-
219 mende techniek in vergelijkbaar onderzoek is Support Vector Machine (SVM).
220 Een andere techniek die gebruikt wordt is logistische regressie. Beide kennen een
221 eigen implementatie in sklearn, maar gezien de grootte van de dataset, duurt
222 dit te lang met een gridsearch. Om deze reden is er in beide gevallen voor ge-
223 kozen om gebruik te maken van de functie SGDClassifier, die beide technieken
224 leert met *stochastic gradient descent learning*. Er is hiervoor gevarieerd met de
225 regularisatie, learning rate en maximum aantal iteraties. Voor regularisatie is
226 hier geëxperimenteerd met Lasso en Ridge regularisatie, en een combinatie van
227 beide genaamd Elasticnet. De andere parameters zijn gelaten op de standaard-
228 waarden van scikit-learn.[7].

229 **Naive Bayes** Een simpelere techniek die gebruikt wordt voor politieke tekst-
230 classificatie is Naive Bayes. Dit algoritme neemt aan dat elke *feature* onafhan-
231 kelijk is ten op zichte van de rest. Dit is bij tekstclassificatie vaak niet het geval
232 omdat het gebruik van sommige woorden gepaard kan gaan met het gebruik
233 van andere woorden. Daarnaast is het gebruik van meerdere n-grams in een
234 classificatie schending van de aanname, want als bijvoorbeeld een bigram er in
235 voorkomt dan komen ook beide unigrams er sowieso in voor. Desalniettemin
236 blijkt Naive Bayes effectief te zijn voor tekstclassificatie[7, 5]. Hiervoor zijn de
237 functies van scikit-learn MultinomialNB en BernoulliNB gebruikt.[7, 5]

238 **Beoordelen van kwaliteit** De meest gebruikte methoden om kwaliteit van
239 politieke tekstclassificatie te beoordelen zijn accuracy en F_1 score, die opge-
240 bouwd is uit recall en precision. Deze scores zijn opgebouwd uit het aantal
241 correct positief (tp), foutief positief (fp), correct negatief (tn) en foutief nega-
242 tief (fn) geclassificeerde waarden.

$$Precision = \frac{tp}{tp + fp} \quad (1)$$

$$Recall = \frac{tp}{tp + tn} \quad (2)$$

$$Accuracy = \frac{tp + tn}{tp + tn + fp + fn} \quad (3)$$

$$F_1 = 2 * \frac{Precision * Recall}{Precision + Recall} \quad (4)$$

Deze waarden worden per klasse bepaald en daar wordt vervolgens een gemiddelde van genomen, gewogen bij documenten behorende tot die klasse. [8, 7].

Voor de classificatiemethoden wordt waar mogelijk gebruik gemaakt van functies van de Python module scikit-learn[7], aangevuld met zelf geschreven code als dit niet reeds beschikbaar is. Bij al deze classificatiemethoden wordt gevarieerd met meerdere parameters door middel van een gridsearch. Hierbij wordt gebruikt gemaakt van 5-fold cross-validation. Daardoor wordt de data gespleten in vijf delen, waarvan steeds één deel als testset wordt gebruikt en de rest voor training.

3.2.2 Deelvraag 2

In het onderzoek van Diermeier et al. worden alle eigennamen weggelaten zodat, volgens hen, namen van personen en partijen niet de classificatie domineren. Aangezien hier bij deelvraag 1 niet voor is gekozen, wordt bij deze deelvraag gekeken hoe groot het effect hiervan is, specifiek gericht op partijnamen en achternamen van kamerleden. Voor deze deelvraag wordt wederom een classificatie gedaan met de classificatiemethode die resulteerde uit deelvraag 1. In deze classificatie worden alle partijnamen vervangen door de tag PARTIJNAAM en alle namen van Kamerleden vervangen door de KAMERLIDNAAM. Deze namen zijn uit de Handelingen gehaald. Voor partijnamen zijn ook lidwoorden toegevoegd, voor achternamen van kamerleden zijn ook verkortingen meegenomen. Dit laatste omdat bijvoorbeeld *Van Haersma Buma* vaak aangesproken wordt als *Buma*. Voornamen van Kamerleden worden zelden tot nooit gebruikt, dus die zijn er niet uitgehaald. Een nadeel van deze aanpak is dat ook namen van niet-kamerleden of andere woorden weggehaald kunnen worden als deze hetzelfde zijn als naam van een kamerlid. Door gebruik van gevoeligheid voor hoofdletters is geprobeerd dit te voorkomen. Een opvallend voorbeeld hiervan is de naam Rutte, die zowel behoort tot het kamerlid Arno Rutte als de premier Mark Rutte.

De nauwkeurigheid en F_1 score worden vervolgens vergeleken met de resultaten uit deelvraag 1. Ook wordt gekeken naar verschillen tussen de meest veelzeggende woorden uit deelvraag 1 en uit deze deelvraag.

3.2.3 Deelvraag 3

Om deze deelvraag te beantwoorden zullen de twee experimenten die Graeme Hirst et al. uitvoerden voor dezelfde vraag gereproduceerd worden op de dataset van de Tweede Kamer. Bij deze deelvraag zal de beste classifier uit deelvraag 1 gebruikt worden. Daarnaast kan men ook naar de confusion matrix kijken of het aantal verkeerde classificaties groter is binnen regering of oppositie dan tussen elkaar.

284 In het eerste experiment zullen de tien meest karakteristieke woorden per
285 partij van de ene zittingsperiode vergeleken worden met de tien meest karak-
286 teristieke woorden per partij van de andere zittingsperiode. Als de classificatie
287 op basis van ideologie is in plaats van partij-status, is het te verwachten dat de
288 woorden bij een partij blijven en niet gekoppeld zijn aan in oppositie of regering
289 zitten.

290 In het tweede experiment worden classifiers getraind op de ene zittingspe-
291 riode en getest op de andere zittingsperiode. Als de classificatie op basis van
292 ideologie is in plaats van partij-status, is de verwachting dat er nog steeds aan-
293 zienlijke voorspellingen gedaan worden, aangezien de ideologie naar verwachting
294 redelijk stabiel is binnen tien jaar (hoewel woordgebruik varieert). Als de scores
295 aanzienlijk lager zijn, kan dit het gevolg zijn van het veranderen van partij-status
296 van partijen.

297 Als vergelijkingsmateriaal is voor deze experimenten een tweede dataset
298 nodig uit een ander kabinet. Hiervoor is het wenselijk dat dit kabinet bestaat
299 uit andere partijen dan kabinet-Rutte II. Daarnaast is het ook wenselijk als het
300 niet te ver terug is, zodat onderwerpen en taalgebruik enigszins overeenkomstig
301 zijn. Omdat kabinet-Rutte I een minderheidskabinet was met een bijzondere
302 partij-status voor de PVV, is ervoor gekozen om de Handelingen van de Tweede
303 Kamer tijdens het missionaire kabinet-Balkenende IV (22 februari 2007 tot 20
304 februari 2010) te gebruiken.

305 3.2.4 Deelvraag 4

306 Voor deze deelvraag vergelijken we de resultaten van de eerdere classificatie
307 per partij met een binaire classificatie op basis van rechts en links. Hiervoor
308 wordt wederom de dataset van kabinet-Rutte 2 gebruikt, met het model wat
309 resulteerde uit deelvraag 1.

310 Voor deze vraag moet vastgesteld worden welke partijen links en rechts
311 zijn. Omdat dit lastig te bepalen is en er meerdere indelingen zijn, wordt hier
312 gebruik gemaakt van twee verschillende indelingen. De indeling op basis van het
313 Kieskompas van Andre Krouwel voor de Kamerverkiezing 2012 en de indeling
314 volgens het Manifesto Project[9] gebaseerd op verkiezingsprogramma's voor de
315 Kamerverkiezing van 2012. In beide gevallen is de nullijn van het politieke
316 spectrum gebruikt om te bepalen of een partij links of rechts is.

Tabel 2: Rechts (R) of link (L) indeling per partij op basis van het Kieskompas en het Manifesto Project.

Partij	Kieskompas	Manifesto Project
SP	L	L
PvdA	L	L
GroenLinks	L	L
PvdD	L	L
50PLUS	L	L
D66	R	L
PVV	-	R
ChristenUnie	R	R
SGP	R	R
VVD	R	R
CDA	R	R

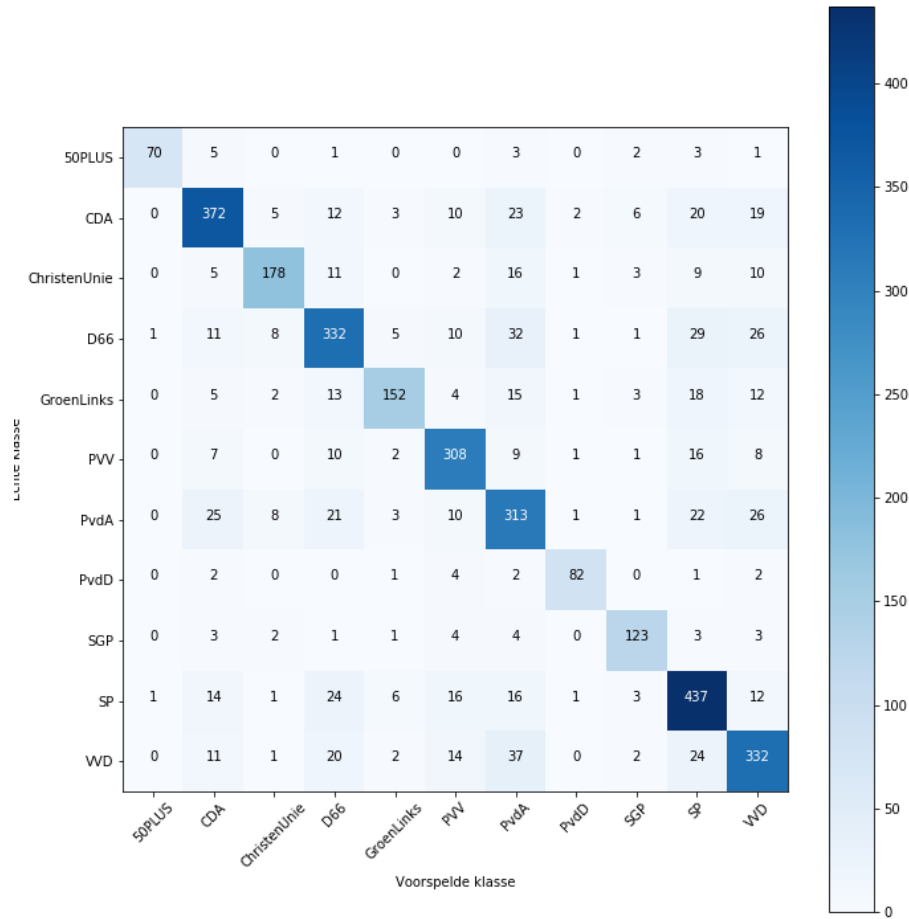
317 4 Evaluatie

318 4.1 Resultaten

319 4.1.1 Deelvraag 1

320 Het beste resultaat werd bereikt met SVM gebruikmakend van *stochastic gra-*
 321 *dient descent learning* en Ridge regularisatie.

Figuur 1: Confusion matrix van beste classificatie.



Tabel 3: Meest relevante woorden per partij op basis van beste classificatie.

50PLUS	CDA	ChristenUnie	D66	GroenLinks
50plus	cda	de christenunie	d66	groenlink
het lid krol	het cda	christenunie	led van veldhov	lid van tonger
lid krol	cda fractie	lid dik	mijn fractie	de led voortman
lid krol nar	de cda fractie	het lid dik	d66 is	led voortman
krol nar mij	de cda	lid dik faber	lid van veldhov	voortman
krol nar	het lid omtzigt	de led dik	van veldhov	led van tonger
krol	lid omtzigt	led dik faber	d66 wil	lid voortman nar
van 50plus	lid omtzigt nar	led dik	veldhov	lid voortman
gepensioneerd	led agnes mulder	led voordewind	d66 vindt	het lid voortman
50plus is	led agnes	de led voordewind	lid van men	van tonger

Tabel 3: Meest relevante woorden per partij op basis van beste classificatie.
(*Vervolg*)

PVV	PvdA	PvdD	SGP	SP	VVD
pvv	de pvda	het lid ouwehand	sgp	sp	de vvd
de pvv	pvda	lid ouwehand nar	de sgp	de sp	vvd
islamitisch	pvda fractie	lid ouwehand	sgp fractie	sp fractie	de vvd fractie
klever	de pvda fractie	de dier	de led dijkgraf	de sp fractie	vvd fractie
graus	de partij van	ouwehand	led dijkgraf	van gerv	de vvd is
miljard	de arbeid	dier	de sgp fractie	gerv	vor de vvd
madlener nar mij	van de arbeid	vor de dier	led van der	het lid smaling	vvd is
madlener nar	partij van de	thiem	de led bisschop	lid smaling	wat de vvd
nederland	arbeid	ouwehand nar	led bisschop	lid smaling nar	de vvd wil
lid madlener	partij van	ouwehand nar mij	dijkgraf	smaling	de vvd betre

322 4.1.2 Deelvraag 2

Tabel 4: Meest relevante woorden per partij op basis van classificatie zonder partij- of kamerlidnamen.

50PLUS	CDA	ChristenUnie	D66	GroenLinks
gepensioneerd	yyyyy fractie	gezinn	mijn fractie	schon
ouder	inwoner	voedselverspill	mijn	kamer hierover te
50 plusser	reger	prostitutie	buitengewon	schon energie
plusser	de reger	rookvrij	hervorm	vergroen
koopkrachtontwikkel	hier	inderdad	natur	bewindsperson
ouderenwerklos	eran	elkar	vandag	vluchtel
50	yyyyy	rechtsstat	daarom	sekswerker
werkend	antwoord	motie	fractie	zou
vor gepensioneerd	middeninkomen	ik heb	het kabinet	werkgeleg
gericht	zer	eurozon	kabinet	hierover te

Tabel 4: Meest relevante woorden per partij op basis van classificatie zonder partij- of kamerlidnamen. (*Vervolg*)

PVV	PvdA	PvdD	SGP	SP	VVD
islamitisch	kinder	dier	punt	huurder	veilig
brussel	jonger	milieu	allerlei	zegt	ondernemer
miljard	circulair	industrie	nadruk	bestuurder	yyyyy
nederland	mijn partij	de bio	bewindslid	herindel	regelgev
belastingbetaler	werk	burger	beantwoord	mens	speelveld
islam	voorzitter yyyyy	bio industrie	vor de beantwoord	armoed	ban
de islam	daarbij	constater dat	zou	voorstell	huis
asielzoeker	beter	bio	vanuit	de bevolk	kader
dit kabinet	sam	de bio industrie	oog	bezuin	instrument
partij vor de	docent	bestrijdingsmiddel	toen	bevolk	wat yyyyyy b

4.2 Discussie

4.2.1 Deelvraag 1

Dit onderzoek heeft zich beperkt tot methoden genoemd in vergelijkbare onderzoeken én waarvan de implementatie beschikbaar is in Python. Een aantal methoden die in gerelateerde literatuur leidden tot goede classificaties zijn daarom niet getest. Ook nieuwe methoden die nog niet gebruikt zijn in een vergelijkbaar onderzoek voor politieke tekst classificatie zijn daarom niet getest. Omdat dus niet alle opties getest zijn, kan geen uitsluitsel gegeven worden dat dit daadwerkelijk het classificatiemodel is. Voor vervolgonderzoek kan daarom gekeken worden naar meer verschillende methoden.

4.2.2 Deelvraag 4

Er zijn verschillende visies op links en rechts, en de indeling van de partijen, ook buiten de twee methoden gekozen in dit onderzoek.

5 Conclusies

Referenties

- [1] Felix Bießmann. Automating political bias prediction. *CoRR*, abs/1608.02195, 2016.
- [2] Graeme Hirst, Yaroslav Riabinin, Jory Graham, and Magali Boizot-roche. Text to ideology or text to party status? *.
- [3] Vasco Ferreira. Using textual transcripts of parliamentary interventions for profiling portuguese politicians. 2016.
- [4] Daniel Diermeier, Jean-François Godbout, Bei Yu, and Stefan Kaufmann. Language and ideology in congress. *British Journal of Political Science*, 42(1):31–55, 2012.

- 347 [5] Conal Sathi Maneesh Bhand, Dan Robinson. Text classifiers for political
348 ideologies, 2009.
- 349 [6] Bjørn Høyland, Jean-François Godbout, Emanuele Lapponi, and Erik Vell-
350 dal. Predicting party affiliations from european parliament debates. In
351 *Proceedings of the ACL 2014 Workshop on Language Technologies and*
352 *Computational Social Science*, pages 56–60. Association for Computatio-
353 nal Linguistics, 2014.
- 354 [7] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Gri-
355 sel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas,
356 A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay.
357 Scikit-learn: Machine learning in Python. *Journal of Machine Learning*
358 *Research*, 12:2825–2830, 2011.
- 359 [8] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. *Intro-*
360 *duction to Information Retrieval*. Cambridge University Press, New York,
361 NY, USA, 2008.
- 362 [9] Andrea Volkens, Pola Lehmann, Theres Matthieß, Nicolas Merz, Sven Re-
363 gel, and Bernhard Weßels. The manifesto data collection. manifesto project
364 (mrg/cmp/marpor). version 2017b, 2017.
- 365 [10] Bei Yu, Stefan Kaufmann, and Daniel Diermeier. Classifying party affilia-
366 tion from political speech. *Journal of Information Technology & Politics*,
367 5(1):33–48, 2008.

368 A Slides