

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 .



UNIVERSITEIT VAN AMSTERDAM

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	11
30	4.2 Discussie	12
31	4.2.1 Deelvraag 1	12
32	4.2.2 Deelvraag 4	12
33	5 Conclusies	12
34	A Slides	13

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 andere politiek stelsel en cultuur heeft, verschillen de resultaten.
47 Daarnaast gebruikt elk onderzoek ook een andere methode voor het classificeren.
48 Daarnaast vinden sommige onderzoeken dat deze classificatie minder het gevolg
49 is van ideologie 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 is deze classificatie afhankelijk van of een partij rechts of links
65 is?

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

71 **Overzicht van scriptie** In sectie 2 zal gerelateerd werk besproken worden,
72 met name vergelijkbare onderzoeken in andere landen. In sectie 3 zal de me-
73 thodologie van de verschillende deelvragen behandeld worden. In sectie 4 zul-
74 len vervolgens de resultaten weergegeven worden. In sectie 5 zal een evaluatie
75 plaatsvinden van zowel de resultaten als de gehanteerde methodologie. In sectie
76 6 wordt ten slotte de onderzoeksvraag beantwoord.

77 2 Gerelateerd werk

78 Diermeier et al. deden onderzoek naar het classificeren op basis van ideologi-
79 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

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

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

138 2.2 Invloed van oppositie of regering

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

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

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

157 3 Methodologie

158 3.1 De data

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

169 Deze dataset bestaat uit een aantal soorten spreekbeurten, zoals speeches,
170 interrupties en antwoorden. Daarnaast ook door verschillende soorten sprekers,
171 zoals de voorzitter, Tweede Kamerleden, leden van het kabinet en gastsprekers.

172 Uit deze dataset is gekozen voor de eerste spreekbeurt nadat een spreker achter
 173 het spreekgestoelte is gaan staan, aangezien deze vaak langer zijn dan de andere
 174 spreekbeurten en naar verwachting meer zeggen over ideologie. In de oorspron-
 175 kelijke xml-bestanden hadden deze spreekbeurten het attribuut *nieuw*="ja".
 176 Daarnaast is alleen gekozen voor sprekers waarvan er een partij-affiliatie ver-
 177 meld staat, dit is niet het geval voor leden van het kabinet, de voorzitter en
 178 gastsprekers (met uitzondering van Nederlandse leden van het Europees Parle-
 179 ment).

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

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

188 3.2 Methoden

189 3.2.1 Deelvraag 1

190 Om deze deelvraag te beantwoorden zullen een aantal classificatiemethoden ver-
 191 geleken worden. Aangezien het onmogelijk is om alle classificatiemethoden te
 192 vergelijken, beperkt dit onderzoek zich tot classificatiemethoden die gebruikt
 193 zijn in vergelijkbare onderzoeken, zoals besproken in 2.1. Daarnaast is omwille
 194 van de tijd ervoor gekozen om alleen gebruik te maken van methoden waar-
 195 van reeds implementaties beschikbaar waren in Python. Hieronder worden de
 196 verschillende onderdelen besproken.

197 **Pre-processing** Voor pre-processing is gebruik gemaakt van tokenisation en
 198 lowercasing. Voor tokenisation is de reguliere expressie
 199 $w+$ gebruikt, die daarmee alleen de letters van het alfabet overhoudt. Deze
 200 woorden zijn vervolgens allemaal omgezet in kleine letters. Vervolgens is er
 201 gevarieerd tussen wel of geen gebruik maken van stemming. In het geval van

stemming is gebruik gemaakt van de Snowball Stemmer via de Python NLTK module.

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

Voor dit onderzoek zijn de volgende wegingen voor woorden getest: *boolean* (wel of niet aanwezig), *tf* (woordfrequentie), *tf-norm* (woordfrequentie genormaliseerd door documentlengte) en *tf-idf*. Daarnaast wordt in dit onderzoek geëxperimenteerd met een minimale of maximale woord- of documentfrequentie. Ook is gekeken naar het effect van combinaties van unigrams, bigrams en trigrams.

Support Vector Machines en Logistische Regressie De meest voorkomende techniek in vergelijkbaar onderzoek is Support Vector Machine (SVM). Een andere techniek die gebruikt wordt is logistische regressie. Beide kennen een eigen implementatie in sklearn, maar gezien de grootte van de dataset, duurt dit te lang met een gridsearch. Om deze reden is er in beide gevallen voor gekozen om gebruik te maken van de functie *SGDClassifier*, die beide technieken leert met *stochastic gradient descent learning*. Er is hiervoor gevarieerd met de regularisatie, learning rate en maximum aantal iteraties. De andere parameters zijn gelaten op de standaardwaarden van scikit-learn[7].

Naive Bayes Een simpelere techniek die gebruikt wordt voor politieke tekstclassificatie is Naive Bayes. Dit algoritme neemt aan dat elke *feature* onafhankelijk is ten op zichte van de rest. Dit is bij tekstclassificatie vaak niet het geval omdat het gebruik van sommige woorden gepaard kan gaan met het gebruik van andere woorden. Daarnaast is het gebruik van meerdere n-grams in een classificatie schending van de aanname, want als bijvoorbeeld een bigram er in voorkomt dan komen ook beide unigrams er sowieso in voor. Desalniettemin blijkt Naive Bayes effectief te zijn voor tekstclassificatie[7, 5].

Er zijn twee frequent gebruikte aannames voor de distributies in tekstclassificatie; *Multinomial* en *Bernoulli*. In gerelateerde werken wordt niet gespecificeerd welke gebruikt wordt. Om deze reden zijn ze allebei gebruikt. Hiervoor zijn respectievelijk de functies van scikit-learn *MultinomialNB* en *BernoulliNB* gebruikt.[7, 5]

Beoordelen van kwaliteit De meest gebruikte methoden om kwaliteit van politieke tekstclassificatie te beoordelen zijn accuracy en F_1 score, die opgebouwd is uit recall en precision. Deze scores zijn opgebouwd uit het aantal correct positief (*tp*), foutief positief (*fp*), correct negatief (*tn*) en foutief negatief (*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. Voor partijnamen zijn ook lidwoorden daarvoor meegenomen, voor achternamen van kamerleden zijn ook verkortingen meegenomen. Dit laatste om dat bijvoorbeeld *Van Nieuwenhuizen-Wijbenga* vaak genoemd wordt als *Van Nieuwenhuizen*. 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. Deze resultaten worden vervolgens vergeleken met de resultaten uit deelvraag 1.

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.

Als vergelijkingsmateriaal is voor deze experiment een tweede dataset nodig uit een ander kabinet. Hiervoor is het wenselijk dat dit kabinet bestaat uit andere partijen dan kabinet-Rutte II. Daarnaast is het ook wenselijk als het niet te ver terug is, zodat onderwerpen en taalgebruik enigszins overeenkomstig zijn. Omdat kabinet-Rutte I een minderheidskabinet was met een bijzondere

partij-status voor de PVV, is ervoor gekozen om de Handelingen van de Tweede Kamer tijdens het missionaire kabinet-Balkenende IV (22 februari 2007 tot 20 februari 2010) te gebruiken.

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

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

3.2.4 Deelvraag 4

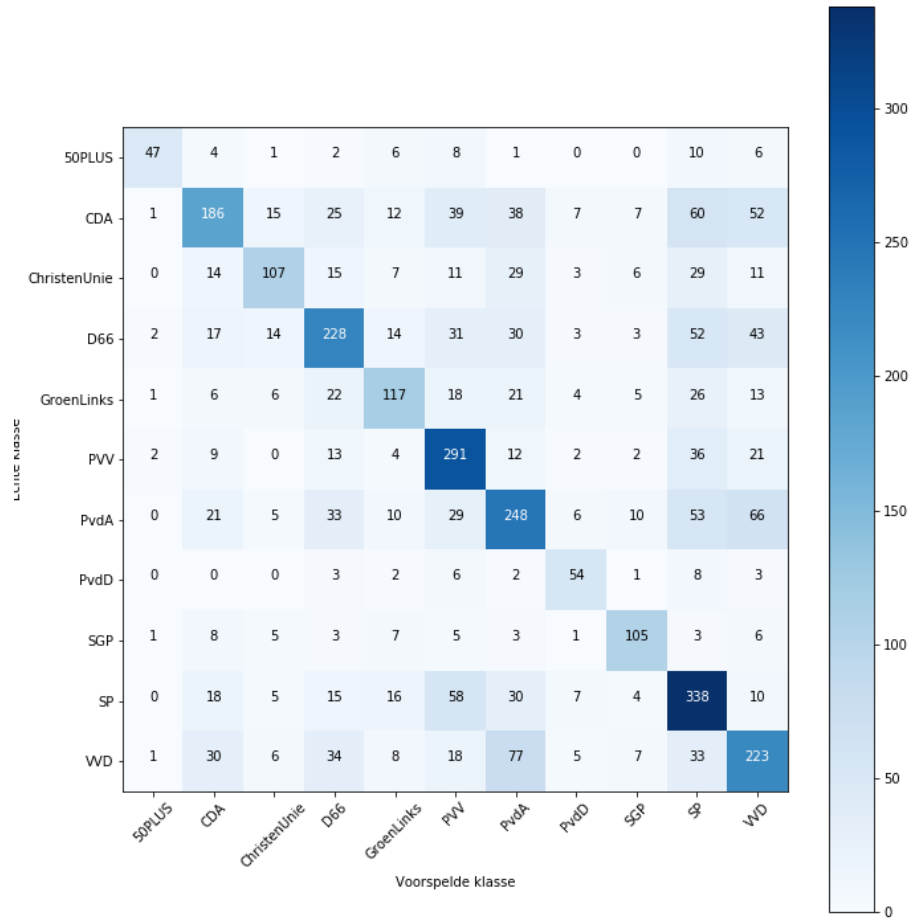
Voor deze deelvraag vergelijken we de resultaten van de eerdere classificatie per partij met een binaire classificatie op basis van rechts en links. Hiervoor wordt wederom de dataset van kabinet-Rutte 2 gebruikt, met het model wat resulteerde uit deelvraag 1.

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

Figuur 1: Confusion matrix van beste classificatie.



4 Evaluatie

4.1 Resultaten

4.1.1 Deelvraag 1

Het beste resultaat werd bereikt met SVM gebruikmakend van *stochastic gradient descent learning*

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

50PLUS	CDA	ChristenUnie	D66	GroenLinks
gepensioneerd	inwoner	prostitutie	mijn fractie	schon
plusser	yyyyy fractie	rookvrij	mijn	schon energie
50 plusser	reger	voedselverspill	hervorm	belastingontwijk
ouder	de regering	rechtsstat	buitengewon	kamer hierover te
ouderenwerklos	yyyyy	ik vraag	vandag	zou
50	antwoord	mensenhandel	natur	in elk geval
yyyyy	hier	publiek belang	kans	elk geval
koopkrachtontwikkel	eran	zull	fractie	werkgeleg
vor gepensioneerd	limburg	schepping	het kabinet	bewindsperson
gericht	onz inwoner	inderdad	kabinet	hierover te informeer

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

PVV	PvdA	PvdD	SGP	SP	VVD
islamitisch	circulair	dier	eenverdiener	huurder	ondernemer
nederland	kinder	milieu	mevrouw de voorzitter	armoed	regelgev
miljard	hun	klimaatverander	allerlei	mens	speelveld
brussel	lager over	burger	vor de beantwoord	voorstell	veilig
belastingbetaler	voorzitter yyyyy	dierenwelzijn	dank zer	segregatie	yyyyy frac
islam	mijn partij	aard	punt	zegt	bedrijfslev
partij	circulair economie	constater dat	mevrouw de	bezuin	de sector
de islam	tevred	de aard	bewindslid	herindel	ban
grenz	verdien	constater	wel	waarbij	huis
asielzoeker	collega	verbod	aandacht	erg	yyyyy

316 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	inwoner	prostitutie	mijn fractie	schon
plusser	yyyyy fractie	rookvrij	mijn	schon energie
50 plusser	reger	voedselverspill	hervorm	belastingontwijk
ouder	de regering	rechtsstat	buitengewon	kamer hierover te
ouderenwerklos	yyyyy	ik vraag	vandag	zou
50	antwoord	mensenhandel	natur	in elk geval
yyyyy	hier	publiek belang	kans	elk geval
koopkrachtontwikkel	eran	zull	fractie	werkgeleg
vor gepensioneerd	limburg	schepping	het kabinet	bewindsperson
gericht	onz inwoner	inderdad	kabinet	hierover te informeer

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

PVV	PvdA	PvdD	SGP	SP	VVD
islamitisch	circulair	dier	eenverdiener	huurder	onderneme
nederland	kinder	milieu	mevrouw de voorzitter	armoed	regelgev
miljard	hun	klimaatverander	allerlei	mens	speelveld
brussel	lager over	burger	vor de beantwoord	voorstell	veilig
belastingbetaler	voorzitter yyyyy	dierenwelzijn	dank zer	segregatie	yyyyy frac
islam	mijn partij	aard	punt	zegt	bedrijfslev
partij	circulair economie	constater dat	mevrouw de	bezuin	de sector
de islam	tevred	de aard	bewindslied	herindel	ban
grenz	verdien	constater	wel	waarbij	huis
asielzoeker	collega	verbod	aandacht	erg	yyyyy

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.

- 341 [5] Conal Sathi Maneesh Bhand, Dan Robinson. Text classifiers for political
342 ideologies, 2009.
- 343 [6] Bjørn Høyland, Jean-François Godbout, Emanuele Laponi, and Erik Vell-
344 dal. Predicting party affiliations from european parliament debates. In
345 *Proceedings of the ACL 2014 Workshop on Language Technologies and*
346 *Computational Social Science*, pages 56–60. Association for Computatio-
347 nal Linguistics, 2014.
- 348 [7] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Gri-
349 sel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas,
350 A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay.
351 Scikit-learn: Machine learning in Python. *Journal of Machine Learning*
352 *Research*, 12:2825–2830, 2011.
- 353 [8] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. *Intro-*
354 *duction to Information Retrieval*. Cambridge University Press, New York,
355 NY, USA, 2008.
- 356 [9] Andrea Volkens, Pola Lehmann, Theres Matthieß, Nicolas Merz, Sven Re-
357 gel, and Bernhard Weßels. The manifesto data collection. manifesto project
358 (mrg/cmp/marpor). version 2017b, 2017.
- 359 [10] Bei Yu, Stefan Kaufmann, and Daniel Diermeier. Classifying party affilia-
360 tion from political speech. *Journal of Information Technology & Politics*,
361 5(1):33–48, 2008.

362 A Slides