

1 IDEOLOGIE EN CLASSIFICATIE IN DE HANDELINGEN
2 VAN DE TWEEDE KAMER

3 INGEDIEND VOOR GEDEELTELIJKE VOLTOOIING VAN GRAAD VAN
4 BACHELOR OF SCIENCE

5 JASPER VAN DER HEIDE
6 10732721

7 BACHELOR INFORMATIEKUNDE
8 FACULTEIT DER NATUURWETENSCHAPPEN, WISKUNDE EN
9 INFORMATICA
10 UNIVERSITEIT VAN AMSTERDAM

11 2018-06-28

12

	Begeleider	Tweede lezer
Titel, Naam	Dr Maarten Marx	Ir Loek Stolwijk
Affiliatie	UvA, FNWI, IvI	UvA, FNWI, IvI
Email	M.J.Marx@uva.nl	A.M.Stolwijk@uva.nl



UNIVERSITEIT VAN AMSTERDAM

14	Inhoudsopgave	
15	1 Introductie	3
16	2 Gerelateerd werk	4
17	2.1 Tekstclassificatie van parlementaire teksten	4
18	2.2 Classificatiemethoden	5
19	2.3 Invloed van partijnamen of sprekersnamen	6
20	2.4 Invloed van oppositie of regering	6
21	3 Methodologie	7
22	3.1 De data	7
23	3.2 Methoden	9
24	3.2.1 DV1: Beste classificatiemethode	9
25	3.2.2 DV2: Invloed van namen	11
26	3.2.3 DV3: Oppositie of regering	12
27	3.2.4 DV4: Links/rechts	14
28	3.2.5 DV5: Woordgebruik van sprekers	15
29	4 Resultaten	15
30	4.1 DV1: Beste classificatiemethode	15
31	4.2 DV2: Invloed van namen	17
32	4.3 DV3: Oppositie of regering	19
33	4.4 DV4: Links/rechts	22
34	4.5 DV5: Woordgebruik van sprekers	22
35	5 Discussie	23
36	5.1 DV1: Beste classificatiemethode	23
37	5.2 DV2: Invloed van namen	24
38	5.3 DV3: Oppositie of regering	25
39	5.4 DV4: Links/rechts	27
40	5.5 DV5: Woordgebruik van sprekers	27
41	5.6 Algemeen	27
42	6 Conclusies	28

43

Samenvatting

44

1 Introductie

Teksten van politieke partijen kunnen dienen als bron voor het bepalen van ideologische positie van andere teksten, aangezien zij zowel tekst bevatten als ook een bekende ideologie in de vorm van een partij van de spreker; de partij-affiliatie. Het classificeren op basis van tekst kan inzichten geven over de relatie tussen ideologie en woordgebruik. Deze informatie kan vervolgens toegepast worden op andere teksten die wellicht ideologisch van aard zijn. Bijvoorbeeld kan men aan de hand van deze informatie teksten uit kranten classificeren op basis van ideologie [1, 2].

In diverse landen zijn al onderzoeken gedaan naar het classificeren naar partij-affiliatie op basis van teksten van politici [1, 3]. Met deze tekstclassificatie naar partij-affiliatie proberen onderzoekers uit te vinden in hoeverre ideologie terug te vinden is in teksten van politici. De resultaten van de tekstclassificaties zijn in alle gevallen ruim boven de baseline. Diverse onderzoeken wijzen daarentegen ook naar redenen dat dit niet alleen het gevolg is van ideologie. Zo suggereren de resultaten van Hirst et al. [2] dat de partij-status (oppositie tegenover regering) van invloed is op de classificatie. Daarnaast laat dit onderzoek ook zien dat partijnamen een grote invloed hebben op de classificatie.

Een onderzoek gericht op het Nederlandse parlement is niet gevonden. Ook beperken veel onderzoeken zich vaak tot één classificatiemethode.

Dit onderzoek richt zich daarom op meer classificatiemethoden en daarnaast 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 naar partij-affiliatie in de Tweede Kamer en wat is het resultaat van dit model?
2. In hoeverre is deze classificatie afhankelijk van achternamen van Kamerleden en partijen?
3. In hoeverre wordt deze classificatie bepaald door partij-status (oppositie of regering)?
4. In hoeverre wordt deze classificatie bepaald door links/rechts positie?
5. In hoeverre wordt deze classificatie bepaald door woordgebruik van sprekers?

Voor de eerste deelvraag zullen Support Vector Machine, Logistische Regressie en Naive Bayes met verschillende parameters vergeleken worden aan de hand van *accuracy* en F_1 score. Bij de tweede deelvraag wordt gekeken naar classificatie zonder achternamen van Kamerleden en partijnamen of met alleen achternamen van Kamerleden en partijnamen. De derde vraag bestaat uit meerdere experimenten, waarin gekeken zal worden naar de hoeveelheid misclassificaties binnen regering of oppositie tegenover tussen regering en oppositie. Daarnaast zal gekeken worden naar overlap in woordgebruik binnen regering en verschil in scores als een partij gewisseld is van partij-status.

89 **Overzicht van scriptie** Sectie 2 bevat vergelijkbare onderzoeken in andere
90 parlementen. Sectie 3 bevat de methodologie van de verschillende deelvragen.
91 Sectie 4 bevat de resultaten. Sectie 5 bevat de evaluatie van zowel de resultaten
92 als de methodologie. Sectie 6 bevat ten slotte het antwoord op de onderzoeks-
93 vraag.

94 2 Gerelateerd werk

95 Toespraken in parlementen worden veel gebruikt in tekstclassificatie, omdat
96 deze veel nette tekst bevatten en vaak gelabeld zijn. Labels zijn bijvoorbeeld
97 naam en partij van de spreker, maar ook daar uit afleidbare labels zoals geslacht,
98 leeftijd en partij-status (oppositie of regering).

99 In dit hoofdstuk zullen verschillende onderzoeken behandeld worden die
100 tekstclassificatie hebben toegepast op parlementaire teksten. Eerst zullen de
101 onderzoeken algemeen besproken worden. Vervolgens zal uitgebreider gekeken
102 worden naar de effecten van verschillende classificatiemethoden. In de latere
103 secties zullen aspecten besproken worden die in vergelijkbare onderzoeken ge-
104 noemd worden als van invloed op de classificatie.

105 2.1 Tekstclassificatie van parlementaire teksten

106 Diermeier et al. deden onderzoek naar het classificeren op basis van ideologische
107 positie in de Amerikaanse Senaat [4]. Ze trainden hun classificatie op de speeches
108 van de 25 meest liberale en de 25 meest conservatieve senatoren van het 101e
109 tot en met het 107e Congres en testten op dezelfde categorieën van het 108e
110 Congres. Een document was in dit onderzoek de verzameling van alle speeches
111 van een senator in een Congres. Deze classificatie resulteerde uiteindelijk in
112 een *accuracy* van 94% (baseline van 50%). Van de 50 senatoren in de test set,
113 kwamen er 44 al voor in de training set, doordat de training op voorgaande
114 Congressen was.

115 Later in het onderzoek vergeleken ze ook de 25 gematigd conservatieve
116 en de 25 gematigd liberale senatoren van dezelfde Congressen. Het resultaat
117 hiervan was 52% (baseline van 50%), dus nauwelijks beter dan de baseline.
118 Als verklaring voor dit verschil ten opzichte van de uitersten zeggen ze dat
119 gematigden een minder duidelijke ideologie hebben.

120 Yu et al. [5] richtte zich vervolgens op zowel het Amerikaanse Huis van
121 Afgevaardigden als de Senaat in 2005. Een document was in dit onderzoek de
122 verzameling van alle speeches van een senator in een Congres en het label de
123 partij. Voor het Huis van Afgevaardigden vonden ze een *accuracy* van 80.1%
124 (baseline van 51.5%) en voor de Senaat 86.0 % (baseline van 55.0%). Ze testten
125 hun classificaties ook op de andere kamer. Van Huis van Afgevaardigden naar
126 Senaat leverde dit een *accuracy* op van 88.0% (baseline van 55.0%) en andersom
127 67.6% (baseline van 51.5%). Hun verklaring voor dit verschil is dat het Huis
128 van Afgevaardigden sterker verdeeld is langs partijlijnen.

129 Vervolgens herhaalden ze de classificaties op het Huis van Afgevaardigden
130 uit 2015, maar testten ditmaal op de Senaat elk jaar tussen 1989 en 2006 afzon-
131 derlijk. Hier zien zij een stijging in *accuracy* van 60% (baseline van 55.0%) in
132 1989 naar 87.0% (baseline van 55.0%) in 2006, maar met twee duidelijke dalen.

133 Ze presenteren twee mogelijke verklaringen voor de trend; het veranderen van
134 de onderwerpen en het sterker verdeeld worden van het Congres.

135 Als een vervolg op deze onderzoeken deden Hirst et al. een vergelijkbaar
136 onderzoek naar het Canadese Parlement [2]. Hierbij werd zowel gekeken naar de
137 Engelse als Franse teksten. Een document werd hier gezien als de samenvoeging
138 van alle spreekbeurten van een spreker. Afhankelijk van taal en dataset vinden
139 zij in dit onderzoek *accuracy* scores van 83.2% en hoger (baseline van 65.5%).

140 Het onderzoek bevat ook een classificatie van het Europees Parlement.
141 Hierbij voegen ze alle teksten van een parlamentslid bij elkaar en delen die
142 op in stukken van gelijke grootte. Zij vinden voor documentgrootte van 267
143 woorden een *accuracy* van 44.0% oplopend tot 61.8% (baseline van 38-39%)
144 voor documentgrootte van 6666.

145 Het onderzoek van Bhand et al. richtte zich op het classificeren van le-
146 den van het Amerikaanse Congres in 2005, op basis van affiliatie (Republikeins
147 of Democratisch) [6]. Een document hierbij was in tegenstelling tot eerderge-
148 noemde onderzoeken een speech. Zij vonden hiervoor uiteindelijk een F_1 score
149 van 0.68 (baseline niet vermeld).

150 Ferreira probeerde interventies van politici te classificeren op basis van
151 geslacht, leeftijdsgroep, partij-affiliatie en oriëntatie in het Portugese parlement
152 [3]. In het geval van classificatie op basis van partij-affiliatie bereikte men een
153 F_1 score van 0.90 (baseline niet vermeld, zes partijen).

154 In het onderzoek van Høyland et al. werd een classificatiemodel voor
155 partij-affiliatie op basis van teksten getraind op het vijfde Europese Parlement
156 (1999-2004) en getest op het zesde Europese Parlement (2004-2009) [7]. Alle
157 teksten van een spreker zijn samengevoegd tot één document. 40% van de
158 sprekers in de test set zaten ook in de training set. Hier verkregen zij een
159 *macro* F_1 score van 0.464 (baseline van 0.097) en *accuracy* van 0.551 (baseline
160 van 0.410). Hun baseline is op basis van altijd classificeren als grootste partij,
161 terwijl voor F_1 score de baseline hoger ligt als hiervoor gekozen wordt voor
162 gokken gewogen bij grootte van een klasse.

163 2.2 Classificatiemethoden

164 Diermeier et al. [4] gebruikten Support Vector Machines. Verder maakten ze
165 gebruik van *tf-idf* met een minimale woordfrequentie van 50 en een minimale
166 documentfrequentie van 10 en *Part-Of-Speech tagging*.

167 Yu et al. [5] maakten gebruik van Support Vector Machines en Naive
168 Bayes, waarvan de varianten multinomial en Bernoulli. De features waren unig-
169 rams, met minimale woordfrequentie van drie en de top 50 meest voorkomende
170 woorden weggelaten. Voor de wegingen van de features bij Support Vector Ma-
171 chines werd geëxperimenteerd met *boolean*, *tf-norm* en *tf-idf*. Het beste resultaat
172 was afhankelijk van welke kamer Voor het huis van afgevaardigden was het Sup-
173 port Vector Machines met als weging *tf-idf* en voor de Senaat Bernoulli Naive
174 Bayes.

175 Hirst et al. maakten gebruik van Support Vector Machines [2]. Ze experi-
176 menteerden met verschillende vormen van pre-processing, inclusief stemmen en
177 het verwijderen van woorden op basis van te hoge of te lage frequentie. Deze
178 variaties maakten in hun onderzoek geen grote verschillen en uiteindelijk is ge-
179 kozen voor het niet stemmen, het weglaten van woorden die in minder dan
180 vijf documenten voorkomen en resultaten van zowel met als zonder de top 500

181 meest frequente woorden. Daarnaast werd geëxperimenteerd met vier wegingen
182 voor woorden: *boolean*, *tf*, *tf-norm* en *tf-idf*, waarvan *tf-idf* het beste resultaat
183 opleverde.

184 Bhand et al. gebruikten verschillende n-grams, inclusief verschillende ma-
185 nieren van *smoothing*[6]. Ze testten als weging voor features zowel *boolean* als
186 *tf*, waarbij ze vonden concludeerden dat *boolean* betere resultaten opleverden.
187 Voor classificatiemodel experimenteerden ze met SVM en Naive Bayes . Voor
188 het selecteren van *features* experimenteerden ze met een minimale frequentie en
189 selectie van woorden op basis van hoogste mutual information. Uiteindelijk was
190 het beste model bij hen een SVM met uni- en bigrams en geselecteerd op basis
191 van mutual information.

192 In het onderzoek van Ferreira werd gebruik gemaakt van twee classifi-
193 catiemethoden: Logistische regressie en MIRA[3]. Logistische regressie werd
194 aangevuld met *group Lasso* regularisatie. Voor wegingen van woorden werd
195 geëxperimenteerd met *tf*, *tf-idf*, Δ -*tf-idf* en Δ -*BM-25*. Daarnaast wordt er ge-
196 bruik gemaakt van woordclustering, *Concise Semantic Analysis* en stylome-
197 trische eigenschappen. Op *Part-Of-Speech tagging* na hadden stylometrische
198 eigenschappen een duidelijke negatieve invloed op de classificatie.

199 Høyland et al. maakten gebruik van Support Vector Machine[7]. Als beste
200 waarde voor de regularisatieterm, de C-parameter, vonden zij 0.8. Daarnaast
201 gebruikten zij *dependency disambiguated stems* wat bij hen een F_1 score van
202 twee procent hoger opleverden dan normale stemming.

203 2.3 Invloed van partijnamen of sprekersnamen

204 Diermeier et al. lieten de namen van de sprekers en verwijzingen naar staten
205 die de senatoren representeren weg, omdat deze volgens hen de classificatie te
206 makkelijk zouden maken [4]. Hirst et al. vinden inderdaad dat partijnamen (en
207 het weglaten daarvan) bij het Europees Parlement een grote invloed hebben op
208 de classificatie [2]. Bij het Europees Parlement zien zij met name het gebruik
209 van de eigen partijnaam door een spreker, terwijl zij in het Canadese parlement
210 vooral zien dat de naam van de andere partij gebruikt wordt door een spreker.

211 2.4 Invloed van oppositie of regering

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

221 In hetzelfde onderzoek trainden ze ook hun classificaties op het ene par-
222 lement en testten deze op het andere parlement. Hierbij vonden zij in beide
223 gevallen een *accuracy* ver onder de baseline. Daarnaast deden ze ook nog een
224 classificatie op de sprekers die in beide parlementen zaten en een andere classi-
225 ficatie op sprekers die niet in beide parlementen zaten. Bij de eerste classificatie

226 vonden ze *accuracy* scores rond de baseline, terwijl in de tweede situatie *accuracy*
227 scores gevonden werden ver boven de baseline.

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

230 3 Methodologie

231 3.1 De data

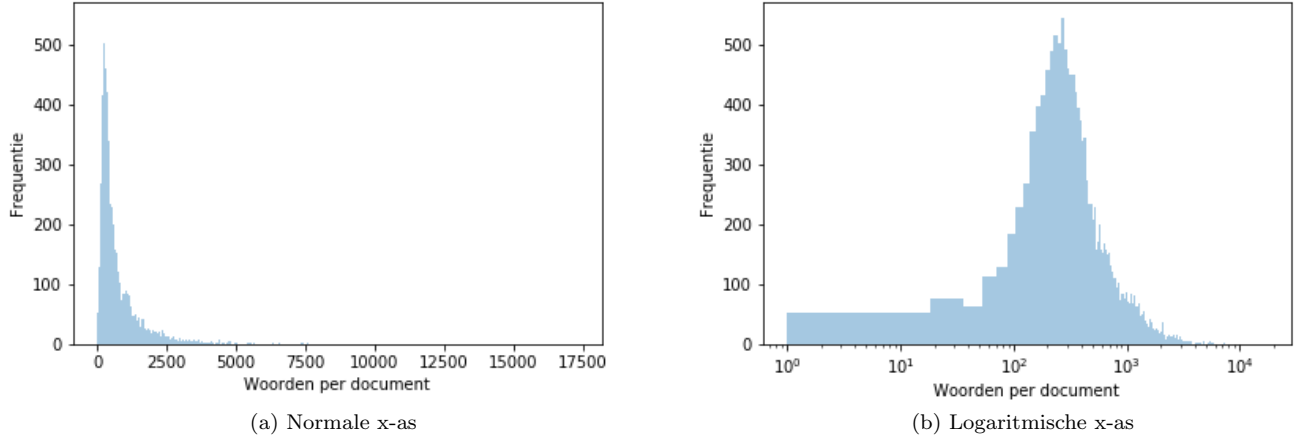
232 De data die gebruikt worden, zijn de Handelingen van de Tweede Kamer gedu-
233 rende het missionaire kabinet-Rutte II (5 november 2012 tot 22 maart 2017).
234 Er is gekozen voor dit kabinet, omdat de data hiervoor makkelijk verkrijgbaar
235 was, het kabinet lang zat - waardoor er veel data is - en het recent is waardoor
236 het makkelijker te interpreteren is. In dit kabinet zaten de PvdA en VVD. Deze
237 data zijn in xml-formaat van de website officielebekendmakingen.nl gehaald sa-
238 men met bijbehorende metadatabestanden. De bestanden van de Handelingen
239 bevatten voornamelijk informatie over spreekbeurten tijdens een debat, waar-
240 onder naam van een spreker, partij-affiliatie, inhoud van de spreekbeurt en het
241 soort spreekbeurt. Deze gegevens zijn samengevoegd tot één tabel.

242 Deze dataset bestaat uit een aantal soorten spreekbeurten; debatbijdra-
243 gen, interrupties en antwoorden. Een debatbijdrage is de eerste onafgebroken
244 spreekbeurt die een spreker geeft achter een spreekgestoelte, aangeduid in de
245 xml-file met het attribuut *nieuw="ja"*. Dit kan een bijdrage in een debat zijn
246 of een vraag tijdens een vragenuur. Interrupties zijn de vragen die andere poli-
247 tici stellen vanachter de interruptiemicrofoon aan een spreker. De antwoorden
248 zijn vervolgens de reactie van een spreker achter het spreekgestoelte op een in-
249 terruptie. Aangezien een debatbijdrage geïnterrupteerd kan worden, kan deze
250 inhoudelijk doorlopen in een antwoord van een spreker. Gerelateerde onderzoe-
251 ken voegen vaak alle teksten van een spreker samen tot één document. Dit is
252 alleen niet mogelijk met de hoeveelheid kleine partijen in de Tweede Kamer,
253 die dan niet altijd in een training of test set zijn vertegenwoordigd. Daarom
254 is in dit onderzoek ervoor gekozen om een debatbijdrage met alle bijbehorende
255 antwoorden samen te voegen tot één document voor de classificatie.

256 Daarnaast zijn er verschillende soorten sprekers; de voorzitter, Tweede
257 Kamerleden, leden van het kabinet en gastsprekers. Hieruit is alleen gekozen
258 voor sprekers waarvan er een partij-affiliatie vermeld staat. Dit is niet het geval
259 voor leden van het kabinet, de voorzitter en gastsprekers (met uitzondering van
260 Nederlandse leden van het Europees Parlement).

261 Deze dataset bevat vervolgens naast de verkozen partijen na de Tweede
262 Kamerverkiezingen van 2012 ook afsplitsingen van die partijen (tien in totaal) en
263 bezoeken van vertegenwoordigingen van Nederlandse partijen uit het Europees
264 Parlement (tien in totaal). Omdat van beide categorieën relatief weinig data
265 is en er overlap zit met hun oorspronkelijke of gelieerde partij, zijn deze er
266 uit gehaald. 50PLUS is in 2014 [8] uiteengevallen in twee fracties die aanspraak
267 maakten op de partij-affiliatie 50PLUS. Vanaf dit moment zijn deze documenten
268 niet meer meegenomen om onduidelijkheid te voorkomen.

269 De documenten verschillen in grootte. De distributie van documentgrootte
270 lijkt op een lognormale verdeling, maar met een Kolmogorov-Smirnov test is hier
271 geen bewijs voor gevonden [9].



Figuur 1: Aantal woorden per document

Om toch de uitschieters er uit te halen, is aangenomen dat de distributie wel lognormaal verdeeld is en zijn daarmee de documenten buiten het betrouwbaarheidsinterval van 95% eruit gehaald. De documenten met een lengte van minimaal 28 en maximaal 1492 woorden bleven daarmee over. De gemiddelde documentlengte is daarna 498 woorden en de mediaan is 386 woorden.

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

	Totaal	Vragenuur	Debat
SP	2284	107	2177
CDA	1901	88	1813
D66	1889	133	1756
PvdA	1821	112	1709
PVV	1700	49	1651
VVD	1694	76	1618
ChristenUnie	1068	32	1036
GroenLinks	1068	47	1021
SGP	655	10	645
PvdD	432	14	418
50PLUS	387	12	375
Totaal	14899	680	14219

Deze 14899 documenten zijn verdeeld over 2984 debatten, waarbij elke vraag tijdens het vragenuur als één debat gezien wordt. Op basis van de aantallen is er voor classificatie een baseline *accuracy* van 0.15 (door altijd grootste partij te kiezen) en baseline F_1 score van 0.11 (door te gokken gewogen bij aantal documenten van een partij).

282 3.2 Methoden

283 3.2.1 DV1: Beste classificatiemethode

284 Om deze deelvraag te beantwoorden zullen een aantal classificatiemethoden ver-
285 geleken worden. Aangezien het niet mogelijk is om alle classificatiemethoden
286 te vergelijken, beperkt dit onderzoek zich tot classificatiemethoden die gebruikt
287 zijn in vergelijkbare onderzoeken, zoals besproken in sectie 2.2. Er is ervoor ge-
288 kozen om alleen gebruik te maken van methoden waarvan reeds implementaties
289 beschikbaar waren in scikit-learn. Voor alle methoden wordt gezocht naar de
290 beste parameters, ook wel bekend als een grid search. Deze grid search wordt
291 gedaan door 5-fold cross-validation, waarbij de training set steeds 80% is en de
292 test set 20% van de totale dataset. Een totaal aantal van 6480 combinaties van
293 methoden en parameters zijn getest. De hypothese is dat de scores lager zijn
294 dan die gevonden in het gerelateerd werk, omdat de documentgrootte kleiner is
295 en de baseline scores lager.

296 **Pre-processing** Voor pre-processing is gebruik gemaakt van tokenisation en
297 lowercasing. Voor tokenisation is de reguliere expressie
298 *w+* gebruikt, waardoor alles behalve letters en cijfers weggehaald wordt. Ver-
299 volgens is er gevarieerd tussen wel of geen gebruik maken van stemming. In
300 het geval van stemming is gebruik gemaakt van de Snowball Stemmer van de
301 Python NLTK module.

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

308 Voor dit onderzoek zijn de volgende wegingen voor woorden getest: *boolean*
309 (wel of niet aanwezig), *tf* (woordfrequentie), *tf-norm* (woordfrequentie genor-
310 maliseerd door documentlengte) en *tf-idf* (woordfrequentie gecompenseerd voor
311 documentfrequentie). Daarnaast wordt in dit onderzoek geëxperimenteerd met
312 een minimale of maximale woord- of documentfrequentie. Ook is gekeken naar
313 het effect van combinaties van de volgende n-grams; unigrams, bigrams en tri-
314 grams. N-grams zijn combinaties van N aantal opeenvolgende woorden. Bij
315 een unigram is elke feature gewoon één woord, terwijl bij een bigram dit twee
316 opvolgende woorden zijn. Dit kan nuttig zijn, want als bijvoorbeeld het woord
317 *asfalt* er in voorkomt, dan is het voor ideologie mogelijk relevant of er *minder*
318 *asfalt* of *meer asfalt* staat.

319 **Support Vector Machines en Logistische Regressie** De meest voorko-
320 mende techniek in vergelijkbaar onderzoek is Support Vector Machine (SVM).
321 Een andere techniek die gebruikt wordt is logistische regressie. Beide kennen een
322 eigen implementatie in scikit-learn, maar deze implementaties zijn niet efficiënt
323 met grote datasets. Om deze reden is er in beide gevallen voor gekozen om
324 gebruik te maken van de functie *SGDClassifier*, die beide technieken leert met
325 *stochastic gradient descent learning*. Voor regularisatie is hier geëxperimenteerd
326 met L1 en L2 regularisatie, en een combinatie van beide genaamd Elasticnet.

De andere parameters zijn gelaten op de standaardwaarden van scikit-learn [10]. Een belangrijke onaangepaste waarde is die van maximaal aantal iteraties, waarvoor de scikit-learn standaard 5 is. Volgens scikit-learn convergeert de SGDClassifier rond de $10^6/n$ iteraties waar n het aantal documenten in de training set is. In het geval van deze dataset zou dat 84 iteraties zijn. Vanwege de grootte van de grid search was het voor dit onderzoek niet mogelijk het maximaal aantal iteraties te verhogen tijdens de grid search. De resultaten buiten de grid search zullen gebaseerd zijn op een maximaal aantal iteraties van 100.

Naive Bayes Een andere 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 in voor. Desalniettemin blijkt Naive Bayes effectief te zijn voor tekstclassificatie [6, 10]. Hiervoor zijn de functies van scikit-learn MultinomialNB en BernoulliNB gebruikt [6, 10].

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 worden berekend op basis van vier variabelen. Deze variabelen geven weer hoeveel documenten wel of niet bij een partij horen, en of deze wel of niet als dusdanig zijn geïdentificeerd [11].

	Behorend tot partij	Niet behorend tot partij
Geïdentificeerd als partij	<i>true positive (tp)</i>	<i>false positive (fp)</i>
Niet geïdentificeerd als partij	<i>false negative (fn)</i>	<i>true negative (tn)</i>

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

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

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

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

Accuracy is het percentage van documenten dat correct geïdentificeerd is. *Accuracy* wordt voor de hele classificatie gedaan en niet per klasse. *Precision* is het percentage van documenten geïdentificeerd als een partij, dat ook bij die partij hoort. *Recall* is het percentage documenten van documenten behorende tot een partij, dat ook als die partij geïdentificeerd is. F_1 is het harmonisch gemiddelde van *recall* en *precision*. *Precision*, *recall* en daarmee F_1 worden per partij berekend. Er zijn drie varianten om deze scores voor de hele classificatie te berekenen.

Allereerst is er *micro*, waarbij alle variabelen bij elkaar opgeteld worden en vervolgens de scores berekend. Dit leidt ertoe dat resultaten van partijen met

362 veel documenten belangrijker zijn. Als een classificatie kleine partijen groten-
363 deels fout classificeert, kan deze score alsnog hoog zijn. In het geval van meer
364 dan twee partijen is dit hetzelfde als *accuracy*.

365 Als tweede is er *macro*, waarbij alle scores per partij berekend worden en
366 daarvan het gemiddelde wordt genomen. Dit leidt er dan weer toe dat resultaten
367 van partijen met weinig documenten net zo belangrijk zijn. Hierdoor kan een
368 classificatie met een laag aantal correct geclassificeerde documenten hoog scoren
369 door vooral kleine partijen goed te classificeren.

370 Als laatste is er dan nog *gewogen*, deze berekent net als *macro* de scores per
371 partij, maar neemt hiervan het gemiddelde gewogen bij het aantal documenten
372 behorend tot een partij. Deze wijkt weinig af van de *micro* variant, tenzij er
373 uitschieters zijn bij partijen.

374 Aangezien *micro* al terugkomt in *accuracy* en het nadeel van *macro* te
375 groot is omdat de partijen nogal variëren in grootte, is gekozen voor *gewogen*
376 F_1 score naast *accuracy*.

377 3.2.2 DV2: Invloed van namen

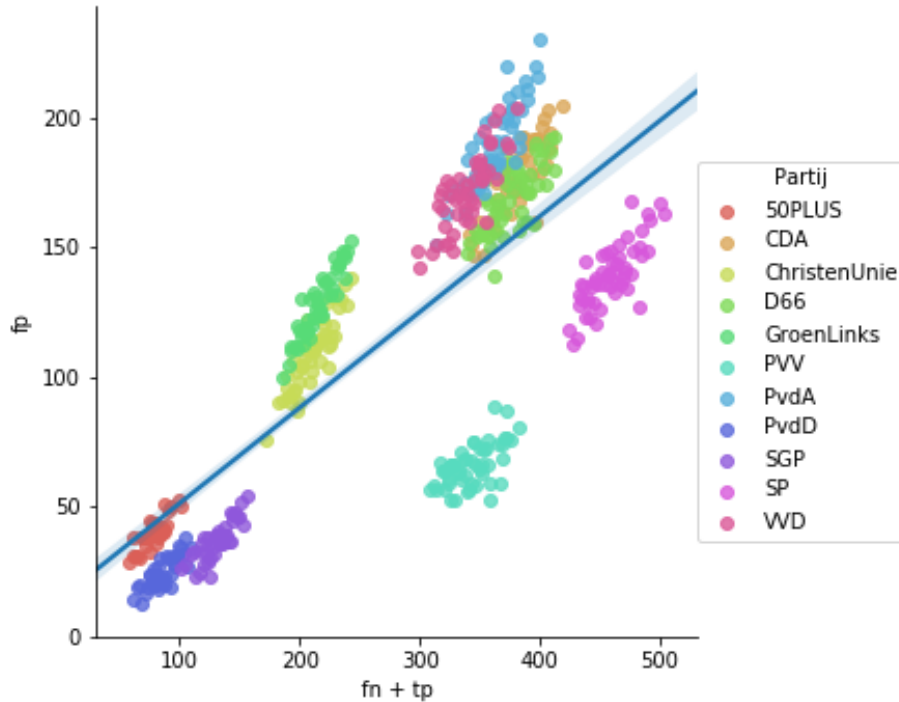
378 In Diermeier et al. [4] wordt aangenomen dat namen een groot effect hebben
379 op de classificatie en Hirst et al. [2] bevestigen dit voor het Europees Par-
380 lement. Aangezien hier bij deelvraag 1 niet voor is gekozen, wordt bij deze
381 deelvraag gekeken hoe groot het effect hiervan is, specifiek gericht op partij-
382 namen en achternamen van Kamerleden. Voor deze deelvraag wordt wederom
383 een classificatie gedaan met de classificatiemethode die resulteerde uit deelvraag
384 1. In deze classificatie worden alle partijnamen vervangen door *PARTIJNAAM*
385 en alle achternamen van Kamerleden vervangen door *KAMERLIDNAAM*. Deze
386 namen zijn uit de Handelingen gehaald. Voor partijnamen zijn ook lidwoorden
387 toegevoegd, voor achternamen van Kamerleden zijn ook verkortingen meegenomen.
388 Dit laatste omdat bijvoorbeeld *Van Haersma Buma* vaak aangesproken
389 wordt als *Buma*. Voornamen van Kamerleden worden zelden tot nooit gebruikt,
390 dus die zijn er niet uitgehaald. Een nadeel van deze aanpak is dat ook namen
391 van niet-Kamerleden of andere woorden weggehaald kunnen worden als deze
392 hetzelfde zijn als naam van een Kamerlid. Door gebruik van gevoeligheid voor
393 hoofdletters is geprobeerd dit te voorkomen. Een opvallend voorbeeld hiervan is
394 de naam Rutte, die zowel behoort tot het Kamerlid Arno Rutte als de premier
395 Mark Rutte. Steekproefgewijs is gekeken of er nog namen achter zijn gebleven,
396 maar die zijn niet gevonden.

397 Ook wordt gekeken naar classificatie met alleen partijnamen en achterna-
398 men van Kamerleden. Alle andere woorden worden weggehaald. Namen van
399 Kamerleden en partijen die niet aan elkaar geschreven worden, zoals *Partij van*
400 *de Arbeid*, worden aan elkaar geschreven zodat het één feature wordt. Doordat
401 alle andere woorden weggehaald zijn, worden de bi- en trigrams combinaties
402 van namen die zinnen uit elkaar kunnen staan, dus die niet meer informatie
403 geven dan unigrams. Daarom wordt er gebruikt van de classificatiemethode uit
404 deelvraag 1, maar dan met alleen unigrams. Hoge scores voor deze classificatie
405 geven aan dat met alleen namen classificatie goed te doen is en dat dit dus een
406 grote bijdrage heeft geleverd aan de resultaten uit deelvraag 1.

3.2.3 DV3: Oppositie of regering

Om deze deelvraag te beantwoorden zal een analyse gedaan worden van de confusion matrix en zullen twee experimenten die gebaseerd zijn op experimenten uit Hirst et al. voor dezelfde vraag uitgevoerd worden op de dataset van de Tweede Kamer. Bij deze deelvraag zal de classificatiemethode uit deelvraag 2 gebruikt worden.

Als er een afhankelijkheid is van partij-status, dan is te verwachten dat het aantal misclassificaties minus verwachte waarde binnen regeringspartijen en binnen oppositiepartijen hoger ligt dan tussen oppositiepartijen en regeringspartijen. De verwachte waarde is afhankelijk van het aantal documenten van een partij in de training set [12]. Aangezien de test set uit dezelfde set als de training is gehaald, is de verwachte waarde ook afhankelijk van het aantal documenten van een partij in de test set. Uit de voorverkenning (op basis van resultaten uit deelvraag 1 en 2) blijkt deze correlatie tussen het aantal *false positives* van een partij en het aantal documenten behorend tot die partij.



Figuur 2: Het aantal *false positives* ten opzichte van het aantal documenten behorend tot die partij (*false negatives* en *true positives*). Dit is op basis van 50 classificaties met verschillende test en train set. De pearson correlatie is 0.78.

Op basis van dit verband is het verwachte aantal documenten

$$V_{i,j} = fn_i * \frac{tp_j + fn_j}{tn_i + fp_i} \quad (5)$$

waar $i \neq j$ met i de echte partij waar een document bijhoort en j de (foutief)

424 voorspelde partij.

425 De error is dan het verschil van de verwachte waarde en het daadwerkelijk
426 aantal documenten

$$e_{i,j} = D_{i,j} - V_{i,j} \quad (6)$$

427 met opnieuw $i \neq j$ en i de voorspelde partij en j de echte partij waar een
428 document bijhoort.

429 Als dit een goede benadering is van de error, dan is het te verwachten
430 dat deze normaal verdeeld is [13]. Om te kijken of er een bias is, worden de
431 distributies binnen regeringspartijen en binnen oppositiepartijen vergeleken met
432 de distributie tussen beide groepen. Om de invloed van variantie door de wil-
433 lekeurige splitsing documenten voor trainen en testen te beperken, wordt de
434 classificatie 100 keer gedaan en worden deze errors bij elkaar in distributie ge-
435 nomen. In het geval dat de distributies normaal verdeeld zijn, zal de statistische
436 test plaatsvinden op basis van een eenzijdige t-toets. Als de distributies niet
437 normaal verdeeld zijn, zal dit plaatsvinden door een Mann-whitneytoets. Het
438 gekozen significantieniveau (α) is 0.05. De nulhypothese is dat er geen verschil
439 is tussen de verdelingen. De alternatieve hypothese is dan dat de distributie
440 van binnen oppositie of regering groter is dan die tussen een regerings- en oppo-
441 sitiepartij. Als de nulhypothese wordt verworpen, kan dus aangenomen worden
442 dat er een verschil is op basis van partij-status.

443 In het eerste experiment uit Hirst et al. zullen de meest karakteristieke
444 woorden per partij van de ene zittingsperiode vergeleken worden met de meest
445 karakteristieke woorden per partij van de andere zittingsperiode. Als de classi-
446 ficatie op basis van ideologie is in plaats van partij-status, is het te verwachten
447 dat de woorden bij een partij blijven en niet gekoppeld zijn aan in oppositie of
448 regering zitten.

449 In het tweede experiment uit Hirst et al. worden classificaties getraind op
450 een zittingsperiode en getest op een andere zittingsperiode. Als de classificatie
451 afhankelijk is van partij-status is de verwachting dat de scores van partijen die
452 gewisseld zijn van oppositie naar regering of andersom lagere scores krijgen dan
453 partijen die niet van partij-status zijn veranderd.

454 Als vergelijkingsmateriaal is voor deze experimenten een tweede dataset
455 nodig uit een ander kabinet. Hiervoor is het wenselijk dat dit kabinet bestaat uit
456 andere partijen dan kabinet-Rutte II. Daarnaast is het ook wenselijk als het niet
457 te ver terug is, zodat onderwerpen en taalgebruik enigszins overeenkomstig zijn.
458 Omdat kabinet-Rutte I een minderheidskabinet was met een bijzondere partij-
459 status voor de PVV, is ervoor gekozen om de Handelingen van de Tweede Kamer
460 tijdens het missionaire kabinet-Balkenende IV (22 februari 2007 tot 20 februari
461 2010) te gebruiken. Dit kabinet bestond uit CDA, PvdA en ChristenUnie.

462 De partij 50PLUS bestond nog niet gedurende kabinet-Balkenende IV,
463 dus documenten van deze partij zijn weggelaten. Verder heeft dezelfde verwer-
464 king van data plaatsgevonden, zoals beschreven in 3.1. Alleen de minimum- en
465 maximumlengte is overgenomen van de dataset van kabinet-Rutte II.

Tabel 2: Aantal documenten per partij gedurende het missionaire kabinet-Balkenende IV.

	Totaal	Vragenuur	Debat
CDA	1039	53	986
ChristenUnie	561	3	558
D66	518	22	496
GroenLinks	760	50	710
PVV	971	56	915
PvdA	903	22	881
PvdD	165	3	162
SGP	507	6	501
SP	1222	70	1152
VVD	1041	55	986
Totaal	7687	340	7347

3.2.4 DV4: Links/rechts

Als de classificatie afhankelijk is van links/rechts positie, dan is het te verwachten dat, net als bij deelvraag 3, de misclassificaties minus de verwachte waarde groter zijn als twee partijen dichtbij elkaar staan op de links/rechts as. Daarvoor zal wederom formule 5 gebruikt worden als verwachte waarde en dus formule 6 als error.

Er zijn verschillende methoden om partijen in te delen op een links/rechts as. Er is hier gekozen voor de indeling van het Manifesto Project [14]. Het Manifesto Project geeft scores op een heel aantal politieke posities, waaronder dus links/rechts, op basis van het verkiezingsprogramma van dat jaar, in dit geval dus van 2012.

Tabel 3: Links/rechts score per partij van het Manifesto Project voor de verkiezingsprogramma's van 2012.

Partij	Score van Manifesto Project
SP	-20.926
GroenLinks	-9.584
PvdA	-6.558
PvdD	-6.465
50PLUS	-6.311
D66	-0.778
ChristenUnie	10.203
PVV	15.642
CDA	17.701
VVD	22.629
SGP	26.6

477 3.2.5 DV5: Woordgebruik van sprekers

478 De vorige classificaties trainden op documenten en werden getest op andere
479 documenten, maar wel van dezelfde sprekers als uit de training set. Naast
480 de ideologie kan de classificatie daarom ook getraind zijn op het taalgebruik
481 van sprekers. Als een Kamerlid bijvoorbeeld een woord regelmatig in speeches
482 gebruikt, maar niet wordt gebruikt door zijn partijgenoten, wordt dit wel gezien
483 als een belangrijk woord voor de classificatie naar partij-affiliatie. Hirst et al.
484 [2] plaatsten al een soortgelijke kanttekening bij de resultaten van Diermeier et
485 al.

486 Om te kijken of dit effect er is, wordt er opnieuw een classificatie gedaan
487 met de methode uit deelvraag 2. Ditmaal worden alleen niet de individuele
488 documenten verdeeld over de training en test set, maar worden de Kamerleden,
489 met bijbehorende documenten, verdeeld over de training en test set. Als taalge-
490 bruik van een spreker in de training set voorheen invloed had op de classificatie,
491 zal dat nu geen effect meer hebben omdat er geen documenten van die spreker
492 meer voorkomen in de test set. De meest karakteristieke woorden uit de resulta-
493 ten van deelvraag 2 suggereren dat woordgebruik van Kamerleden invloed heeft
494 (zie tabel 5). De hypothese is daarom ook dat deze nieuwe classificatie lagere
495 scores vindt.

496 4 Resultaten

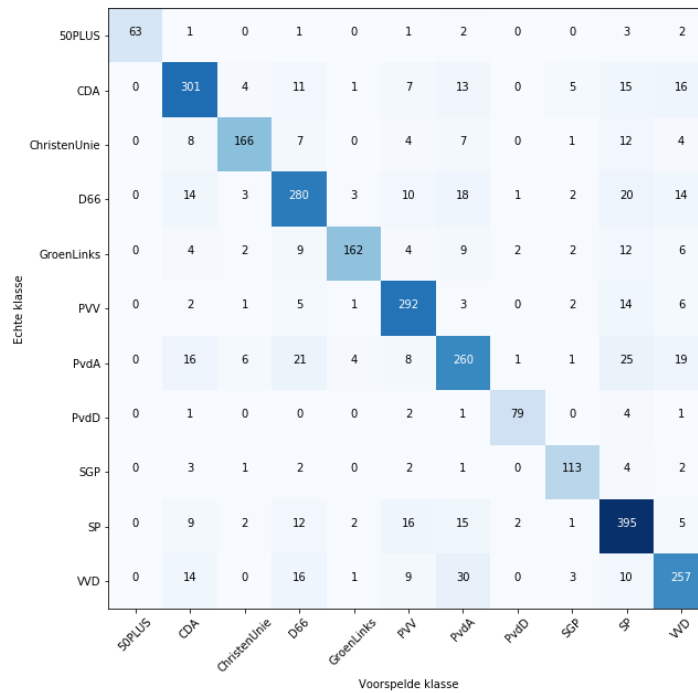
497 4.1 DV1: Beste classificatiemethode

498 Het beste resultaat werd bereikt met Support Vector Machines gebruikmakend
499 van *stochastic gradient descent learning* en Elasticnet regularisatie. De woorden
500 waren hierbij gestemd. De features waren zowel unigrams, bigrams als trigrams.
501 Geen features zijn hierin weggelaten door minimale of maximale documentfre-
502 quenties. Het maximum aantal iteraties was 5 voor de grid search, maar alle
503 resultaten zijn op basis van 100.

504 Tabel 4 laat de scores zien per partij met het aantal documenten in de
505 test set. De *accuracy* voor deze classificatie is 0.80. De F_1 scores per partij
506 liggen tussen de 0.7 en 0.9. De partijen met een sterke focus op één onderwerp,
507 50PLUS, PVV en PvdD, als ook de SGP hebben hoge scores, terwijl de coa-
508 litiepartijen, VVD en PvdA, lagere scores hebben. Figuur 3 laat zien waar de
509 fouten in deze classificatie zitten. De meest karakteristieke features per partij
510 zijn te zien in tabel 5. Met meest karakteristiek worden de n-grams bedoeld die
511 de hoogste coëfficiënt hebben in de classificatie en die dus relatief het meeste
512 belangrijk zijn voor de classificatie van een partij. Hierin is te zien dat vrijwel
513 alle n-grams achternamen van Kamerleden of partijnamen bevatten.

Tabel 4: Classificatie scores per partij van beste classificatiemethode (SVM). Gemiddelde van vijf splitsingen van training en test set. Maximum aantal iteraties is 100.

	Precision	Recall	F_1 score	Documenten
50PLUS	0.97	0.86	0.91	78
PvdD	0.89	0.88	0.89	83
SGP	0.89	0.86	0.88	137
PVV	0.83	0.89	0.86	339
ChristenUnie	0.86	0.78	0.82	220
CDA	0.83	0.80	0.81	376
GroenLinks	0.89	0.73	0.81	203
SP	0.75	0.86	0.80	448
D66	0.76	0.76	0.76	385
VVD	0.75	0.72	0.74	340
PvdA	0.73	0.73	0.73	371
Totaal	0.80	0.80	0.80	2980



Figuur 3: Confusion matrix van beste classificatiemethode (SVM). Gemiddelde van vijf splitsingen van training en test set.

Tabel 5: Meest karakteristieke n-grams per partij op basis van beste classificatie gedurende kabinet-Rutte II. N-grams die niet achternamen van Kamerleden of partijnamen bevatten, zijn dikgedrukt.

50PLUS	CDA	ChristenUnie	D66	GroenLinks
50plus	cda	de christenunie	d66	groenlinks
het lid krol	het cda	christenunie	mijn fractie	lid van tongeren
lid krol naar	cda fractie	lid dik	leden van veldhoven	lid voortman naar
lid krol	de cda fractie	het lid dik	van veldhoven	het lid voortman
krol naar mij	de cda	lid dik faber	veldhoven	lid voortman
krol naar	lid omtzigt naar	dik faber	lid van veldhoven	voortman
krol	lid omtzigt	faber	lid van meenen	tongeren
van 50plus	het lid omtzigt	leden voordewind	d66 wil	van tongeren
gepensioneerden	het cda is	de leden voordewind	d66 is	tongeren naar mij
ouderen	cda is	dik	de leden schouw	van tongeren naar

Tabel 5: Meest karakteristieke n-grams per partij op basis van beste classificatie gedurende kabinet-Rutte II. N-grams die niet achternamen van Kamerleden of partijnamen bevatten, zijn dikgedrukt. (*Vervolg*)

PVV	PvdA	PvdD	SGP	SP	VVD
pvv	de pvda	lid ouwehand	sgp	sp	de vvd
de pvv	pvda	het lid ouwehand	de sgp	de sp	vvd
nederland	de partij van	lid ouwehand naar	sgp fractie	sp fractie	de vvd fractie
islamitische	van de arbeid	ouwehand naar	de sgp fractie	de sp fractie	vvd fractie
miljard	de arbeid	ouwehand naar mij	de leden dijkgraaf	van gerven	de vvd is
het lid graus	partij van de	ouwehand	leden dijkgraaf	gerven	vvd is
lid graus	partij van	dieren	leden van der	lid van gerven	voor de vvd
graus	arbeid	voor de dieren	mevrouw de voorzitter	smaling	wat de vvd
lid graus naar	de partij	de dieren	mevrouw de	leijten naar mij	vvd vindt
klever	pvda fractie	thieme	de leden bisschop	leijten naar	de vvd vindt

514 4.2 DV2: Invloed van namen

515 In tabel 5 was al te zien dat de meest karakteristieke n-grams voornamelijk ach-
516 ternamen van Kamerleden of partijnamen bevatten. In tabel 6 zijn de scores te
517 zien voor een classificatie met alleen achternamen van Kamerleden en partijna-
518 men. De *accuracy* is 0.61. De scores zijn gedaald ten opzichte van de resultaten
519 van deelvraag 1, maar ruim hoger dan de baseline scores.

Tabel 6: Classificatierapport van beste classificatie met alleen achternamen van Kamerleden en partijnamen. Hiervoor is alleen gebruikgemaakt van unigrams. Gemiddelde van vijf splitsingen van training en test set.

	Precision	Recall	F_1 score
50PLUS	0.82	0.88	0.85
PvdD	0.68	0.78	0.69
GroenLinks	0.71	0.66	0.68
PVV	0.66	0.71	0.67
CDA	0.67	0.65	0.66
ChristenUnie	0.66	0.58	0.62
SP	0.61	0.64	0.62
VVD	0.68	0.57	0.62
SGP	0.69	0.54	0.60
D66	0.56	0.53	0.54
PvdA	0.56	0.51	0.52
Totaal	0.64	0.62	0.62

520 In tabel 7 zijn de scores te zien van classificatie met achternamen van
521 Kamerleden en partijnamen vervangen. Deze zijn aanzienlijk lager dan de scores
522 uit deelvraag 1 en ook nog lager dan de scores met alleen namen. Wel zijn de
523 scores nog ruim hoger dan de baseline. In tabel 8 is vervolgens te zien welke
524 n-grams het meest karakteristiek zijn per partij voor deze classificatie.

Tabel 7: Classificatie scores per partij van beste classificatie zonder achternamen van Kamerleden en partijnamen met het relatieve verschil ten opzichte van tabel 4. Gemiddelde van vijf splitsingen van training en test set.

	Precision	Recall	F_1 score	ΔF_1 score (%)
SGP	0.71	0.73	0.72	-18
PvdD	0.75	0.70	0.72	-19
PVV	0.63	0.80	0.70	-19
ChristenUnie	0.68	0.46	0.55	-21
CDA	0.52	0.53	0.52	-23
SP	0.54	0.71	0.61	-24
D66	0.55	0.55	0.55	-28
VVD	0.54	0.49	0.52	-30
50PLUS	0.86	0.49	0.62	-32
PvdA	0.51	0.48	0.50	-32
GroenLinks	0.64	0.38	0.48	-41
Totaal	0.59	0.58	0.57	-29

Tabel 8: Meest karakteristieke n-grams per partij op basis van classificatie uit deelvraag 1 zonder achternamen van Kamerleden en partijnamen gedurende kabinet-Rutte II.

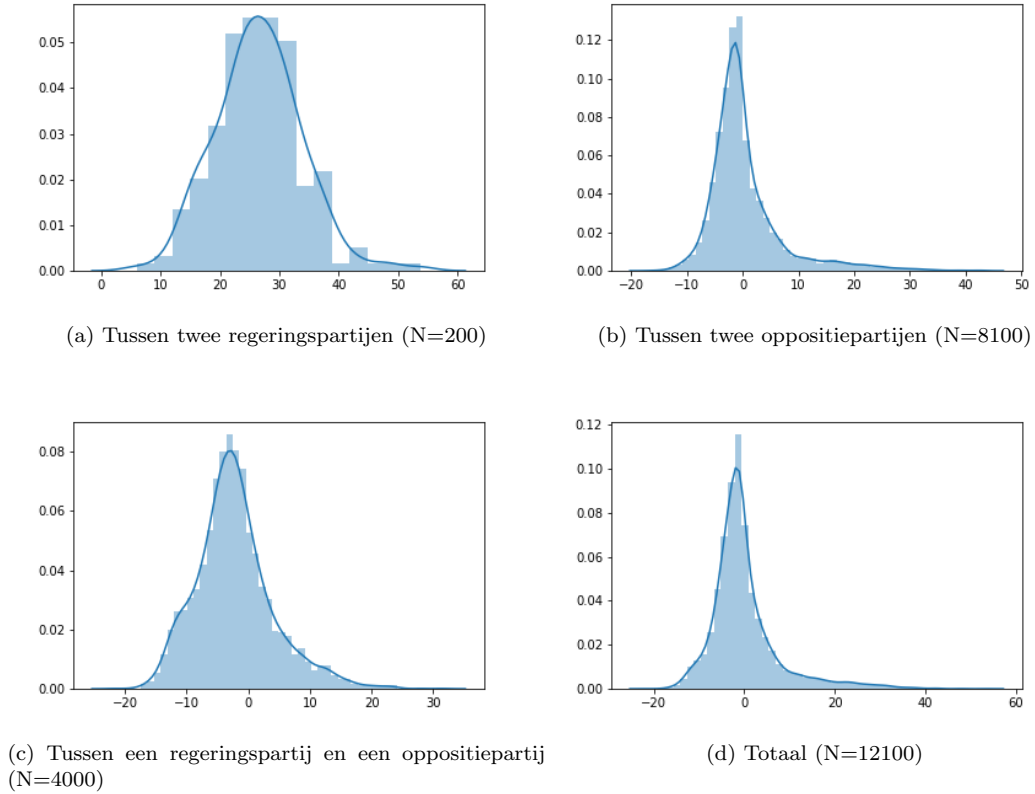
50PLUS	CDA	ChristenUnie	D66	GroenLinks
gepensioneerden	PARTIJ fractie	gezinnen	mijn fractie	zou
ouderen	inwoners	mensenhandel	mijn	kamer hierover te
koopkrachtontwikkeling	regering	inderdaad	natuurlijk	persoonsgebonden
oudere	PARTIJ	onder	fractie	schone energie
plussers	echt	zullen	buitengewoon	in elk geval
50	de regering	horeca	belangrijk	hierover te
werkenden	hier	begeleiding	het kabinet	elk geval
50 plussers	fractie	motie	vandaag	hierover te informeren
voor gepensioneerden	zorginstellingen	gezinnen met	minister	in elk
overwegende dat	wij	ik constateer	kabinet	vluchtelingen

Tabel 8: Meest relevante n-grams per partij op basis van classificatie uit deelvraag 1 zonder achternamen van Kamerleden en partijnamen gedurende kabinet-Rutte II. (*Vervolg*)

PVV	PvdA	PvdD	SGP	SP	VVD
islamitische	en	dieren	mevrouw de voorzitter	zegt	PARTIJ
nederland	jongeren	natuur	mevrouw de	mening dat	speelveld
miljard	vragen	industrie	beantwoording	huurders	aangegeven
islam	open standaarden	bio industrie	punt	armoede	regelgeving
natuurlijk	die	constaterende	voor de beantwoording	van mening dat	volgens mij
al	collega	constaterende dat	de beantwoording	de bevolking	PARTIJ fractie
de islam	daarbij	bio	wel	mensen	PARTIJ is
brussel	kinderen	milieu	allereerst	voorstellen	banen
miljarden	toezeggingen	aarde	bewindslieden	bevolking	ondernemers
dit kabinet	de regering tevens	de bio	je	segregatie	voor PARTIJ

4.3 DV3: Oppositie of regering

In figuur 4 zijn de distributies van de errors, zoals gedefinieerd in formule 6 te zien van combinaties van regerings- en oppositiepartijen.



Figuur 4: Genormaliseerde distributie van de error uit formule 6 voor de verschillende combinaties.

528 Voor alle distributies kan de nulhypothese verworpen worden dat deze
 529 normaal verdeeld zijn, hoewel dit wel verwacht was. In tabel 9 is vervolgens te
 530 zien dat er een significant verschil is tussen de distributies binnen regering en
 531 oppositie tegenover de distributie tussen regering en oppositiepartij.

Tabel 9: Uitslagen van eenzijdige Mann-whitneytoets tussen de distributie tussen een regeringspartij en oppositiepartij en twee distributies. α is 0.05.

	p-waarde	U-waarde
Tussen twee regeringspartijen	3×10^{-124}	717653
Tussen twee oppositiepartijen	1×10^{-93}	16090205

532 In tabel 10 zijn de meest karakteristieke n-grams te zien voor classificatie
 533 van kabinet-Balkenende IV. Hierin zijn geen opvallende overlappen te zien van
 534 regeringspartijen met de classificatie van kabinet-Rutte II in tabel 8.

Tabel 10: Meest karakteristieke n-grams per partij op basis van classificatie uit deelvraag 2 gedurende kabinet-Balkenende IV.

CDA	ChristenUnie	D66	GroenLinks	PVV
PARTIJ fractie	fractie van PARTIJ	premier	PARTIJfractie	burgers
wij	de fractie van	ik hoop	premier	door
fractie	de fractie	de premier	fractie van PARTIJ	deze
wij hebben	fractie van	arbeidsmarkt	de fractie	gewoon
KAMERLID	mijn fractie	hoop	de fractie van	burger
dank	beantwoording	de arbeidsmarkt	fractie van	immigratie
aangegeven	geweest	hij	politieke	niet
zorgvuldige	verschillende	ik	ik	belastinggeld
overleg	van PARTIJ	dadelijk	de premier	onze
ons	moment	schone energie	een beetje	natuurlijk

Tabel 10: Meest karakteristieke n-grams per partij op basis van classificatie uit deelvraag 2 gedurende kabinet-Balkenende IV. (*Vervolg*)

PvdA	PvdD	SGP	SP	VVD
wij	dieren	mijn fractie	mensen	PARTIJ
vrouwen	bio industrie	beantwoording	zegt	PARTIJ fractie
belangrijk	bio	wel	niet	onze fractie
achtergrond	de bio industrie	toch	leraren	fractie
goed	de bio	enkele	is	ondernemers
volgens mij	natuur	bewindslieden	vandaar	want
mbo	dierenwelzijn	de bewindslieden	leerlingen	voorzitter PARTIJ fractie
groep	dierproeven	helder	waarom	justitie
ben	veehouderij	diverse	militaire	antwoorden
alle	industrie	de voorzitter	onderwijs	in elk

535 In tabel 11 zijn de resultaten van de classificatiescores te zien waarbij de
 536 classificatie getraind is op een zittingsperiode, maar getest op een andere. De
 537 resultaten zijn sterk gedaald, maar nog boven de baseline. De daling verschilt
 538 enorm per partij en zittingsperiode met dalingen van F_1 scores tussen 12 en
 539 92%.

Tabel 11: F_1 scores van de classificatie getraind op ene zittingsperiode en getest op andere zittingsperiode. Scores van een classificatie getraind en getest op kabinet-Rutte II zonder 50PLUS zijn bijgevoegd ter referentie, als ook de relatieve daling. De classificatiemethode uit deelvraag 1 is gebruikt zonder achternamen van Kamerleden en partijnamen. Partijen met een asterisk zijn gewisseld van partij-status.

		Training set → Test set			
Rutte II		Balkenende IV → Rutte II Baseline = 0.11		Rutte II → Balkenende IV Baseline = 0.12	
	F_1	F_1	ΔF_1 score (%)	F_1	ΔF_1 score (%)
SGP	0.74	0.56	-24	0.49	-34
PvdD	0.73	0.64	-12	0.45	-38
PVV	0.70	0.50	-29	0.60	-14
SP	0.61	0.41	-33	0.53	-13
ChristenUnie*	0.55	0.37	-33	0.22	-60
D66	0.54	0.16	-70	0.28	-48
CDA*	0.53	0.28	-47	0.43	-19
PvdA	0.52	0.29	-44	0.27	-48
VVD*	0.51	0.18	-65	0.10	-80
GroenLinks	0.49	0.31	-37	0.04	-92
Totaal	0.58	0.34	-41	0.35	-40

540 4.4 DV4: Links/rechts

541 4.5 DV5: Woordgebruik van sprekers

542 In tabel 12 staan de scores van classificatie waarbij de Kamerleden verdeeld zijn
543 over de training en test set. De scores zijn hierbij amper hoger dan de baseline.

Tabel 12: Classificatierapport van beste classificatie met de Kamerleden verdeeld over training en test set. Gemiddelde van vijf splitsingen van training en test set.

	Precision	Recall	F_1 score	ΔF_1 score (%)
50PLUS	0.29	0.06	0.09	
CDA	0.12	0.20	0.14	
ChristenUnie	0.08	0.14	0.09	
D66	0.22	0.22	0.22	
GroenLinks	0.16	0.04	0.05	
PVV	0.29	0.50	0.37	
PvdA	0.25	0.19	0.21	
PvdD	0.46	0.17	0.22	
SGP	0.17	0.05	0.07	
SP	0.34	0.33	0.33	
VVD	0.31	0.26	0.24	
Totaal	0.31	0.24	0.24	

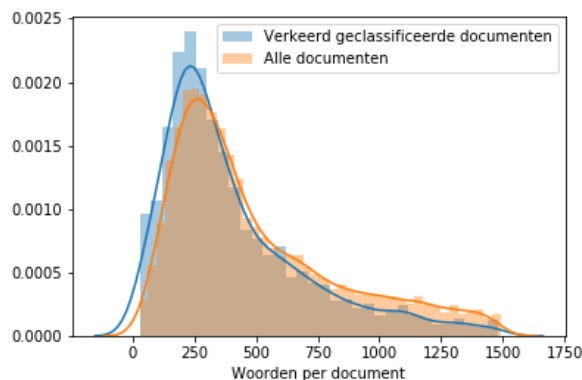
5 Discussie

5.1 DV1: Beste classificatiemethode

Het onderzoek behaalt resultaten in lijn der verwachting op basis van gerelateerd werk en daarnaast ruim boven de baseline scores. De lage scores voor de coalitiepartijen steunen de hypothese van een afhankelijkheid van partij-status, zoals besproken wordt in deelvraag 3. Het bijna alleen voorkomen van namen van partijen en Kamerleden in de meest karakteristieke n-grams per partij in tabel 5 steunt daarnaast het vermoeden dat deze classificatie sterk afhankelijk is van die namen, zoals besproken wordt in deelvraag 2.

Dit onderzoek heeft zich beperkt tot methoden genoemd in vergelijkbare onderzoeken en waarvan de implementatie beschikbaar is in scikit-learn. 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. Daarnaast richtte zich dit ook maar op een beperkt aantal parameterwaarden. Voor vervolgonderzoek kan daarom dit onderdeel uitgebreid worden. Het effect van het beperkte maximum iteraties was bij de beste classificatiemethode 2%.

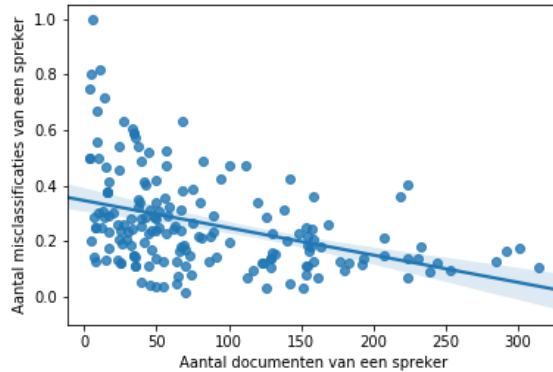
Het onderzoek van Hirst et al. vond dat resultaten afhankelijk kunnen zijn van documentgrootte. Alle documenten in dit onderzoek zijn kleiner dan de grootste documentgrootte uit het onderzoek van Hirst et al. en ook de minimale documentgrootte ligt lager dan de kleinste documentgrootte uit dat onderzoek. Het effect wat zij vinden tussen documentgrootte van 267 en 6666 is een verschil in *accuracy* van 19.8%. Dit onderzoek vindt inderdaad dat kleinere documenten vaker foutief geclassificeerd worden.



Figuur 5: Genormaliseerde distributie van documentlengtes van foutief geclassificeerde documenten en alle documenten. Totaal van 5-fold cross-validation, waardoor documenten vaker voor kunnen komen. Mediaan documentlengte van foutief geclassificeerde documenten is 321 en voor alle documenten 386.

Voor een vervolgonderzoek kan uitgebreider gekeken worden naar dit effect en wat dit betekent voor de resultaten. Het percentage documenten van een vragenuur is tweemaal zo hoog bij foutief geclassificeerde documenten, maar dit lijkt te komen doordat deze documenten vaak kleiner zijn (mediaan is 286).

Er is verder nog gekeken naar andere verbanden tussen documenten die verkeerd zijn geclassificeerd. Daarbij is nog te zien dat sprekers met weinig documenten relatief iets meer voorkomen in verkeerd geclassificeerde documenten.



Figuur 6: Aantal misclassificaties gedeeld door totaal aantal documenten per spreker tegenover totaal aantal documenten van een spreker. Misclassificaties zijn totaal van 5-fold cross-validation, waardoor documenten vaker mee kunnen tellen. De pearson correlatie is -0.28 en de p-waarde 1.07×10^{-4} .

Dit versterkt het vermoeden dat de classificatie mede plaatsvindt op basis van woordgebruik van individuele sprekers, zoals besproken wordt in deelvraag 5.

5.2 DV2: Invloed van namen

De resultaten laten zien dat de classificatie sterk afhankelijk is van partijnamen en achternamen van Kamerleden. De hogere scores voor de classificatie met alleen namen dan zonder namen in combinatie met de woorden in tabel 5 suggereert dat dit het belangrijkste was in de classificatie van deelvraag 1. Deze daling was te verwachten op basis van gerelateerd werk.

De n-grams in tabel 8 komen bij veel partijen overeen met hun ideologie, vooral bij de partijen met een sterke focus op één onderwerp; PVV, PvdD en 50PLUS. Daarnaast zijn er ook n-grams die niet veel over ideologie lijken te zeggen, zoals; *volgens mij*, *ik constateer* en *in elk geval*. Vooral de SGP heeft n-grams die niet veel lijken te zeggen over de ideologie, hoewel deze partij desalniettemin de hoogste F_1 score heeft. Met name opvallend hierbij is *mevrouw de voorzitter*, aangezien deze woorden door alle partijen gebruikt worden om via de voorzitter te praten. Voor een vervolgonderzoek kan gekeken naar waarom deze n-grams zo karakteristiek zijn voor partijen. Een hypothese is dat deze n-grams eigen zijn aan een individueel Kamerlid.

De classificatiemethode die gebruikt is in deze deelvraag, is gebaseerd op de beste methode voor de dataset uit deelvraag 1. Hierin was gevonden dat een combinatie van uni-, bi- en trigrams het beste resultaat opleverde. In tabel 5 is te zien dat trigrams behoren tot de meest karakteristieke n-grams, hoewel de woorden in trigrams vaak overlappen met uni- en bigrams. In tabel 8 daarentegen zijn er nog maar een paar trigrams, welke grotendeels procedurele zinnen

600 zijn of toevoeging van een lidwoord op een uni- of bigram. Dit verschil sugge-
601 reert dat trigrams minder belangrijk zijn in de classificatie zonder de namen,
602 dus de classificatiemethode uit deelvraag 1 niet het beste is voor deze classifica-
603 tie. In vervolgonderzoek kan de opzet van deelvraag 1 toegepast worden op de
604 classificatie zonder de namen, om zo te komen tot een classificatiemethode die
605 het beste resultaat oplevert op de classificatie zonder namen.

606 Er is ook gekeken naar andere namen in de lijst van 100 meest karakte-
607 ristieke woorden per partij, zoals van gebieden, bedrijven of bewindspersonen.
608 Bewindspersonen komen hier niet in voor. Er komen een aantal gebieden in
609 voor, zoals *aruba*, *limburg* en *saoedi arabië*. Ook komen er organisaties als *gvo*
610 *hvo* en *monsanto* in voor. Deze woorden lijken in sommige gevallen een weer-
611 spiegeling te zijn voor ideologie, dus voor vervolgonderzoek lijkt het niet nodig
612 te zijn deze te verwijderen.

613 5.3 DV3: Oppositie of regering

614 In tabel 4 is het opvallend dat de coalitiepartijen lage scores krijgen. Daarnaast
615 laat figuur 3 zien dat er een hoge overlap zit tussen deze twee partijen.

616 De statistische toetsresultaten in tabel 9 laten zien dat inderdaad de error
617 groter is binnen oppositie of regering dan tussen een regerings- en oppositiepar-
618 tij. Dit suggereert dat inderdaad partij-status invloed heeft op de classificatie.

619 De verwachting was dat de error normaal verdeeld zou zijn. De verde-
620 lingen uit figuur 4 hebben globaal wel de vorm van een normaal verdeling. In
621 figuur 2 is het daarnaast opvallend dat partijen zoals SP en PVV ruim onder de
622 regressielijn zitten, terwijl andere partijen er een stuk boven zitten. Dit geeft
623 aanleiding te vermoeden dat er naast het aantal documenten van een partij
624 nog meer factoren van invloed zijn op het aantal misclassificaties en daarmee
625 de verwachte waarde. En deze verwachte waarde en de daar uit volgende er-
626 ror zijn een belangrijke aanname van deze methode. Voor deze methode is het
627 dus belangrijk uit te vinden of dit een goede benadering is van de verwachte
628 waarde. In deelvraag 4 wordt gekeken of links/rechts positie hier nog invloed
629 heeft. Voor een vervolgonderzoek kan nog verder gekeken worden naar invloeden
630 op verwachte waarde of andere confounding biases.

631 De overlap van 100 meest karakteristieke n-grams tussen regeringspartijen
632 die niet voorkomen bij oppositiepartijen gedurende kabinet-Rutte II beperkt
633 zich tot de woorden *en* en *blij*, als ook *toezegging* voor VVD en *toezeggingen*
634 voor PvdA.

Tabel 13: N-grams die bij minimaal één regeringspartij in beide kabinetten voorkomen in de 100 meest karakteristieke n-grams, maar niet voor één van de twee partijen tijdens het andere kabinet.

		Kabinet-Rutte II	
		PvdA	VVD
Kabinet-Balkenende IV	CDA	<i>toezeggingen hun collega KAMERLID in aanpak collega</i>	<i>algemeen algemeen overleg toezegging helder overleg aangegeven voor voor PARTIJ</i>
	ChristenUnie	<i>mijn waarop blij collega KAMERLID erg</i>	<i>gaan termijn blij met de volgens volgens mij blij beantwoording</i>
	PvdA		<i>volgens volgens mij</i>

Hoewel er een aantal overeenkomsten zijn qua meest karakteristieke n-grams tussen regeringspartijen van de twee kabinetten, lijkt dit beperkt. De meeste overeenkomsten lijken daarnaast niet heel inhoudelijk gerelateerd aan partij-status. Deze resultaten suggereren daarom ook maar een beperkte invloed van partij-status op de classificatie. Voor een vervolgonderzoek kan uitgebreider gekeken worden naar de overlappende meest karakteristieke n-grams en wat deze zeggen over een regeringspartij.

De scores in tabel 11 laten een duidelijke daling zien ten opzichte van een classificatie van alleen kabinet-Rutte II. Deze algemene daling kan verklaard worden door verschuiving in ideologie, verschil in woordgebruik, verandering van onderwerpen en/of verandering in aantal documenten per partij. De daling is het grootst bij VVD, maar valt mee bij de twee andere partijen die gewisseld zijn van partij-status, ChristenUnie en CDA. Daarnaast is de daling ook heel sterk bij oppositiepartijen GroenLinks en D66, alsook de regeringspartij in beide kabinetten, PvdA. Dat de daling niet consequent groter is bij partijen die gewisseld zijn van partij-status, suggereert dat de invloed van partij-status beperkt is op de classificatie.

Dat de experimenten uit Hirst et al. in hun onderzoek wel invloed vinden, maar in dit onderzoek niet kan komen doordat hun onderzoek zich richt op binaire classificatie, terwijl dit onderzoek meerdere partijen heeft. Zo kan het ontbreken van gemeenschappelijke n-grams komen doordat regeringspartijen zich ook van elkaar moeten onderscheiden in dit onderzoek, waarvoor n-grams die relevant zijn voor partij-status weinig effect hebben, terwijl in het onderzoek van Hirst et al. de regeringspartij alleen onderscheiden hoeft te worden van de oppositiepartij. Daarnaast verklaren zij dat een daling tussen twee zittingsperi-

660 oden met een wisseling van partij-status het gevolg is van deze wisseling, terwijl
661 in dit onderzoek gekeken kan worden naar dit effect voor partijen die wel en
662 niet gewisseld zijn.

663 **5.4 DV4: Links/rechts**

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

666 **5.5 DV5: Woordgebruik van sprekers**

667 De resultaten uit tabel 12 zijn laag, amper hoger dan de baseline. Dit sugge-
668 reert inderdaad dat eerdere classificaties in grote mate toch afhankelijk waren
669 van het woordgebruik van sprekers. Dit is opmerkelijk aangezien vergelijkbare
670 onderzoeken dit effect niet vinden. De meest karakteristieke n-grams van deze
671 classificatie wijken daarnaast grotendeels niet af van die uit tabel 8.

672 Een alternatieve verklaring is dat de classificatie nu mede op basis van
673 woordvoerderschap is. Per onderwerp heeft een partij vaak maar één woord-
674 voerder, met uitzonderingen van wijzigingen in de fractie. Het is aannemelijk
675 dat het taalgebruik afhankelijk is van woordvoerderschap, aangezien er andere
676 termen gebruikt worden bij bijvoorbeeld een debat over zorg dan bij een debat
677 over onderwijs. Als een woordvoerder op een bepaald onderwerp van een partij
678 in de test set voorkomt, is er een grote kans dat geen enkele spreker van die
679 partij eerder over dat onderwerp heeft gepraat, want de woordvoerder gaat nou
680 eenmaal daarover. Daardoor heeft deze spreker veel n-grams die ook voorko-
681 men bij andere woordvoerders over dat onderwerp, maar van andere partij. Als
682 deze n-grams ook belangrijk zijn voor de classificatie kan het zijn dat de woord-
683 voerder geclassificeerd wordt bij een partij van een andere woordvoerder. Een
684 vervolgonderzoek kan kijken of dit een verklaring is.

685 Vergelijkbare onderzoeken vermijden dit mogelijke probleem door alle spreek-
686 beurten van een spreker samen te voegen tot één document. Zoals al eerder
687 vermeld is dit onpraktisch voor de kleinere partijen. Voor een vervolgonderzoek
688 kan desalniettemin gekeken worden naar deze methode om te kijken of dat wel
689 een weerspiegeling is van ideologische verschillen.

690 **5.6 Algemeen**

691 Het vergelijken van deze resultaten met vergelijkbaar werk is ingewikkeld, aan-
692 gezien de keuzes en eigenschappen van die onderzoeken het niet een één-op-één
693 vergelijking maken. Voorbeelden hiervan zijn de taal, het parlement, de do-
694 cumentgrootte, baselines, behouden of weglaten van namen, een spreker als
695 document zien en het trainen en testen op dezelfde spreker. Hoewel de re-
696 sultaten in sommige gevallen lager zijn dan die uit vergelijkbaar werk, is het
697 belangrijk hier rekening mee te houden. Een vervolgonderzoek zou daarom dit
698 onderzoek kunnen reproduceren op een ander parlement om daarmee te kunnen
699 vergelijken.

700 Dit onderzoek richtte zich hoofdzakelijk op de Handelingen gedurende
701 kabinet-Rutte II. Om te kijken in hoeverre het mogelijk is om deze conclu-
702 sie door te trekken naar de algemene Handelingen van de Tweede Kamer, kan

er in vervolgonderzoek gekeken worden naar meerdere zittingsperioden. Ook kan gekeken worden naar veranderingen als een kabinet demissionair is.

Dit onderzoek heeft een aantal beperkingen die in dit hoofdstuk besproken zijn. Het uitvoeren van deze aanbevelingen kan de validiteit en betrouwbaarheid van dit onderzoek vergroten. Ook is dit onderzoek moeilijk te vergelijken met andere onderzoeken om diverse redenen, maar vooral ook omdat het toegepast is op een ander parlement. Desalniettemin geeft dit onderzoek reden om te twijfelen aan de bruikbaarheid van tekstclassificatie van de Handelingen van de Tweede Kamer voor een relatie tussen woordgebruik en ideologie. Daarnaast levert dit onderzoek ook kritieken op een aantal vergelijkbare onderzoeken.

6 Conclusies

Dit onderzoek vindt een *accuracy* en F_1 score van 0.80 voor het classificeren van spreekbeurten in de Tweede Kamer naar partij-affiliatie. De baseline scores zijn respectievelijk 0.11 en 0.15. Als rekening wordt gehouden met partijnamen en achternamen Kamerleden daalt de *accuracy* naar 0.58 en de F_1 score naar 0.57. Dit onderzoek vindt aanwijzingen dat deze classificatie afhankelijk is van de partij-status (oppositie of regering). Als rekening wordt gehouden met woordgebruik van individuele Kamerleden, daalt de nauwkeurigheid verder naar.... Hoewel dit onderzoek hoge scores vindt voor classificatie, lijken deze in grote mate afhankelijk te zijn van andere factoren dan ideologie.

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? In Bertie Kaal, Isa Maks, and Annemarie van Elfrinkhof, editors, *From Text to Political Positions*, chapter 5, pages 93–115. John Benjamins Publishing Company, Amsterdam, 2014.
- [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.
- [5] Bei Yu, Stefan Kaufmann, and Daniel Diermeier. Classifying party affiliation from political speech. *Journal of Information Technology & Politics*, 5(1):33–48, 2008.
- [6] Maneesh Bhand, Dan Robinson, and Conal Sathi. Text classifiers for political ideologies, 2009.
- [7] Bjørn Høyland, Jean-François Godbout, Emanuele Lapponi, and Erik Vell-dal. Predicting party affiliations from european parliament debates. In *Proceedings of the ACL 2014 Workshop on Language Technologies and*

- 743 *Computational Social Science*, pages 56–60. Association for Computatio-
744 nal Linguistics, 2014.
- 745 [8] Laura Klompenhouwer. Extra ledenvergadering 50plus om splitsing. *NRC*
746 *Handelsblad*, June 2014.
- 747 [9] Eric Jones, Travis Oliphant, Pearu Peterson, et al. SciPy: Open source
748 scientific tools for Python, 2001.
- 749 [10] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Gri-
750 sel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas,
751 A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay.
752 Scikit-learn: Machine learning in Python. *Journal of Machine Learning*
753 *Research*, 12:2825–2830, 2011.
- 754 [11] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. *Intro-*
755 *duction to Information Retrieval*. Cambridge University Press, New York,
756 NY, USA, 2008.
- 757 [12] Mahendra Sahare and Hitesh Gupta. A review of multi-class classifica-
758 tion for imbalanced data. *International Journal of Advanced Computer*
759 *Research*, 2(3), 2012.
- 760 [13] *NIST/SEMATECH e-Handbook of Statistical Methods*. NIST/SEMA-
761 TECH, April 2012.
- 762 [14] Andrea Volkens, Pola Lehmann, Theres Matthieß, Nicolas Merz, Sven Re-
763 gel, and Bernhard Weßels. The manifesto data collection. manifesto project
764 (mrg/cmp/marpor). version 2017b, 2017.