

1 INVLOED VAN IDEOLOGIE BEPERKT OP
2 TEKSTCLASSIFICATIE IN 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

13

Samenvatting

In verschillende onderzoeken zijn parlementaire teksten geclassificeerd naar partij-affiliatie. Dit onderzoek heeft eerst gezocht naar de beste classificatiemethode voor het Nederlandse parlement. Vervolgens is gekeken in hoeverre de classificatie het gevolg is van ideologie. Hiervoor is gekeken naar de invloed van namen, in regering of oppositie zitten, positie op de links-rechts as en woordgebruik van sprekers.

De beste classificatiemethode met een nauwkeurigheid van 0.80 is Support Vector Machines. Dit daalt naar 0.58 als achternamen van Kamerleden en partijnamen weggehaald worden. Het onderzoek vond ook aanwijzingen dat de classificatie afhankelijk is van of een partij in regering of oppositie zit. Aanwijzingen voor afhankelijkheid van positie op links-rechts as zijn daarentegen niet gevonden. Als laatste daalt de nauwkeurigheid verder naar 0.27 als Kamerleden verdeeld worden over de training en test set, wat suggereert dat de oorspronkelijke classificatie afhankelijk was van woordgebruik van sprekers. Dit leidt tot de conclusie dat in grote mate de classificatie niet het gevolg is van ideologie.

31	Inhoudsopgave	
32	1 Introductie	3
33	2 Gerelateerd werk	4
34	2.1 Tekstclassificatie van parlementaire teksten	4
35	2.2 Classificatiemethoden	6
36	2.3 Invloed van partijnamen of sprekersnamen	6
37	2.4 Invloed van oppositie of regering	7
38	3 Methodologie	7
39	3.1 De data	7
40	3.2 Methoden	9
41	3.2.1 DV1: Beste classificatiemethode	9
42	3.2.2 DV2: Invloed van namen	11
43	3.2.3 DV3: Oppositie of regering	12
44	3.2.4 DV4: Links-rechts as	15
45	3.2.5 DV5: Woordgebruik van sprekers	16
46	4 Resultaten	16
47	4.1 DV1: Beste classificatiemethode	16
48	4.2 DV2: Invloed van namen	19
49	4.3 DV3: Oppositie of regering	21
50	4.4 DV4: Links-rechts as	24
51	4.5 DV5: Woordgebruik van sprekers	25
52	5 Discussie	26
53	5.1 DV1: Beste classificatiemethode	26
54	5.2 DV2: Invloed van namen	28
55	5.3 DV3: Oppositie of regering	29
56	5.4 DV4: Links-rechts as	31
57	5.5 DV5: Woordgebruik van sprekers	31
58	5.6 Algemeen	32
59	6 Conclusies	32

60 1 Introductie

61 Teksten van politieke partijen kunnen dienen als bron voor het bepalen van
62 ideologische positie van andere teksten, aangezien zij zowel tekst bevatten als
63 ook een bekende ideologie in de vorm van een partij van de spreker; de partij-
64 affiliatie. Het classificeren op basis van tekst kan inzichten geven over de relatie
65 tussen ideologie en woordgebruik. Deze informatie kan vervolgens toegepast
66 worden op andere teksten die wellicht ideologisch van aard zijn. Bijvoorbeeld
67 kan men aan de hand van deze informatie teksten uit kranten classificeren op
68 basis van ideologie [2, 6].

69 In diverse landen zijn al onderzoeken gedaan naar het classificeren naar
70 partij-affiliatie op basis van speeches in parlementen [1, 2, 3, 4, 6, 7, 15]. Met deze
71 tekstclassificatie naar partij-affiliatie proberen onderzoekers zo goed mogelijk te
72 classificeren. Daarnaast proberen ze vaak ook uit te vinden in hoeverre ideologie
73 terug te vinden is in teksten van politici. De resultaten van de classificaties zijn
74 in de meeste gevallen ruim boven de baseline. Hirst et al. [6] vonden voor het
75 Europese Parlement aanwijzingen dat dit het gevolg was van afstand van op de
76 links-rechts as. Er zijn in deze onderzoeken ook redenen die suggereren dat dit
77 niet alleen het gevolg is van ideologie. Zo suggereren de resultaten van Hirst
78 et al. op het Canadese parlement dat de partij-status (oppositie of regering)
79 van invloed is op de classificatie. Daarnaast laat hun onderzoek naar Europese
80 parlement ook zien dat partijnamen een grote invloed hebben op de classificatie.

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

83 Dit onderzoek richt zich daarom op meerdere classificatiemethoden. Daar-
84 naast zal dit onderzoek zich richten op de Tweede Kamer. De onderzoeksvraag
85 luidt daarom dus ook: "In hoeverre is classificatie naar partij-affiliatie aan de
86 hand van spreekbeurten in de Tweede Kamer het gevolg van ideologie?"

87 Deze vraag wordt beantwoord door de antwoorden te vinden op de vol-
88 gende deelvragen:

- 89 1. Wat is de beste classificatiemethode voor classificatie naar partij-affiliatie
90 in de Tweede Kamer en wat is het resultaat van dit model?
- 91 2. In hoeverre is deze classificatie afhankelijk van achternamen van Kamer-
92 leden en partijnamen?
- 93 3. In hoeverre wordt deze classificatie bepaald door partij-status (oppositie
94 of regering)?
- 95 4. In hoeverre wordt deze classificatie bepaald door positie op de links-rechts
96 as?
- 97 5. In hoeverre wordt deze classificatie bepaald door woordgebruik van spre-
98 kers?

99 Hirst et al. [6] vonden dat voor het Canadese parlement de partij-status van
100 invloed was op de classificatie. In datzelfde onderzoek werd bij het Europese
101 parlement geconstateerd dat ook partijnamen en positie op links-rechts as be-
102 palend zijn. Ook levert dit onderzoek kritiek op een onderzoek van Diermeier
103 et al. [3] waar getraind wordt op dezelfde sprekers als waar op getest wordt. Op

104 basis van dit onderzoek is de hypothese dat al deze factoren van invloed zijn op
105 de classificatie.

106 Voor de eerste deelvraag is Support Vector Machine, logistische regressie
107 en Naive Bayes met verschillende parameters vergeleken aan de hand van *accu-*
108 *racy* en F_1 score. Bij de tweede deelvraag is gekeken naar classificatie zonder
109 achternamen van Kamerleden en partijnamen of met alleen achternamen van
110 Kamerleden en partijnamen. De derde vraag bestaat uit drie experimenten.
111 In de eerste is gekeken naar de hoeveelheid misclassificaties binnen regerings-
112 partijen of binnen oppositiepartijen tegenover tussen een regeringspartij en een
113 oppositiepartij. In de tweede is gekeken naar overlap in woordgebruik binnen
114 regering. In de derde is gekeken naar verschil in scores als een partij gewisseld
115 is van partij-status. Bij de vierde vraag is gekeken naar een verband tussen
116 misclassificaties en afstand tussen twee partijen op de links-rechts as. Bij de
117 vijfde vraag is de classificatie herhaald met Kamerleden verdeeld over training
118 en test set.

119 **Overzicht van scriptie** Sectie 2 bevat vergelijkbare onderzoeken in andere
120 parlementen. Sectie 3 bevat de methodologie van de verschillende deelvragen.
121 Sectie 4 bevat de resultaten. Sectie 5 bevat de evaluatie van zowel de resultaten
122 als de methodologie. Sectie 6 bevat ten slotte het antwoord op de onderzoeks-
123 vraag.

124 2 Gerelateerd werk

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

129 In dit hoofdstuk zullen verschillende onderzoeken behandeld worden die
130 tekstclassificatie hebben toegepast op parlementaire teksten. Eerst zullen de
131 onderzoeken algemeen besproken worden. Vervolgens is uitgebreider gekeken
132 worden de effecten van verschillende classificatiemethoden. In de latere secties
133 worden de aspecten besproken die in vergelijkbare onderzoeken genoemd worden
134 als van invloed op de classificatie.

135 2.1 Tekstclassificatie van parlementaire teksten

136 Diermeier et al. deden onderzoek naar het classificeren op basis van ideologische
137 positie in de Amerikaanse Senaat [3]. Ze trainden hun classificatie op de speeches
138 van de 25 meest liberale en de 25 meest conservatieve senatoren van het 101e
139 tot en met het 107e Congres en testten op dezelfde categorieën van het 108e
140 Congres. Een document was in dit onderzoek de verzameling van alle speeches
141 van een senator in een Congres. Deze classificatie resulteerde uiteindelijk in een
142 nauwkeurigheid van 94% (baseline van 50%). Van de 50 senatoren in de test
143 set, kwamen er 44 al voor in de training set, doordat de training op voorgaande
144 Congressen was.

145 Later in het onderzoek vergeleken ze ook de 25 gematigd conservatieve
146 en de 25 gematigd liberale senatoren van dezelfde Congressen. Het resultaat
147 hiervan was 52% (baseline van 50%), dus nauwelijks beter dan de baseline.

148 Als verklaring voor dit verschil ten opzichte van de uitersten zeggen ze dat
149 gematigden een minder duidelijke ideologie hebben.

150 Yu et al. [15] richtten zich vervolgens op zowel het Amerikaanse Huis van
151 Afgevaardigden als de Senaat in 2005. Een document was in dit onderzoek de
152 verzameling van alle speeches van een congreslid en het label de partij. Voor
153 het Huis van Afgevaardigden vonden ze een nauwkeurigheid van 80.1% (baseline
154 van 51.5%) en voor de Senaat 86.0 % (baseline van 55.0%). Ze testten hun
155 classificaties ook op de andere kamer. Van Huis van Afgevaardigden naar Senaat
156 leverde dit een nauwkeurigheid op van 88.0% (baseline van 55.0%) en andersom
157 67.6% (baseline van 51.5%). Hun verklaring voor dit verschil was dat het Huis
158 van Afgevaardigden sterker verdeeld is langs partijlijnen.

159 Vervolgens herhaalden ze de classificaties op het Huis van Afgevaardigden
160 uit 2015, maar testten ditmaal op de Senaat elk jaar tussen 1989 en 2006 af-
161 zonderlijk. Hierin was een stijging in nauwkeurigheid van 60.0% (baseline van
162 55.0%) in 1989 naar 87.0% (baseline van 55.0%) in 2006 te zien, maar met twee
163 duidelijke dalen. Ze presenteren twee mogelijke verklaringen voor de trend; het
164 veranderen van de onderwerpen en het sterker verdeeld worden van het Congres.

165 Als een vervolg op deze onderzoeken deden Hirst et al. een vergelijkbaar
166 onderzoek naar het Canadese Parlement [6]. Hierbij werd zowel gekeken naar de
167 Engelse als Franse teksten. Een document werd hier gezien als de samenvoeging
168 van alle spreekbeurten van een spreker. Afhankelijk van taal en dataset vonden
169 zij in dit onderzoek nauwkeurigheid scores van 83.2% en hoger (baseline van
170 65.5%).

171 Het onderzoek bevat ook een classificatie van het Europees Parlement.
172 Hierbij voegden ze alle teksten van een parlements lid bij elkaar en deelden die
173 op in documenten van gelijke grootte. Voor documentgrootte van 267 woorden
174 werd een nauwkeurigheid van 44.0% gevonden oplopend tot 61.8% (baseline van
175 38-39%) voor documentgrootte van 6666.

176 Bhand et al. [1] richtten zich op het classificeren van leden van het Ameri-
177 kaanse Congres in 2005, op basis van partij-affiliatie (Republikeins of Democra-
178 tisch). Een document hierbij was in tegenstelling tot eerdergenoemde onderzoe-
179 ken een speech. Zij vonden hiervoor uiteindelijk een F_1 score van 0.68 (baseline
180 niet vermeld).

181 Ferreira [4] probeerde interventies van politici te classificeren op basis van
182 geslacht, leeftijdsgroep, partij-affiliatie en oriëntatie in het Portugese parlement.
183 In het geval van classificatie op basis van partij-affiliatie bereikte men een F_1
184 score van 0.90 (baseline niet vermeld, zes partijen).

185 Høyland et al. trainden een classificatie voor partij-affiliatie op basis van
186 teksten van het vijfde Europese Parlement (1999-2004) en testten vervolgens
187 op het zesde Europese Parlement (2004-2009) [7]. Alle teksten van een spreker
188 waren samengevoegd tot één document. 40% van de sprekers in de test set zaten
189 ook in de training set. Hier werd een macro F_1 score van 0.464 (baseline van
190 0.097) en nauwkeurigheid van 0.551 (baseline van 0.410) verkregen. De baseline
191 is in dit onderzoek op basis van altijd classificeren als grootste partij, terwijl
192 voor F_1 score de baseline hoger ligt als hiervoor gekozen wordt voor gokken
193 gewogen bij grootte van een klasse.

194 2.2 Classificatiemethoden

195 Diermeier et al. [3] gebruikten Support Vector Machines. Verder maakten ze
196 gebruik van *tf-idf* met een minimale woordfrequentie van 50 en een minimale
197 documentfrequentie van 10 en *Part-Of-Speech tagging*.

198 Yu et al. [15] maakten gebruik van Support Vector Machines en Naive
199 Bayes, waarvan de varianten multinomial en Bernoulli. De features waren unig-
200 rams, met minimale woordfrequentie van drie en de top 50 meest voorkomende
201 woorden weggelaten. Voor de wegingen van de features bij Support Vector Ma-
202 chines werd geëxperimenteerd met *boolean*, *tf-norm* en *tf-idf*. De beste classifica-
203 tiemethode was afhankelijk van de dataset. Voor het Huis van Afgevaardigden
204 was het Support Vector Machines met als weging *tf-idf* en voor de Senaat Ber-
205 nouilli Naive Bayes.

206 Hirst et al. [6] maakten gebruik van Support Vector Machines. Ze experi-
207 menteerden met verschillende vormen van pre-processing, inclusief stemmen en
208 het verwijderen van woorden op basis van te hoge of te lage frequentie. Deze
209 variaties maakten in hun onderzoek geen grote verschillen en uiteindelijk is ge-
210 kozen voor het niet stemmen, het weglaten van woorden die in minder dan
211 vijf documenten voorkomen en resultaten van zowel met als zonder de top 500
212 meest frequente woorden. Daarnaast werd geëxperimenteerd met vier wegingen
213 voor woorden: *boolean*, *tf*, *tf-norm* en *tf-idf*, waarvan *tf-idf* het beste resultaat
214 opleverde.

215 Bhand et al. [1] gebruikten verschillende n-grams, inclusief verschillende
216 manieren van *smoothing*. Ze testten als weging voor features zowel *boolean* als
217 *tf*, waarbij ze vonden dat *boolean* betere resultaten opleverden. Voor classifica-
218 tiemodel experimenteerden ze met SVM en Naive Bayes. Voor het selecteren
219 van *features* experimenteerden ze met een minimale frequentie en selectie van
220 woorden op basis van hoogste *mutual information*. Uiteindelijk was het beste
221 model bij hen een SVM met uni- en bigrams en geselecteerd op basis van *mutual*
222 *information*.

223 Ferreira maakten gebruik van twee classificatiemethoden: Logistische re-
224 gressie en *margin-infused relaxed algorithm* (MIRA) [4]. Logistische regressie
225 werd aangevuld met *group Lasso* regularisatie, wat het beste resultaat opleverde.
226 Voor wegingen van woorden werd geëxperimenteerd met *tf*, *tf-idf*, Δ -*tf-idf* en
227 Δ -*BM-25*. Daarnaast wordt er gebruik gemaakt van woordclustering, *Concise*
228 *Semantic Analysis* en stylometrische eigenschappen. Op *Part-Of-Speech tagging*
229 na hadden stylometrische eigenschappen een duidelijke negatieve invloed op de
230 classificatie.

231 Høyland et al. maakten gebruik van Support Vector Machine [7]. Als
232 beste waarde voor de regularisatieterm, de C-parameter, vonden zij 0.8. Daar-
233 naast gebruikten zij *dependency disambiguated stems*, wat een F_1 score van twee
234 procent hoger opleverde dan gebruik van normale stemming.

235 2.3 Invloed van partijnamen of sprekersnamen

236 Diermeier et al. [3] lieten de namen van de sprekers en verwijzingen naar staten
237 die de senatoren representeren weg, omdat deze volgens hen de classificatie te
238 makkelijk zouden maken. Hirst et al. [6] vonden inderdaad dat partijnamen -
239 en het weglaten daarvan - bij het Europees Parlement een grote invloed hebben
240 op de classificatie. Bij het Europees Parlement was te zien dat een spreker de

eigenaam gebruikte, terwijl in het Canadese parlement vooral te zien was dat de naam van de andere partij gebruikt wordt door een spreker.

2.4 Invloed van oppositie of regering

Hirst et al. [6] vonden in hun onderzoek dat de classificatie van spreker in het Canadese parlement op basis van partij-affiliatie meer zegt over de status van de partij (regering of oppositie). Zo vergeleken zij de top tien karakteristieke woorden van de liberalen en conservatieven in het 36e parlement (liberalen in de regering) en het 39e parlement (conservatieven in de regering. Hier vonden zij dat vier van de tien woorden van de liberalen (regering) in het 36e parlement bij het 39e parlement bij de conservatieven (regering) te vinden waren. Andersom gebeurde hetzelfde met één van de tien woorden van de conservatieven (oppositie) in het 36e parlement naar liberalen (oppositie) in het 39e parlement.

In hetzelfde onderzoek trainden ze ook hun classificaties op het ene parlement en testten deze op het andere parlement. Hierbij vonden zij in beide gevallen een nauwkeurigheid ver onder de baseline.

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

3 Methodologie

3.1 De data

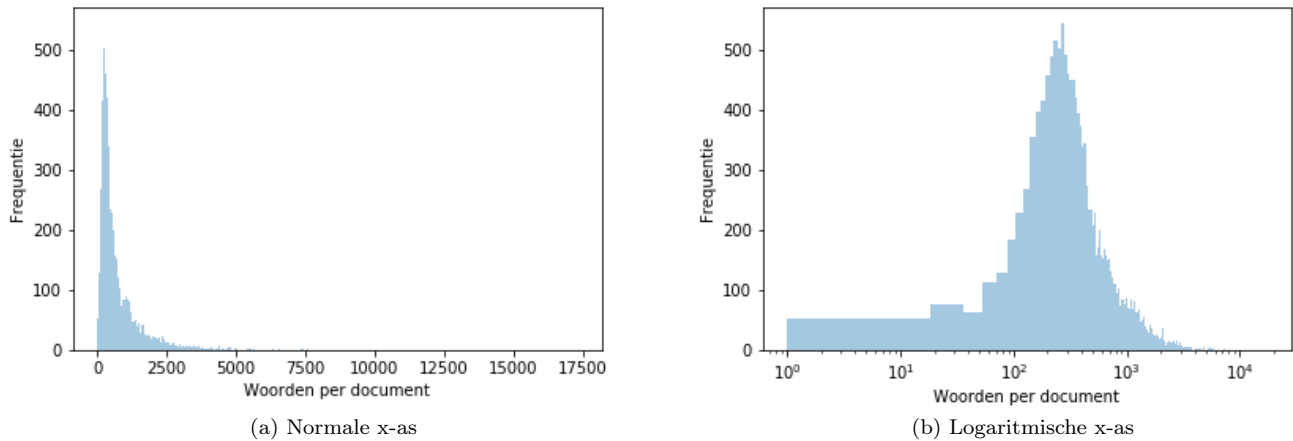
De data die gebruikt zijn, zijn de Handelingen van de Tweede Kamer gedurende het missionaire kabinet-Rutte II (5 november 2012 tot 22 maart 2017). Er was gekozen voor dit kabinet, omdat de data hiervoor makkelijk verkrijgbaar was, het kabinet lang zat - waardoor er veel data is - en het recent is waardoor het makkelijker te interpreteren is. In dit kabinet zaten de PvdA en VVD. Deze data zijn in xml-formaat van de website officiële bekendmakingen.nl gehaald samen met bijbehorende metadatabestanden. De bestanden van de Handelingen bevatten voornamelijk informatie over spreekbeurten tijdens een debat, waaronder naam van een spreker, partij-affiliatie, inhoud van de spreekbeurt en het soort spreekbeurt. Deze gegevens zijn samengevoegd tot één tabel.

Deze dataset bestaat uit een aantal soorten spreekbeurten; debatbijdragen, interrupties en antwoorden. Een debatbijdrage is de eerste onafgebroken spreekbeurt die een spreker geeft achter een spreekgestoelte, aangeduid in de xml-file met het attribuut *nieuw*="ja". Dit kan een bijdrage in een debat zijn of een vraag tijdens een vragenuur. Interrupties zijn de vragen die andere politici stellen vanachter de interruptiemicrofoon aan een spreker. De antwoorden zijn vervolgens de reactie van een spreker achter het spreekgestoelte op een interruptie. Aangezien een debatbijdrage geïnterrupteerd kan worden, kan deze inhoudelijk doorlopen in een antwoord van een spreker. Vergelijkbare onderzoeken voegen vaak alle teksten van een spreker samen tot één document. Dit was alleen niet mogelijk voor dit onderzoek met de hoeveelheid kleine partijen in de Tweede Kamer, die dan niet altijd in een training of test set zijn vertegenwoordigd. Daarom was in dit onderzoek ervoor gekozen om een debatbijdrage samengevoegd met alle bijbehorende antwoorden te beschouwen als één document.

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

Deze dataset bevatte vervolgens naast de verkozen partijen na de Tweede Kamerverkiezingen van 2012 ook afsplitsingen van die partijen (tien in totaal) en bezoeken van vertegenwoordigingen van Nederlandse partijen uit het Europees Parlement (tien in totaal). Omdat van beide categorieën relatief weinig data was en er overlap zat met hun oorspronkelijke of gelieerde partij, waren deze er uit gehaald. 50PLUS is in 2014 [9] uiteengevallen in twee fracties die aanspraak maakten op de partij-affiliatie 50PLUS. Vanaf dit moment zijn deze documenten niet meer meegenomen om onduidelijkheid te voorkomen.

De documenten verschilden in grootte (aantal woorden). De distributie van documentgrootte lijkt op een lognormale verdeling, maar met een Kolmogorov-Smirnov test was hier geen bewijs voor gevonden [8].



Figuur 1: Aantal woorden per document

Om toch de uitschieters er uit te halen, was aangenomen dat de distributie wel lognormaal verdeeld is en waren 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 nauwkeurigheid 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).

3.2 Methoden

3.2.1 DV1: Beste classificatiemethode

Om deze deelvraag te beantwoorden zijn een aantal classificatiemethoden vergeleken. Aangezien het niet mogelijk was om alle classificatiemethoden te vergelijken, beperkte dit onderzoek zich tot classificatiemethoden die gebruikt zijn in vergelijkbare onderzoeken, zoals besproken in sectie 2.2. Er was ervoor gekozen om alleen gebruik te maken van methoden waarvan reeds implementaties beschikbaar waren in scikit-learn. Voor alle methoden werd gezocht naar de beste parameters, ook wel bekend als een grid search. Deze grid search werd gedaan door vijfmaal kruisvalidatie (*cross-validation*), waarbij de training set steeds 80% was en de test set 20% van de totale dataset. Een totaal aantal van 6480 combinaties van methoden en parameters zijn getest. De verwachting was dat de scores lager zijn dan die gevonden in het gerelateerd werk, omdat de documentgrootte kleiner is en de baseline scores lager zijn.

Pre-processing Voor pre-processing is gebruik gemaakt van tokenisation en lowercasing. Voor tokenisation is de reguliere expressie $w+$ gebruikt, waardoor alles behalve letters en cijfers weggehaald wordt. Ook is er gevarieerd tussen wel of geen gebruik maken van stemming. In het geval van stemming is gebruik gemaakt van de Snowball Stemmer van de Python NLTK module.

Bag-of-words model Bag-of-words model is de meest gebruikte representatie van data in vergelijkbare onderzoeken. Deze is ook gebruikt in dit onderzoek.

333 Bij het bag-of-words model wordt elk document gerepresenteerd als een vector,
334 waarbij elke kolom een woord is met een bijbehorende waarde. Voornaamste
335 beperking van dit model is dat het geen rekening houdt met de volgorde van
336 woorden, wat een groot effect kan hebben op de betekenis van een document.

337 Voor dit onderzoek waren de volgende wegen voor woorden getest:
338 *boolean* (wel of niet aanwezig), *tf* (woordfrequentie), *tf-norm* (woordfrequentie
339 genormaliseerd door documentlengte) en *tf-idf* (woordfrequentie gecompenseerd
340 voor documentfrequentie). Daarnaast werd in dit onderzoek geëxperimenteerd
341 met een minimale of maximale woord- of documentfrequentie. Ook is gekeken
342 naar het effect van combinaties van de volgende n-grams; unigrams, bigrams en
343 trigrams. N-grams zijn combinaties van N aantal opeenvolgende woorden. Bij
344 een unigram is elke feature gewoon één woord, terwijl bij een bigram dit twee
345 opvolgende woorden zijn. Dit kan van belang zijn, want als bijvoorbeeld het
346 woord *asfalt* er in voorkomt, dan is het voor ideologie mogelijk relevant of er
347 *minder asfalt* of *meer asfalt* staat.

348 **Support Vector Machine en Logistische Regressie** De meest voorko-
349 mende techniek in vergelijkbaar onderzoek is Support Vector Machine (SVM).
350 Een andere techniek die gebruikt wordt, is logistische regressie. Beide hebben
351 een eigen implementatie in scikit-learn, maar deze implementaties zijn niet ef-
352 ficiënt met grote datasets. Om deze reden is er in beide gevallen voor gekozen
353 om gebruik te maken van de functie *SGDClassifier*, die beide technieken leert
354 met *stochastic gradient descent learning*. Voor regularisatie was hier geëxpe-
355 rimenteerd met L1 en L2 regularisatie en een combinatie van beide genaamd
356 Elasticnet. De andere parameters zijn gelaten op de standaardwaarden van
357 scikit-learn [12]. Een belangrijke onaangepaste waarde was die van maximaal
358 aantal iteraties, waarvoor de scikit-learn standaard 5 is. Volgens scikit-learn
359 convergeert de *SGDClassifier* rond de $10^6/n$ iteraties waar n het aantal docu-
360 menten in de training set is. In het geval van deze dataset zou dat 84 iteraties
361 zijn. Vanwege de grootte van de grid search was het voor dit onderzoek niet
362 mogelijk het maximaal aantal iteraties te verhogen tijdens de grid search. De
363 resultaten buiten de grid search zullen gebaseerd zijn op een maximaal aantal
364 iteraties van 100.

365 **Naive Bayes** Een andere techniek die gebruikt wordt voor politieke tekstclas-
366 sificatie is Naive Bayes. Dit algoritme neemt aan dat elke *feature* onafhankelijk
367 is ten op zichte van de rest. Dit is bij tekstclassificatie vaak niet het geval
368 omdat het gebruik van sommige woorden gepaard kan gaan met het gebruik
369 van andere woorden. Daarnaast is het gebruik van meerdere n-grams in een
370 classificatie schending van de aanname, want als bijvoorbeeld een bigram er in
371 voorkomt dan komen ook beide unigrams er in voor. Desalniettemin blijkt Naive
372 Bayes effectief te zijn voor tekstclassificatie [1, 12]. Hiervoor zijn de functies van
373 scikit-learn *MultinomialNB* en *BernoulliNB* gebruikt [1, 12].

374 **Beoordelen van kwaliteit** De meest gebruikte methoden om kwaliteit van
375 politieke tekstclassificatie te beoordelen zijn nauwkeurigheid en F_1 score, die
376 opgebouwd is uit sensitiviteit en precisie. Deze scores worden berekend op
377 basis van vier hoeveelheden van mogelijke resultaten van een classificatie. Deze

resultaten geven weer hoeveel documenten wel of niet bij een partij horen, en of deze wel of niet als dusdanig zijn geassocieerd [10] .

Tabel 2: Mogelijke resultaten van een classificatie.

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

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

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

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

$$F_1 = 2 * \frac{Precisie * Sensitiviteit}{Precisie + Sensitiviteit} \quad (4)$$

Nauwkeurigheid (*accuracy*) is het percentage van documenten dat correct geassocieerd is. Nauwkeurigheid wordt voor de hele classificatie gedaan en niet per klasse. Precisie (*precision*) is het percentage van documenten geassocieerd als een partij, dat ook bij die partij hoort. Sensitiviteit (*recall*) is het percentage documenten van documenten behorende tot een partij, dat ook als die partij geassocieerd is. F_1 is het harmonisch gemiddelde van sensitiviteit en precisie. Precisie, sensitiviteit 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 hoeveelheden van mogelijke resultaten bij elkaar opgeteld worden en vervolgens de scores berekend. Dit leidt ertoe dat resultaten van partijen met veel documenten belangrijker zijn. Als een classificatie kleine partijen grotendeels fout associeert, kan deze score alsnog hoog zijn. In het geval van meer dan twee partijen is dit hetzelfde als nauwkeurigheid.

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

Als laatste is er gewogen. Hierbij wordt net als macro de scores per partij berekend, maar neemt hiervan het gemiddelde gewogen bij het aantal documenten behorend tot een partij. Deze wijkt weinig af van de micro variant, tenzij er uitschieters zijn bij partijen.

Aangezien micro al terugkomt in nauwkeurigheid en het nadeel van macro te groot is omdat de partijen nogal variëren in grootte, was gekozen voor gewogen F_1 score naast nauwkeurigheid.

3.2.2 DV2: Invloed van namen

In Diermeier et al. [3] werd aangenomen dat namen een groot effect hebben op de classificatie. Hirst et al. [6] bevestigden dit voor het Europees Parlement. Aangezien hier bij deelvraag 1 niet voor was gekozen, werd bij deze

412 deelvraag gekeken hoe groot het effect hiervan is, specifiek gericht op partijna-
413 men en achternamen van Kamerleden. Op basis van vergelijkbaar onderzoek is
414 de hypothese dat de achternamen van Kamerleden en partijnamen van invloed
415 zijn.

416 Voor deze deelvraag werd wederom een classificatie gedaan met de classi-
417 ficatiemethode die resulteerde uit deelvraag 1. In deze classificatie werden alle
418 partijnamen vervangen door *PARTIJNAAM* en alle achternamen van Kamerle-
419 den vervangen door *KAMERLIDNAAM*. Deze namen waren uit de Handelingen
420 gehaald. Voor partijnamen waren ook lidwoorden toegevoegd en voor achter-
421 namen van Kamerleden zijn ook verkortingen meegenomen. Dit laatste omdat
422 bijvoorbeeld *Van Haersma Buma* vaak aangesproken wordt als *Buma*. Voorna-
423 men van Kamerleden worden zelden tot nooit gebruikt, dus die waren er niet
424 uitgehaald. Een nadeel van deze aanpak is dat ook namen van niet-Kamerleden
425 of andere woorden weggehaald kunnen worden als deze hetzelfde zijn als naam
426 van een Kamerlid. Door gebruik van gevoeligheid voor hoofdletters was gepro-
427 beerd dit te voorkomen. Een opvallend voorbeeld hiervan is de naam Rutte, die
428 zowel behoort tot het Kamerlid Arno Rutte als de premier Mark Rutte. Steek-
429 proefgewijs was gekeken of er nog namen achter zijn gebleven, maar die waren
430 niet gevonden.

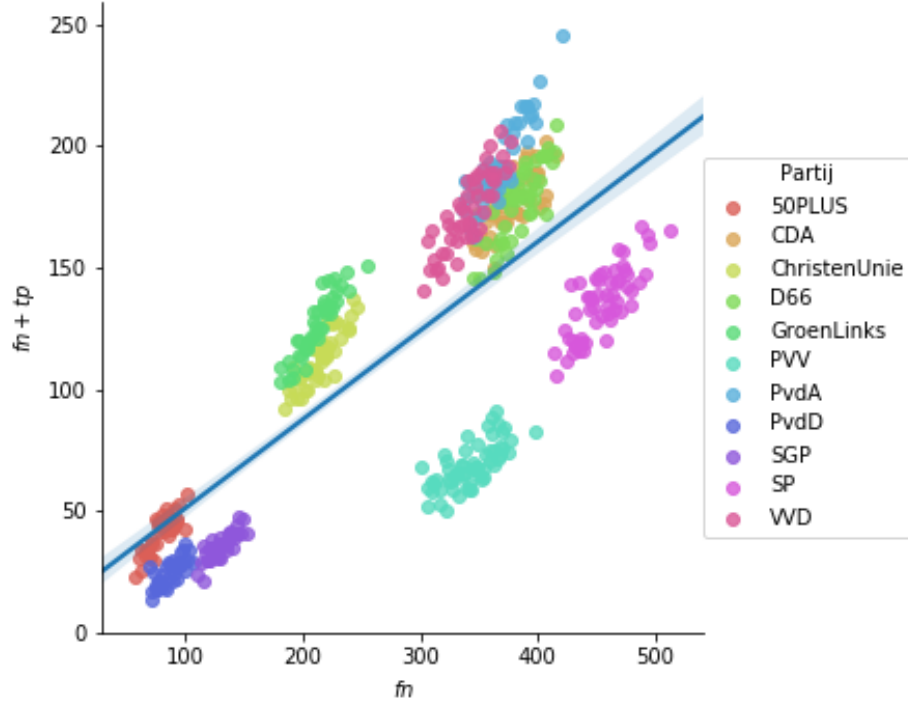
431 Ook werd gekeken naar classificatie met alleen partijnamen en achterna-
432 men van Kamerleden. Alle andere woorden worden weggehaald. Namen van
433 Kamerleden en partijen die niet aan elkaar geschreven worden, zoals *Partij van*
434 *de Arbeid*, zijn aan elkaar geschreven zodat het één feature is. Doordat alle
435 andere woorden weggehaald zijn, waren de bi- en trigrams combinaties van na-
436 men die zinnen uit elkaar kunnen staan, dus die niet meer informatie geven dan
437 unigrams. Daarom werd er gebruikt van de classificatiemethode uit deelvraag
438 1, maar met alleen unigrams.

439 Op basis van de hypothese is de verwachting dat voor de classificatie zon-
440 der namen de scores een stuk lager zijn dan deelvraag 1 en de scores van de
441 classificatie met alleen namen aanzienlijk hoger zijn dan de baseline scores.

442 3.2.3 DV3: Oppositie of regering

443 Om deze deelvraag te beantwoorden zijn drie experimenten uitgevoerd. Twee
444 daarvan zijn gebaseerd op experimenten uit Hirst et al. [6] voor dezelfde vraag.
445 De derde is ontwikkeld voor dit onderzoek. Met deze laatste wordt begonnen.
446 Bij deze deelvraag is de classificatiemethode uit deelvraag 1 zonder achternamen
447 van Kamerleden en partijnamen gebruikt. De hypothese is op basis van de
448 bevindingen van Hirst et al. dat de classificatie inderdaad afhankelijk is van
449 partij-status.

450 Als er een afhankelijkheid is van partij-status, dan is het te wachten dat
451 het aantal misclassificaties minus verwachte waarde binnen regeringspartijen en
452 binnen oppositiepartijen hoger ligt dan tussen een oppositiepartij en regerings-
453 partij. De verwachte waarde is afhankelijk van het aantal documenten van een
454 partij in de training set [13]. Aangezien de test set uit dezelfde set als de training
455 werd gehaald, is de verwachte waarde ook afhankelijk van het aantal documen-
456 ten van een partij in de test set. Uit de voorverkenning (op basis van resultaten
457 uit deelvraag 1 en 2) bleek deze correlatie tussen het aantal *false positives* van
458 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 100 classificaties met verschillende train en test set. De Pearson correlatie is 0.77 en de p-waarde 5.40×10^{-101} .

Op basis van dit verband is het verwachte aantal documenten ($V_{i,j}$) van partij i die foutief geassocieerd worden als partij j gedefinieerd als

$$V_{i,j} = fn_i * \frac{D_j}{D - D_i} \quad (5)$$

waar $i \neq j$ en D het totaal aantal documenten en D_i en D_j het aantal documenten van respectievelijk partij i en j . De teller van de breuk is het aantal documenten die bij partij j horen en de noemer het aantal documenten die niet bij partij i horen. Op deze manier is $\sum_{j=0}^n (V_{i,j}) = fn_i$ waar n het aantal partijen is minus partij i .

De error ($e_{i,j}$) is dan het verschil van het daadwerkelijk aantal misclassificaties ($D_{i,j}$) en de verwachte waarde ($V_{i,j}$)

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

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

Als dit een goede benadering is van de error, dan is het te verwachten dat deze normaal verdeeld is [11]. Om te kijken of er een bias is, werden de distributies binnen regeringspartijen en binnen oppositiepartijen vergeleken met

473 de distributie tussen beide groepen. Om de invloed van variantie door de wil-
474 lekeurige splitsing documenten voor trainen en testen te beperken, werd de
475 classificatie 100 keer gedaan. Met behulp van normaalheidstoets is gekeken of
476 de distributies normaal verdeeld zijn (α is 0.01). Als de distributies normaal
477 verdeeld zijn, vond de statistische test plaats op basis van een eenzijdige t-toets.
478 Als de distributies niet normaal verdeeld zijn, vond dit plaats door een Mann-
479 whitneytoets. Het gekozen significantieniveau (α) is 0.01. De nulhypothese is
480 dat er geen verschil is tussen de verdelingen. De alternatieve hypothese is dan
481 dat de distributie van binnen oppositie of regering groter is dan die tussen een
482 regerings- en oppositiepartij. Op basis van de bevindingen van Hirst et al. was
483 de hypothese dat de nulhypothese verworpen kan worden.

484 In het eerste experiment gebaseerd op Hirst et al. zijn de meest karakte-
485 ristieke woorden per partij van de ene zittingsperiode vergeleken met de meest
486 karakteristieke woorden per partij van de andere zittingsperiode, waar een ka-
487 binet uit andere partijen bestond. De verwachting is dat als de classificatie
488 niet het gevolg is van partij-status dat de woorden bij een partij blijven en niet
489 gekoppeld zijn aan in oppositie of regering zitten. Aansluitend bij de hypo-
490 these is dus de verwachting dat woorden wel wisselen van partij wanneer ze van
491 partij-status gewisseld zijn.

492 In het tweede experiment gebaseerd op Hirst et al. zijn de classificaties
493 getraind op een zittingsperiode en getest op een andere zittingsperiode, waar
494 wederom een kabinet uit andere partijen bestond. Als de classificatie afhankelijk
495 was van partij-status was de verwachting dat de scores van partijen die gewisseld
496 zijn van partij-status sterker gedaald zijn dan partijen die niet van partij-status
497 zijn veranderd. Op basis van de hypothese is dan ook de verwachting dat bij de
498 partijen die gewisseld zijn partij-status een sterkere daling te zien is.

499 Als vergelijkingsmateriaal was voor deze experimenten een tweede dataset
500 nodig uit een ander kabinet. Hiervoor is het wenselijk dat dit kabinet bestaat uit
501 andere partijen dan kabinet-Rutte II. Daarnaast is het ook wenselijk als het niet
502 te ver terug is, zodat onderwerpen en taalgebruik enigszins overeenkomstig zijn.
503 Omdat kabinet-Rutte I een minderheidskabinet was met een bijzondere partij-
504 status voor de PVV, is ervoor gekozen om de Handelingen van de Tweede Kamer
505 tijdens het missionaire kabinet-Balkenende IV (22 februari 2007 tot 20 februari
506 2010) te gebruiken. Dit kabinet bestond uit CDA, PvdA en ChristenUnie.

507 De partij 50PLUS bestond nog niet gedurende kabinet-Balkenende IV, dus
508 documenten van deze partij zijn weggelaten uit de dataset van kabinet-Rutte
509 II. Verder heeft dezelfde verwerking van data plaatsgevonden, zoals beschreven
510 in 3.1. Alleen de minimum- en maximumlengte is overgenomen van de dataset
511 van kabinet-Rutte II.

Tabel 3: 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

512 3.2.4 DV4: Links-rechts as

513 Als de classificatie afhankelijk is van positie op de links-rechts as dan is het te
514 verwachten dat, net als bij deelvraag 3, de misclassificaties minus de verwachte
515 waarde groter zijn als twee partijen dichterbij elkaar staan op de links-rechts as.
516 Daarvoor werd wederom formule 5 gebruikt als verwachte waarde en formule 6
517 als error. De hypothese is dat de classificatie deels afhankelijk is van positie op
518 de links-rechts as.

519 Er zijn verschillende methoden om partijen in te delen op een links-rechts
520 as. Er is hier gekozen voor de indeling van het Manifesto Project [14]. Het
521 Manifesto Project geeft scores op een heel aantal politieke posities, waaronder
522 de links-rechts as, op basis van het verkiezingsprogramma van dat jaar. Voor
523 de dataset van kabinet-Rutte II is gebruikt gemaakt van de scores op basis van
524 de verkiezingsprogramma's voor de verkiezingen van 2012.

Tabel 4: Scores op de links-rechts as 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

525 Er wordt vervolgens gekeken door middel van een Pearson correlatie toets

526 of er een correlatie is tussen de error van twee partijen en de afstand op de
527 links-rechts as van die partij. Het significantieniveau (α) hiervoor is opnieuw
528 0.01. De nulhypothese is dat er geen negatieve correlatie is tussen de error en
529 de afstand op de links-rechts as. De alternatieve hypothese is dat er wel een
530 negatieve correlatie is tussen de error en de afstand op de links-rechts as.

531 Als uit deelvraag 3 blijkt dat partij-status invloed heeft op de error, zal
532 bovenstaande methode ook uitgevoerd worden voor de aparte combinaties; bin-
533 nen oppositie en tussen regeringspartij en oppositiepartij. Binnen regering is
534 dit niet mogelijk aangezien er maar één afstand is, die tussen PvdA en VVD.

535 De voorspelling op basis van de hypothese is dat de nulhypothese verwor-
536 pen kan worden.

537 3.2.5 DV5: Woordgebruik van sprekers

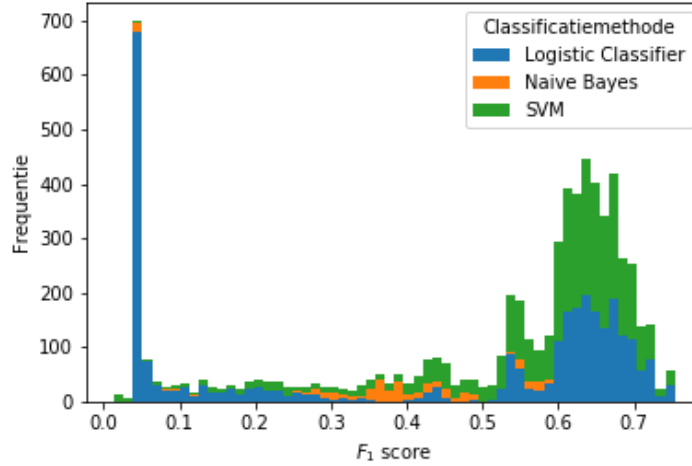
538 De vorige classificaties trainden op documenten en werden getest op andere
539 documenten, maar wel van dezelfde sprekers als uit de training set. Naast
540 de ideologie kan de classificatie daarom ook getraind zijn op het taalgebruik
541 van sprekers. Als een Kamerlid bijvoorbeeld een woord regelmatig in speeches
542 gebruikt, maar niet wordt gebruikt door zijn partijgenoten, wordt dit wel gezien
543 als een belangrijk woord voor de classificatie naar partij-affiliatie. Hirst et al.
544 [6] plaatsten al een soortgelijke kanttekening bij de resultaten van Diermeier et
545 al. [3]. De hypothese is dat de classificatie afhankelijk is van woordgebruik van
546 sprekers

547 Om te kijken of dit effect er is, wordt er opnieuw een classificatie gedaan
548 met de classificatiemethode uit deelvraag 1 zonder achternamen van Kamerleden
549 en partijnamen. Ditmaal worden alleen niet de individuele documenten verdeeld
550 over de training en test set, maar worden de Kamerleden, met bijbehorende
551 documenten, verdeeld over de training en test set. Als taalgebruik van een
552 spreker in de training set voorheen invloed had op de classificatie, zal dat nu geen
553 effect meer hebben omdat er geen documenten van die spreker meer voorkomen
554 in de test set. De verwachting is daarom ook dat deze classificatie lagere scores
555 vindt dan die van deelvraag 2.

556 4 Resultaten

557 4.1 DV1: Beste classificatiemethode

558 In figuur 3 zijn de uitslagen van de grid search te zien. Hierin is te zien dat SVM
559 en logistische regressie beide hoge scores behalen, maar dat logistische regressie
560 ook veel lage scores haalt tussen 0 en 0.1. Naive Bayes zit tussen de 0.25 en 0.6.



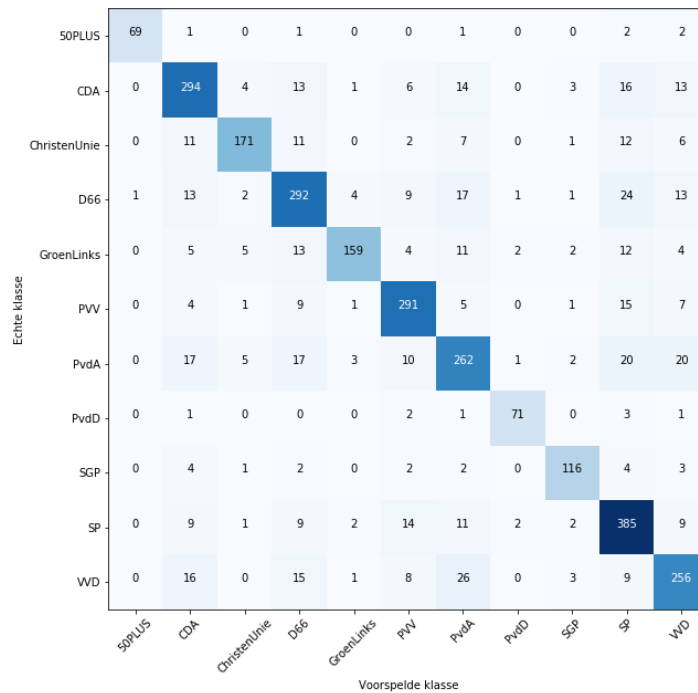
Figuur 3: Histogram van de grid search met de F_1 scores van de classificatiemethoden

561 Het beste resultaat werd bereikt met Support Vector Machines gebruik-
562 makend van *stochastic gradient descent learning* en L2 regularisatie. In de grid
563 search behaalde deze methode een F_1 score en nauwkeurigheid van 0.75. Voor
564 beide scores was dit het hoogste van de grid search. De woorden waren hier-
565 bij gestemd. De features waren zowel unigrams, bigrams als trigrams. Geen
566 features zijn weggelaten door minimale of maximale documentfrequenties. De
567 waarden van deze features waren *tf-idf* scores. Het maximum aantal iteraties
568 was 5 voor de grid search, maar de rest van resultaten zijn op basis van 100
569 iteraties.

570 Tabel 5 laat de scores zien per partij met het aantal documenten in de
571 test set. De nauwkeurigheid voor deze classificatie is 0.80. De F_1 scores per
572 partij liggen tussen de 0.7 en 0.9. De partijen met een sterke focus op één
573 onderwerp, 50PLUS, PVV en PvdD, als ook de SGP hebben hoge scores. De
574 coalitiepartijen, VVD en PvdA, daarentegen hebben lagere scores. Figuur 4
575 laat zien waar de fouten in deze classificatie zitten. De meest karakteristieke
576 n-grams per partij zijn te zien in tabel 6. Met meest karakteristiek worden de
577 n-grams bedoeld die de hoogste coëfficiënt hebben in de classificatie en die dus
578 relatief het meeste belangrijk zijn voor de classificatie van een partij. Hierin is
579 te zien dat vrijwel alle n-grams achternamen van Kamerleden of partijnamen
580 bevatten.

Tabel 5: Classificatie scores per partij van de beste classificatiemethode (SVM). Gemiddelde van vijfmaal kruisvalidatie.

	Precisie	Sensitiviteit	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 4: Confusion matrix van de beste classificatiemethode (SVM). Gemiddelde van vijfmaal kruisvalidatie.

Tabel 6: 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 6: 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

581 4.2 DV2: Invloed van namen

582 In tabel 6 was al te zien dat de meest karakteristieke n-grams voornamelijk ach-
583 ternamen van Kamerleden of partijnamen bevatten. In tabel 7 zijn de scores
584 te zien voor een classificatie met alleen achternamen van Kamerleden en partij-
585 namen. De nauwkeurigheid is 0.61. De scores zijn gedaald ten opzichte van de
586 resultaten van deelvraag 1, maar hoger dan de baseline scores.

Tabel 7: Classificatierapport van beste classificatie met alleen achternamen van Kamerleden en partijnamen. Hiervoor is alleen gebruikgemaakt van unigrams. Gemiddelde van vijfmaal kruisvalidatie.

	Precisie	Sensitiviteit	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

587 In tabel 8 zijn de F_1 scores te zien van classificatie met achternamen van
588 Kamerleden en partijnamen vervangen. De nauwkeurigheid hiervan is 0.58. De
589 scores zijn lager dan die uit deelvraag 1 en lager dan van de classificatie met
590 alleen namen. Wel zijn de scores nog steeds hoger dan de baseline. In tabel 9
591 is vervolgens te zien welke n-grams het meest karakteristiek zijn per partij voor
592 deze classificatie.

Tabel 8: Classificatie scores per partij van beste classificatiemethode (SVM) uit deelvraag 1 zonder achternamen van Kamerleden en partijnamen met het relatieve verschil in F_1 score ten opzichte van tabel 5. Gemiddelde van vijfmaal kruisvalidatie.

	Precisie	Sensitiviteit	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 9: Meest karakteristieke n-grams per partij op basis van de classificatiemethode (SVM) uit deelvraag 1 zonder achternamen van Kamerleden en partijnamen gedurende kabinet-Rutte II.

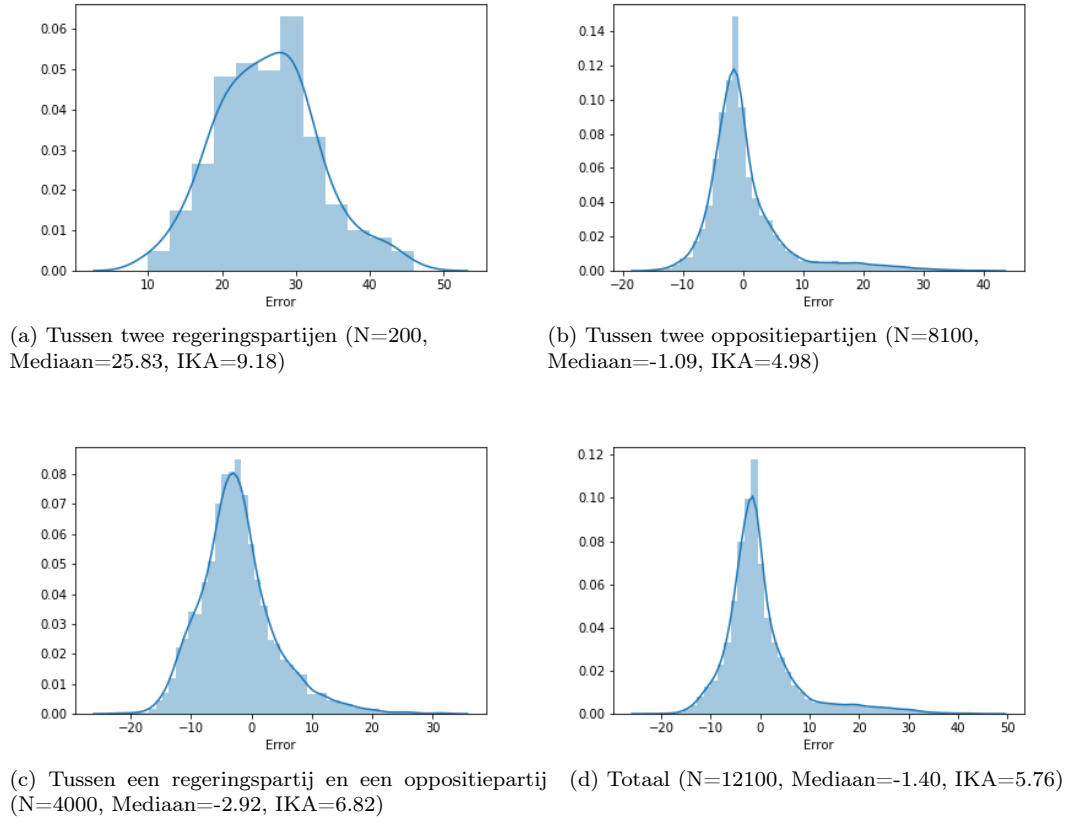
50PLUS	CDA	ChristenUnie	D66	GroenLinks
ouderen	PARTIJ fractie	dementie	mijn fractie	belastingontwijking
gepensioneerden	inwoners	gezinnen	mijn	zou
plussers	regering	zullen	natuurlijk	kamer hierover te
50 plussers	PARTIJ	vluchtelingen	fractie	persoonsgebonden
oudere	de regering	ik hoop	het kabinet	in elk geval
koopkrachtontwikkeling	diverse	inderdaad	buitengewoon	elk geval
50	hier	motie	belangrijk	vluchtelingen
werkenden	echt	hoop	vandaag	in elk
PARTIJ	een aantal	begeleiding	kabinet	hierover te informeren
overwegende dat	fractie	horeca	daarom	budget

Tabel 9: Meest relevante n-grams per partij op basis van de classificatiemethode (SVM) 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	volgens mij
miljard	daarbij	constaterende	beantwoording	van mening dat	PARTIJ fractie
natuurlijk	tevens	constaterende dat	voor de beantwoording	bezuinigingen	aruba
islam	vragen	dierenwelzijn	bewindslieden	mensen	regelgeving
de islam	wij	bio industrie	de beantwoording	huurders	aangegeven
al	beter	industrie	wel	voorstellen	speelveld
dit kabinet	kinderen	de bio	punt	segregatie	volgens
brussel	samen	dierproeven	nadrukkelijk	van mening	essentieel
asielzoekers	toezeggingen	de bio industrie	je	bestuurders	en

593 4.3 DV3: Oppositie of regering

594 In figuur 5 zijn de distributies van de errors, zoals gedefinieerd in formule 6 te
595 zien van combinaties van regerings- en oppositiepartijen. Bijgevoegd zijn het
596 aantal combinaties (N), het gemiddelde (μ) en de standaarddeviatie (σ).



Figuur 5: Genormaliseerde distributie van de error uit formule 6 voor de verschillende combinaties.

Voor alle distributies was de nulhypothese verworpen worden dat deze normaal verdeeld zijn ($p < 0.01$). In tabel 10 is vervolgens te zien dat er een significant verschil is tussen de distributies binnen regering en binnen oppositie tegenover de distributie tussen regering en oppositiepartij. Tussen regeringspartijen zijn er gemiddeld 26.11 misclassificaties meer dan verwacht en tussen oppositiepartijen gemiddeld 0.43.

Tabel 10: Uitslagen van eenzijdige Mann-whitneytoets tussen de distributie tussen een regeringspartij en oppositiepartij en twee distributies. α is 0.01.

	p -waarde	U -waarde
Tussen twee regeringspartijen	7.04×10^{-124}	717042
Tussen twee oppositiepartijen	4.4×10^{-108}	16328471

In tabel 11 zijn de meest karakteristieke n-grams te zien voor classificatie van kabinet-Balkenende IV. Hierin zijn geen opvallende overlappen te zien van regeringspartijen met de classificatie van kabinet-Rutte II in tabel 9.

Tabel 11: Meest karakteristieke n-grams per partij op basis van beste classificatiemethode uit deelvraag 1 zonder achternamen van Kamerleden en partijnamen gedurende kabinet-Balkenende IV.

CDA	ChristenUnie	D66	GroenLinks	PVV
PARTIJ fractie	fractie van PARTIJ	premier	PARTIJfractie	burgers
wij	de fractie van	de premier	fractie van PARTIJ	onze
fractie	de fractie	hij	de fractie	burger
wij hebben	fractie van	ik hoop	de fractie van	gewoon
aangegeven	moment	arbeidsmarkt	fractie van	natuurlijk
PARTIJ fractie heeft	mijn fractie	plannen	premier	de burgers
dank	verschillende	hoop	mij	door
overleg	beantwoording	de arbeidsmarkt	ik	politie
KAMERLID	PARTIJfractie	dadelijk	politieke	land
buitengewoon	blij	ministerie	en	niet

Tabel 11: Meest karakteristieke n-grams per partij op basis van beste classificatiemethode uit deelvraag 1 zonder achternamen van Kamerleden en partijnamen gedurende kabinet-Balkenende IV. (Vervolg)

PvdA	PvdD	SGP	SP	VVD
vrouwen	dieren	mijn fractie	mensen	PARTIJ
wij	bio industrie	wel	zegt	PARTIJ fractie
belangrijk	dierenwelzijn	beantwoording	leerlingen	onze fractie
kinderen	bio	voorzitter ik wil	is	fractie
goed	de bio industrie	toch	niet	ondernemers
vragen	de bio	diverse	vandaar	je
antwoorden	natuur	de bewindslieden	verdrag	praten
medewerkers	dierproeven	allerlei	personeel	markt
ben	veehouderij	natuurlijk	problemen	dat
iedereen	industrie	bewindslieden	waarom	voorzitter PARTIJ fractie

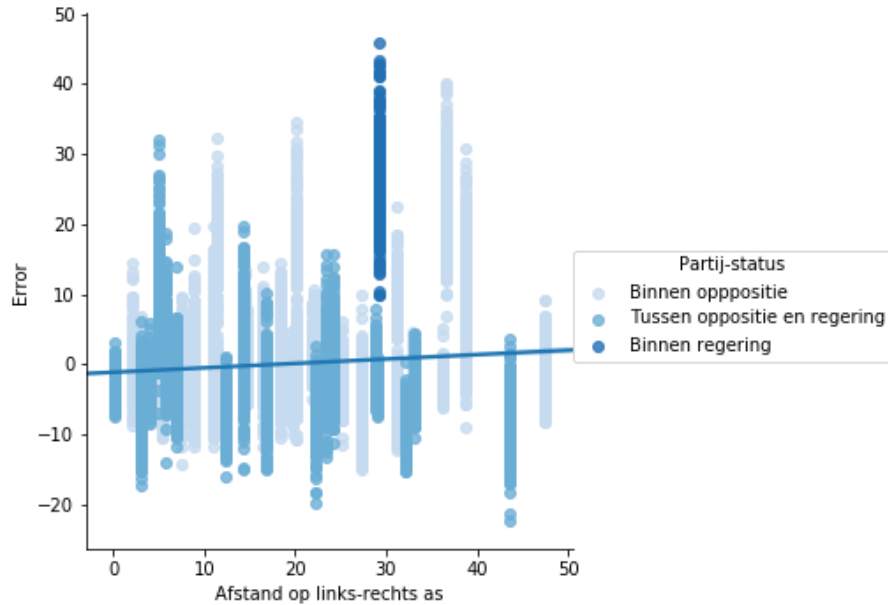
In tabel 12 zijn de resultaten van de classificatiescores te zien waarbij de classificatie getraind is op een zittingsperiode, maar getest op een andere. De resultaten zijn sterk gedaald, maar nog boven de baseline. De daling verschilt enorm per partij en zittingsperiode met dalingen van F_1 scores tussen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

Tabel 12: F_1 scores van de classificatie getraind op de dataset van Balkenende IV of Rutte II (minus 50PLUS) en getest op de ander. 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

611 4.4 DV4: Links-rechts as

612 In tabel 6 is de error te zien ten opzichte van de afstand op de links-rechts as.



Figuur 6: Error ten opzichte van de afstand op de links-rechts as van twee partijen. Gebaseerd op 100 classificaties met verschillende test en train set. De Pearson correlatie is 0.09 en de p -waarde 2.39×10^{-20} .

De Pearson correlatie van 0.09 is daarmee met een p -waarde van 2.39×10^{-20} significant op het significantieniveau van 0.01, maar wel positief gecorreleerd. Uit deelvraag 3 bleek dat de error binnen oppositie of regering significant afweek van de error tussen regering en oppositie. Dit effect lijkt ook zichtbaar in figuur 6. Daarom is er ook gekeken naar de correlatie tussen afstand op de links-rechts as en error binnen oppositie en tussen regerings- en oppositiepartij. De resultaten zijn te zien in tabel 13. Beide correlaties zijn statistische significant op het significantieniveau van 0.01, maar in tegengestelde richting.

Tabel 13: Pearson correlatie tussen error en afstand op de links-rechts as voor combinaties van partij-status.

	Pearson correlatie	p -waarde
Tussen oppositie- en regeringspartij	-0.29	3.44×10^{-69}
Tussen twee oppositiepartijen	0.18	1.76×10^{-55}

4.5 DV5: Woordgebruik van sprekers

In tabel 14 staan de scores van classificatie waarbij de Kamerleden verdeeld zijn over de training en test set. De scores zijn hierbij nauwelijks hoger dan de baseline.

Tabel 14: Classificatierapport van beste classificatie met de Kamerleden verdeeld over training en test set. Gemiddelde van tienmaal kruisvalidatie.

	Precisie	Sensitiviteit	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	

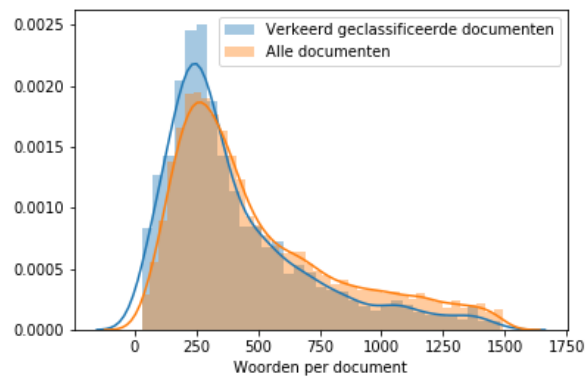
5 Discussie

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 6 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 0.05.

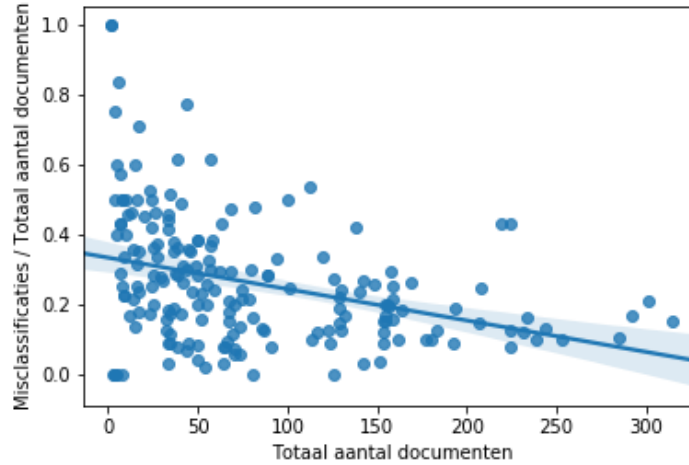
Het onderzoek van Hirst et al. vond dat resultaten afhankelijk kunnen zijn van documentgrootte. Alle documenten in dit onderzoek zijn kleiner dan de grootste documentgrootte uit het onderzoek van Hirst et al. en ook de minimale documentgrootte ligt lager dan de kleinste documentgrootte uit dat onderzoek.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



Figuur 7: Genormaliseerde distributie van documentlengtes van misclassificaties en alle documenten. Totaal van vijfmaal kruisvalidatie, waardoor documenten vaker voor kunnen komen. Mediaan documentlengte van foutief geclassificeerde documenten is 321 en voor alle documenten 386.

649 Voor een vervolgonderzoek kan uitgebreider gekeken worden naar dit effect
 650 en wat dit betekent voor de resultaten. Het percentage documenten van een vra-
 651 genuur was tweemaal zo hoog bij misclassificaties. De mediaan documentlengte
 652 van deze documenten is 286. De hoge aanwezigheid van deze documenten bij
 653 misclassificaties lijkt daarmee het gevolg van kleinere documentlengte.

654 Er is verder nog gekeken naar andere verbanden tussen sprekers wiens
 655 documenten vaak verkeerd geclassificeerd zijn. Daarbij is gevonden dat sprekers
 656 met weinig documenten relatief iets meer voorkomen in misclassificaties.



Figuur 8: Aantal misclassificaties gedeeld door totaal aantal documenten per spreker tegenover totaal aantal documenten van een spreker. Misclassificaties zijn totaal van 5-fold cross-validation. Hierdoor kunnen documenten vaker voorkomen in misclassificaties en ook meerdere keren mee tellen voor het totaal. De Pearson correlatie is -0.28 en de p -waarde 1.07×10^{-4} .

657 Dit versterkt het vermoeden dat de classificatie mede plaatsvond op basis
 658 van woordgebruik van individuele sprekers, zoals besproken wordt in deelvraag
 659 5.

660 5.2 DV2: Invloed van namen

661 De resultaten laten zien dat de classificatie sterk afhankelijk is van partijnamen
 662 en achternamen van Kamerleden. De hogere scores voor de classificatie met
 663 alleen namen dan de scores van de classificaties zonder namen in combinatie met
 664 de woorden in tabel 6, suggereert dat dit het belangrijkste was in de classificatie
 665 van deelvraag 1. Deze daling was te verwachten op basis van gerelateerd werk.

666 De n-grams in tabel 9 komen bij veel partijen overeen met hun ideologie,
 667 vooral bij de partijen met een sterke focus op één onderwerp; PVV, PvdD
 668 en 50PLUS. Daarnaast zijn er ook n-grams die niet veel over ideologie lijken
 669 te zeggen, zoals; *volgens mij*, *ik constateer* en *in elk geval*. Vooral de SGP
 670 heeft n-grams die niet veel lijken te zeggen over de ideologie, hoewel deze partij
 671 desalniettemin de hoogste F_1 score heeft. Met name opvallend hierbij is *mevrouw*
 672 *de voorzitter*, aangezien deze woorden door alle partijen gebruikt worden om via
 673 de voorzitter te praten. Voor een vervolgonderzoek kan gekeken naar waarom
 674 deze n-grams zo karakteristiek zijn voor partijen.

675 De classificatiemethode die gebruikt is in deze deelvraag, is gebaseerd op
 676 de beste methode voor de dataset uit deelvraag 1. Hiervoor was gevonden dat
 677 een combinatie van uni-, bi- en trigrams het beste resultaat opleverde. In tabel
 678 6 is te zien dat trigrams behoren tot de meest karakteristieke n-grams, hoewel
 679 de woorden in trigrams vaak overlappen met uni- en bigrams. In tabel 9 daaren-
 680 tegen zijn er nog maar een paar trigrams, welke grotendeels procedurele zinnen

zijn of een lidwoord toegevoegd aan een uni- of bigram. Dit verschil suggereert dat trigrams minder belangrijk zijn in de classificatie zonder de namen, dus de classificatiemethode uit deelvraag 1 niet het beste is voor deze classificatie of latere classificaties. In vervolgonderzoek kan de opzet van deelvraag 1 toegepast worden op de classificatie van deze deelvraag of latere deelvragen, om zo te komen tot een classificatiemethode die het beste resultaat oplevert voor de deelvraag.

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

5.3 DV3: Oppositie of regering

In tabel 5 is te zien dat de coalitiepartijen lage scores krijgen. Daarnaast laat figuur 4 zien dat er een relatief grote overlap zit tussen deze twee partijen. De resultaten van het eerste experiment, ontwikkeld voor dit onderzoek, vinden een afhankelijkheid van partij-status. De twee andere experimenten versterken deze bevindingen niet, zoals wel het geval was bij Hirst et al. [6]. Hieronder worden de resultaten nader besproken.

De statistische toetsresultaten in tabel 10 laten zien dat inderdaad de error groter is binnen oppositie of regering dan tussen een regerings- en oppositiepartij. Met name regeringspartijen lijken lastiger uit elkaar te halen. Dit suggereert dat inderdaad partij-status invloed heeft op de classificatie.

De verwachting was dat de error normaal verdeeld zou zijn. De verdelingen uit figuur 5 hebben globaal wel de vorm van een normaal verdeling. In figuur 2 is het daarnaast opvallend dat partijen zoals SP en PVV ruim onder de regressielijn zitten, terwijl andere partijen er een stuk boven zitten. Dit geeft het vermoeden dat er naast het aantal documenten van een partij nog meer factoren van invloed zijn op het aantal misclassificaties en daarmee de verwachte waarde. Deze verwachte waarde en de daar uit volgende error waren een belangrijke aanname van deze methode. Voor deze methode is het dus belangrijk uit te vinden of dit een goede benadering is van de verwachte waarde. In deelvraag 4 wordt gekeken of links-recht as positie hier nog invloed heeft. Voor een vervolgonderzoek kan nog verder gekeken worden naar invloeden op verwachte waarde of andere *confounding biases*.

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

Tabel 15: 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</i>	<i>algemeen</i>
		<i>hun</i>	<i>algemeen overleg</i>
		<i>collega KAMERLID</i>	<i>toezegging</i>
		<i>in</i>	<i>helder</i>
		<i>aanpak</i>	<i>overleg</i>
		<i>collega</i>	<i>aangegeven</i>
	ChristenUnie		<i>voor</i>
			<i>voor PARTIJ</i>
		<i>mijn</i>	<i>gaan</i>
		<i>waarop</i>	<i>termijn</i>
	PvdA	<i>blij</i>	<i>blij met de</i>
		<i>collega KAMERLID</i>	<i>volgens</i>
		<i>erg</i>	<i>volgens mij</i>
			<i>blij</i>
			<i>beantwoording</i>
			<i>volgens</i>
			<i>volgens mij</i>

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

729 De scores in tabel 12 laten een duidelijke daling zien ten opzichte van een
730 classificatie van alleen kabinet-Rutte II. Deze algemene daling zou verklaard
731 kunnen worden door veranderingen in ideologie, woordgebruik, onderwerpen
732 en/of aantal documenten per partij. De daling is het grootst bij VVD, maar valt
733 mee bij de twee andere partijen die gewisseld zijn van partij-status, ChristenUnie
734 en CDA. Daarnaast is de daling ook heel sterk bij oppositiepartijen GroenLinks
735 en D66, alsook de regeringspartij in beide kabinetten, PvdA. Dat de daling niet
736 consequent groter is bij partijen die gewisseld zijn van partij-status, suggereert
737 dat de invloed van partij-status beperkt is op de classificatie.

738 Dat de experimenten uit Hirst et al. in hun onderzoek wel invloed vonden,
739 maar in dit onderzoek niet, kan komen doordat hun onderzoek zich richtte op
740 binaire classificatie, terwijl dit onderzoek meerdere partijen had. Zo kan het ont-
741 breken van gemeenschappelijke n-grams komen doordat regeringspartijen zich
742 ook van elkaar moesten onderscheiden in dit onderzoek. Daarvoor hebben n-
743 grams die relevant zijn voor partij-status weinig effect. In het onderzoek van
744 Hirst et al. daarentegen de regeringspartij alleen onderscheiden hoefde te wor-
745 den van de oppositiepartij. Daarnaast verklaarden zij dat de daling tussen twee
746 zittingsperioden het gevolg was van die wisseling van partij-status. In dit on-

derzoek kon daarentegen gekeken worden naar effecten op partijen niet die niet van partij-status zijn gewisseld. Hierin was te zien dat de daling ook aanwezig was bij partijen die niet gewisseld zijn van partij-status.

5.4 DV4: Links-rechts as

De correlatie was tegen de verwachting in positief, waardoor de nulhypothese niet verworpen kan worden. Een deel van deze positieve correlatie lijkt het gevolg van de error tussen de twee regeringspartijen. Daarnaast is het opvallend dat tussen oppositiepartijen de correlatie ook positief is, maar tussen oppositie en regeringspartij juist, zoals eigenlijk verwacht, negatief. Een verklaring hiervoor is niet gevonden.

Alle correlaties zijn statistisch significant, maar de Pearson correlatie en daarmee effectgrootte is klein. Daarnaast is het ook opvallend dat de twee combinaties van partij-statussen een andere correlatierichting hebben. Dit suggereert dat de statistische significantie het gevolg is van de grote steekproef en maar een klein effect [5].

Er zijn verschillende visies op links en rechts en de indeling van partijen op die as. Daarnaast zijn er nog meerdere assen waarlangs partijen vergeleken kunnen worden. Bijvoorbeeld op basis van conservatief en progressief. Een vervolgonderzoek kan uitgebreider kijken naar welke assen relevant zijn voor partijen in de Tweede Kamer en in hoeverre deze invloed hebben op de classificatie.

5.5 DV5: Woordgebruik van sprekers

De resultaten uit tabel 14 zijn laag, amper hoger dan de baseline. Dit suggereert inderdaad dat eerdere classificaties in grote mate toch afhankelijk waren van het woordgebruik van sprekers. Dit is opmerkelijk aangezien vergelijkbare onderzoeken dit effect niet vinden. De meest karakteristieke n-grams van deze classificatie wijken daarnaast grotendeels niet af van die uit tabel 9.

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

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

792 5.6 Algemeen

793 Het vergelijken van deze resultaten met vergelijkbaar werk is ingewikkeld, aan-
794 gezien de keuzes en eigenschappen van die onderzoeken het niet een één-op-één
795 vergelijking maken. Voorbeelden hiervan zijn de taal, het parlement, de do-
796 cumentgrootte, baselines, behouden of weglaten van namen, een spreker als
797 document zien en het trainen en testen op dezelfde spreker. Hoewel de resul-
798 taten in sommige gevallen lager waren dan die uit vergelijkbaar werk, is het
799 belangrijk hier rekening mee te houden. Een vervolgonderzoek zou daarom dit
800 onderzoek kunnen reproduceren op een ander parlement om daarmee te kunnen
801 vergelijken.

802 Dit onderzoek richtte zich hoofdzakelijk op de Handelingen gedurende
803 kabinet-Rutte II. Om te kijken in hoeverre het mogelijk is om deze conclusie
804 door te trekken naar de algemene Handelingen van de Tweede Kamer, kan er
805 in vervolgonderzoek gekeken worden naar meerdere zittingsperioden. Ook kan
806 gekeken worden naar veranderingen als een kabinet demissionair is.

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

815 6 Conclusies

816 Dit onderzoek vindt een nauwkeurigheid en F_1 score van 0.80 voor het classi-
817 ficeren van spreekbeurten in de Tweede Kamer naar partij-affiliatie. De beste
818 classificatiemethode maakt gebruik van Support-Vector Machines. De baseline
819 scores zijn respectievelijk 0.11 en 0.15. Als rekening wordt gehouden met par-
820 tijnamen en achternamen Kamerleden daalt de nauwkeurigheid naar 0.58 en
821 de F_1 score naar 0.57. Dit onderzoek vindt aanwijzingen dat deze classificatie
822 afhankelijk is van de partij-status (oppositie of regering). Daarnaast vindt dit
823 onderzoek geen aanwijzingen dat de classificatie afhankelijk is van positie op de
824 links-rechts as. Als rekening wordt gehouden met woordgebruik van individuele
825 Kamerleden, dalen de nauwkeurigheid en F_1 verder naar 0.27. Daarmee lijkt
826 de classificatie naar partij-affiliatie in grote mate niet het gevolg van ideologie.
827 Deze conclusie trekt daarmee de bruikbaarheid van tekstclassificatie voor het
828 vinden van een relatie tussen woordgebruik en ideologie in twijfel. Op een aantal
829 punten wijken de bevindingen van dit onderzoek af van vergelijkbare onderzoe-
830 ken. Voor een vervolgonderzoek kan dit onderzoek uitgebreid worden met een
831 aantal aanbevelingen.

832 Referenties

- 833 [1] Bhand, M., Robinson, D., and Sathi, C. (2009). Text classifiers for political
834 ideologies.

- [2] Bießmann, F. (2016). Automating political bias prediction. *CoRR*, abs/1608.02195.
- [3] Diermeier, D., Godbout, J.-F., Yu, B., and Kaufmann, S. (2012). Language and ideology in congress. *British Journal of Political Science*, 42(1):31–55.
- [4] Ferreira, V. (2016). Using textual transcripts of parliamentary interventions for profiling portuguese politicians.
- [5] Hair, Jr., J. F., Anderson, R. E., Tatham, R. L., and Black, W. C. (2006). *Multivariate Data Analysis (6th Ed.)*. Prentice-Hall, Inc., Upper Saddle River, NJ, USA.
- [6] Hirst, G., Riabinin, Y., Graham, J., and Boizot-roche, M. (2014). Text to ideology or text to party status? In Kaal, B., Maks, I., and van Elfrinkhof, A., editors, *From Text to Political Positions*, chapter 5, pages 93–115. John Benjamins Publishing Company, Amsterdam.
- [7] Høyland, B., Godbout, J.-F., Lapponi, E., and Velldal, E. (2014). Predicting party affiliations from european parliament debates. In *Proceedings of the ACL 2014 Workshop on Language Technologies and Computational Social Science*, pages 56–60. Association for Computational Linguistics.
- [8] Jones, E., Oliphant, T., Peterson, P., et al. (2001). SciPy: Open source scientific tools for Python.
- [9] Klompenhouwer, L. (2014). Extra ledenvergadering 50plus om splitsing. *NRC Handelsblad*.
- [10] Manning, C. D., Raghavan, P., and Schütze, H. (2008). *Introduction to Information Retrieval*. Cambridge University Press, New York, NY, USA.
- [11] NIST/SEMATECH (2012). *e-Handbook of Statistical Methods*. NIST/SEMATECH.
- [12] Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830.
- [13] Sahare, M. and Gupta, H. (2012). A review of multi-class classification for imbalanced data. *International Journal of Advanced Computer Research*, 2(3).
- [14] Volkens, A., Lehmann, P., Matthieß, T., Merz, N., Regel, S., and Weßels, B. (2017). The manifesto data collection. manifesto project (mrg/cmp/marpor). version 2017b.
- [15] Yu, B., Kaufmann, S., and Diermeier, D. (2008). Classifying party affiliation from political speech. *Journal of Information Technology & Politics*, 5(1):33–48.