

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 as	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 as	22
34	4.5 DV5: Woordgebruik van sprekers	23
35	5 Discussie	24
36	5.1 DV1: Beste classificatiemethode	24
37	5.2 DV2: Invloed van namen	26
38	5.3 DV3: Oppositie of regering	26
39	5.4 DV4: Links-rechts as	28
40	5.5 DV5: Woordgebruik van sprekers	28
41	5.6 Algemeen	29
42	6 Conclusies	29

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 meerdere classificatiemethoden en daarnaast specifiek op de Nederlandse politiek. De onderzoeksvraag luidt daarom dus ook: "In hoeverre is classificatie naar 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 positie op de links-rechts as?
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. Bij de vierde vraag zal

gekeken worden naar een verband tussen misclassificaties en afstand tussen twee partijen op de links-rechts as. Als laatste zal voor de vijfde vraag de classificatie herhaald worden met Kamerleden verdeeld over training en test set.

Overzicht van scriptie Sectie 2 bevat vergelijkbare onderzoeken in andere parlementen. Sectie 3 bevat de methodologie van de verschillende deelvragen. Sectie 4 bevat de resultaten. Sectie 5 bevat de evaluatie van zowel de resultaten als de methodologie. Sectie 6 bevat ten slotte het antwoord op de onderzoeksvraag.

2 Gerelateerd werk

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

In dit hoofdstuk zullen verschillende onderzoeken behandeld worden die tekstclassificatie hebben toegepast op parlementaire teksten. Eerst zullen de onderzoeken algemeen besproken worden. Vervolgens zal uitgebreider gekeken worden naar de effecten van verschillende classificatiemethoden. In de latere secties zullen aspecten besproken worden die in vergelijkbare onderzoeken genoemd worden als van invloed op de classificatie.

2.1 Tekstclassificatie van parlementaire teksten

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

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

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

Vervolgens herhaalden ze de classificaties op het Huis van Afgevaardigden uit 2015, maar testten ditmaal op de Senaat elk jaar tussen 1989 en 2006 afzonderlijk. Hier zien zij een stijging in *accuracy* van 60.0% (baseline van 55.0%) in 1989 naar 87.0% (baseline van 55.0%) in 2006, maar met twee duidelijke dalen. Ze presenteren twee mogelijke verklaringen voor de trend; het veranderen van de onderwerpen en het sterker verdeeld worden van het Congres.

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

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

Bhand et al. [6] richtten zich op het classificeren van leden van het Amerikaanse Congres in 2005, op basis van affiliatie (Republikeins of Democratisch). Een document hierbij was in tegenstelling tot eerdergenoemde onderzoeken een speech. Zij vonden hiervoor uiteindelijk een F_1 score van 0.68 (baseline niet vermeld).

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

Høyland et al. trainden een classificatie voor partij-affiliatie op basis van teksten van het vijfde Europese Parlement (1999-2004) en testten vervolgens op het zesde Europese Parlement (2004-2009) [7]. Alle teksten van een spreker zijn samengevoegd tot één document. 40% van de sprekers in de test set zaten ook in de training set. Hier verkregen zij een *macro* F_1 score van 0.464 (baseline van 0.097) en *accuracy* van 0.551 (baseline van 0.410). Hun baseline is op basis van altijd classificeren als grootste partij, terwijl voor F_1 score de baseline hoger ligt als hiervoor gekozen wordt voor gokken gewogen bij grootte van een klasse.

2.2 Classificatiemethoden

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

Yu et al. [5] maakten gebruik van Support Vector Machines en Naive Bayes, waarvan de varianten multinomial en Bernoulli. De features waren unigrams, met minimale woordfrequentie van drie en de top 50 meest voorkomende woorden weggelaten. Voor de wegingen van de features bij Support Vector Machines werd geëxperimenteerd met *boolean*, *tf-norm* en *tf-idf*. Het beste resultaat was afhankelijk van welke kamer. Voor het Huis van Afgevaardigden was het Support Vector Machines met als weging *tf-idf* en voor de Senaat Bernoulli Naive Bayes.

Hirst et al. [2] maakten gebruik van Support Vector Machines. Ze experimenteerden met verschillende vormen van pre-processing, inclusief stemmen en het verwijderen van woorden op basis van te hoge of te lage frequentie. Deze

181 variaties maakten in hun onderzoek geen grote verschillen en uiteindelijk is ge-
 182 kozen voor het niet stemmen, het weglaten van woorden die in minder dan
 183 vijf documenten voorkomen en resultaten van zowel met als zonder de top 500
 184 meest frequente woorden. Daarnaast werd geëxperimenteerd met vier wegingen
 185 voor woorden: *boolean*, *tf*, *tf-norm* en *tf-idf*, waarvan *tf-idf* het beste resultaat
 186 opleverde.

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

195 Ferreira maakten gebruik van twee classificatiemethoden: Logistische re-
 196 gressie en MIRA [3]. Logistische regressie werd aangevuld met *group Lasso*
 197 regularisatie. Voor wegingen van woorden werd geëxperimenteerd met *tf*, *tf-idf*,
 198 Δ -*tf-idf* en Δ -*BM-25*. Daarnaast wordt er gebruik gemaakt van woordcluste-
 199 ring, *Concise Semantic Analysis* en stylometrische eigenschappen. Op *Part-Of-*
 200 *Speech tagging* na hadden stylometrische eigenschappen een duidelijke negatieve
 201 invloed op de classificatie.

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

206 2.3 Invloed van partijnamen of sprekersnamen

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

214 2.4 Invloed van oppositie of regering

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

224 In hetzelfde onderzoek trainden ze ook hun classificaties op het ene par-
 225 lement en testten deze op het andere parlement. Hierbij vonden zij in beide
 226 gevallen een *accuracy* ver onder de baseline. Daarnaast deden ze ook nog een

classificatie op de sprekers die in beide parlementen zaten en een andere classificatie op sprekers die niet in beide parlementen zaten. Bij de eerste classificatie vonden ze *accuracy* scores rond de baseline, terwijl in de tweede situatie *accuracy* scores gevonden werden ver boven 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 worden, zijn de Handelingen van de Tweede Kamer gedurende het missionaire kabinet-Rutte II (5 november 2012 tot 22 maart 2017). Er is 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.

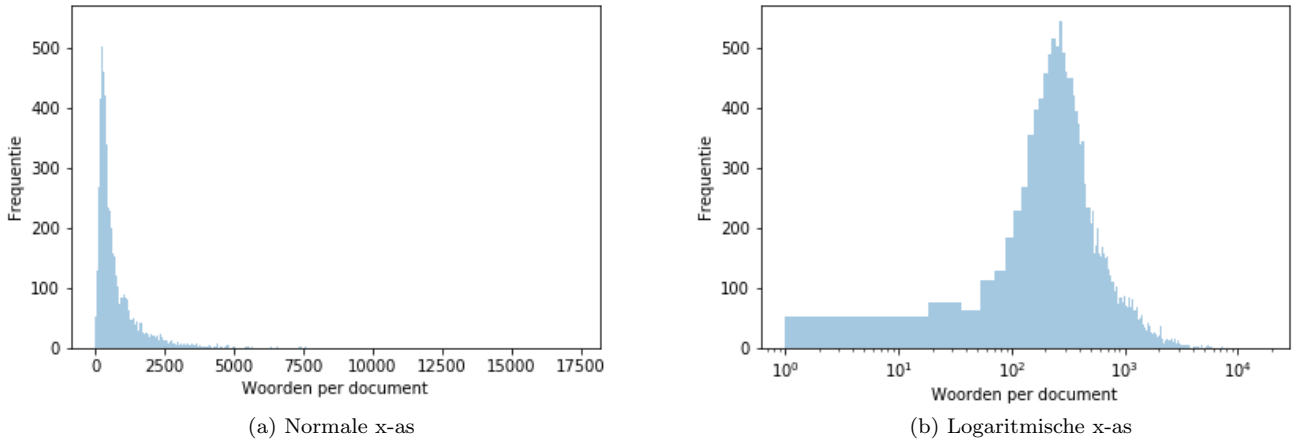
Deze dataset bestaat uit een aantal soorten spreekbeurten; debatbijdragen, interrupties en antwoorden. Een debatbijdrage is de eerste onafgebroken spreekbeurt die een spreker geeft achter een spreekgestoelte, aangeduid in de xml-file met het attribuut *nieuw="ja"*. Dit kan een bijdrage in een debat zijn of een vraag tijdens een vragen uur. 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. Gerelateerde onderzoeken voegen vaak alle teksten van een spreker samen tot één document. Dit is alleen niet mogelijk met de hoeveelheid kleine partijen in de Tweede Kamer, die dan niet altijd in een training of test set zijn vertegenwoordigd. Daarom is in dit onderzoek ervoor gekozen om een debatbijdrage met alle bijbehorende antwoorden samen te voegen tot één document voor de classificatie.

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

Deze dataset bevat 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 is en er overlap zit met hun oorspronkelijke of gelieerde partij, zijn deze er uit gehaald. 50PLUS is in 2014 [8] 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 verschillen in grootte. De distributie van documentgrootte

273 lijkt op een lognormale verdeling, maar met een Kolmogorov-Smirnov test is hier
 274 geen bewijs voor gevonden [9].



Figuur 1: Aantal woorden per document

275 Om toch de uitschieters er uit te halen, is aangenomen dat de distributie
 276 wel lognormaal verdeeld is en zijn daarmee de documenten buiten het betrouw-
 277 baarheidsinterval van 95% eruit gehaald. De documenten met een lengte van
 278 minimaal 28 en maximaal 1492 woorden bleven daarmee over. De gemiddelde
 279 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

280 Deze 14899 documenten zijn verdeeld over 2984 debatten, waarbij elke
 281 vraag tijdens het vragenuur als één debat gezien wordt. Op basis van de aan-
 282 tallen 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).

3.2 Methoden

3.2.1 DV1: Beste classificatiemethode

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

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

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

Support Vector Machines en Logistische Regressie De meest voorkomende techniek in vergelijkbaar onderzoek is Support Vector Machine (SVM). Een andere techniek die gebruikt wordt is logistische regressie. Beide kennen een eigen implementatie in scikit-learn, maar deze implementaties zijn niet efficiënt met grote datasets. Om deze reden is er in beide gevallen voor gekozen om

gebruik te maken van de functie `SGDClassifier`, die beide technieken leert met *stochastic gradient descent learning*. Voor regularisatie is hier geëxperimenteerd 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 geclassificeerd [11].

	Behorend tot partij	Niet behorend tot partij
Geclassificeerd als partij	<i>true positive (tp)</i>	<i>false positive (fp)</i>
Niet geclassificeerd 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 geclassificeerd is. *Accuracy* wordt voor de hele classificatie gedaan en niet per klasse. *Precision* is het percentage van documenten geclassificeerd 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 geclassificeerd 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 veel documenten belangrijker zijn. Als een classificatie kleine partijen groten-deels fout classificeert, kan deze score alsnog hoog zijn. In het geval van meer dan twee partijen is dit hetzelfde als *accuracy*.

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 geclassificeerde documenten hoog scoren door vooral kleine partijen goed te classificeren.

Als laatste is er dan nog *gewogen*, deze berekent net als *macro* de scores per partij, 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 *accuracy* en het nadeel van *macro* te groot is omdat de partijen nogal variëren in grootte, is gekozen voor *gewogen F_1* score naast *accuracy*.

3.2.2 DV2: Invloed van namen

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

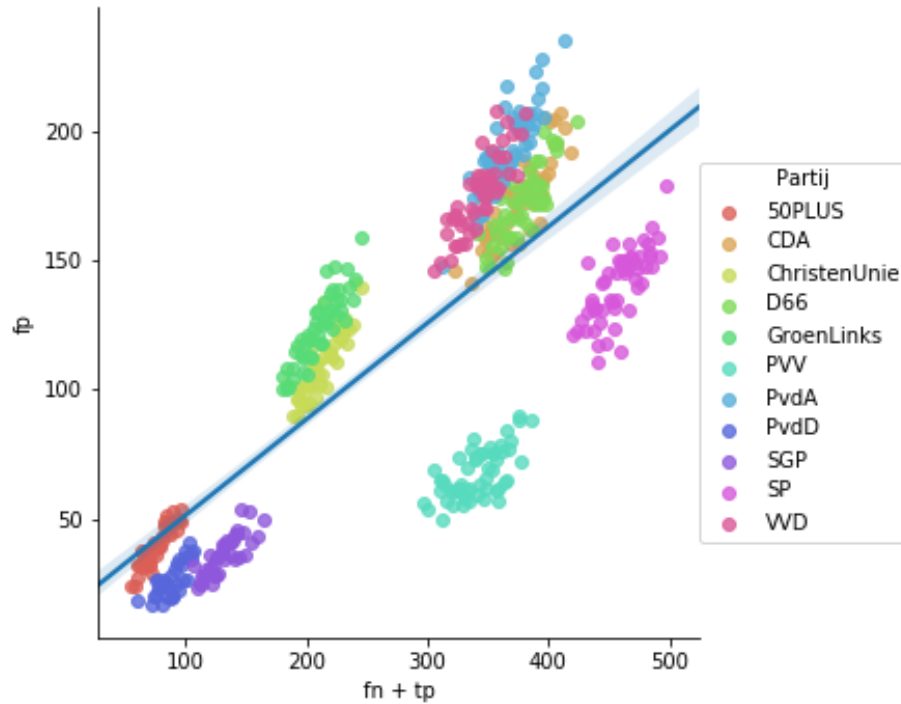
Ook wordt gekeken naar classificatie met alleen partijnamen en achternamen van Kamerleden. Alle andere woorden worden weggehaald. Namen van Kamerleden en partijen die niet aan elkaar geschreven worden, zoals *Partij van de Arbeid*, worden aan elkaar geschreven zodat het één feature wordt. Doordat alle andere woorden weggehaald zijn, worden de bi- en trigrams combinaties van namen die zinnen uit elkaar kunnen staan, dus die niet meer informatie geven dan unigrams. Daarom wordt er gebruikt van de classificatiemethode uit deelvraag 1, maar dan met alleen unigrams. Hoge scores voor deze classificatie

408 geven aan dat met alleen namen classificatie goed te doen is en dat dit dus een
 409 grote bijdrage heeft geleverd aan de resultaten uit deelvraag 1.

410 3.2.3 DV3: Oppositie of regering

411 Om deze deelvraag te beantwoorden zijn drie experimenten uitgevoerd. Twee
 412 daarvan zijn gebaseerd op experimenten uit Hirst et al. [2] voor dezelfde vraag.
 413 De derde is ontwikkeld voor dit onderzoek. Met deze laatste wordt begonnen.
 414 Bij deze deelvraag is de classificatiemethode uit deelvraag 2 gebruikt.

415 Als er een afhankelijkheid is van partij-status, dan is te verwachten dat
 416 het aantal misclassificaties minus verwachte waarde binnen regeringspartijen en
 417 binnen oppositiepartijen hoger ligt dan tussen oppositiepartijen en regeringspar-
 418 tijen. De verwachte waarde is afhankelijk van het aantal documenten van een
 419 partij in de training set [12]. Aangezien de test set uit dezelfde set als de training
 420 is gehaald, is de verwachte waarde ook afhankelijk van het aantal documenten
 421 van een partij in de test set. Uit de voorverkenning (op basis van resultaten uit
 422 deelvraag 1 en 2) blijkt deze correlatie tussen het aantal *false positives* van een
 423 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 test en train set. De Pearson correlatie is 0.77 en de p-waarde 5.40×10^{-101} .

424 Op basis van dit verband definiëren we het verwachte aantal documenten

425 van partij i die foutief geclassificeerd worden als partij j als

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

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

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

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

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

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

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

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

457 Als vergelijkingsmateriaal is voor deze experimenten een tweede dataset
 458 nodig uit een ander kabinet. Hiervoor is het wenselijk dat dit kabinet bestaat uit
 459 andere partijen dan kabinet-Rutte II. Daarnaast is het ook wenselijk als het niet
 460 te ver terug is, zodat onderwerpen en taalgebruik enigszins overeenkomstig zijn.
 461 Omdat kabinet-Rutte I een minderheidskabinet was met een bijzondere partij-
 462 status voor de PVV, is ervoor gekozen om de Handelingen van de Tweede Kamer
 463 tijdens het missionaire kabinet-Balkenende IV (22 februari 2007 tot 20 februari
 464 2010) te gebruiken. Dit kabinet bestond uit CDA, PvdA en ChristenUnie.

465 De partij 50PLUS bestond nog niet gedurende kabinet-Balkenende IV,
 466 dus documenten van deze partij zijn weggelaten. Verder heeft dezelfde verwer-
 467 king van data plaatsgevonden, zoals beschreven in 3.1. Alleen de minimum- en
 468 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 as

Als de classificatie afhankelijk is van positie op de links-rechts as dan is het te verwachten dat, net als bij deelvraag 3, de misclassificaties minus de verwachte waarde groter zijn als twee partijen dichterbij elkaar staan op de links-rechts as. Daarvoor zal wederom formule 5 gebruikt worden als verwachte waarde en 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 de links-rechts as, op basis van het verkiezingsprogramma van dat jaar, in dit geval dus van 2012.

Tabel 3: 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

Er wordt vervolgens gekeken door middel van een Pearson correlatie toets of er een correlatie is tussen de error van twee partijen en de afstand op de links-rechts as van die partij. Het significantieniveau (α) hiervoor is opnieuw

483 0.05. De nulhypothese is dat er geen negatieve correlatie is tussen de error en
484 de afstand op de links-rechts as. De alternatieve hypothese is dat er wel een
485 negatieve correlatie is tussen de error en de afstand op de links-rechts as.

486 Als uit deelvraag 3 blijkt dat partij-status invloed heeft op de error, zal
487 bovenstaande methode ook uitgevoerd worden voor de aparte combinaties; bin-
488 nen oppositie en tussen regeringspartij en oppositiepartij. Binnen regering is
489 niet mogelijk aangezien dat maar één afstand is.

490 3.2.5 DV5: Woordgebruik van sprekers

491 De vorige classificaties trainden op documenten en werden getest op andere
492 documenten, maar wel van dezelfde sprekers als uit de training set. Naast
493 de ideologie kan de classificatie daarom ook getraind zijn op het taalgebruik
494 van sprekers. Als een Kamerlid bijvoorbeeld een woord regelmatig in speeches
495 gebruikt, maar niet wordt gebruikt door zijn partijgenoten, wordt dit wel gezien
496 als een belangrijk woord voor de classificatie naar partij-affiliatie. Hirst et al.
497 [2] plaatsten al een soortgelijke kanttekening bij de resultaten van Diermeier et
498 al.

499 Om te kijken of dit effect er is, wordt er opnieuw een classificatie gedaan
500 met de methode uit deelvraag 2. Ditmaal worden alleen niet de individuele
501 documenten verdeeld over de training en test set, maar worden de Kamerleden,
502 met bijbehorende documenten, verdeeld over de training en test set. Als taalge-
503 bruik van een spreker in de training set voorheen invloed had op de classificatie,
504 zal dat nu geen effect meer hebben omdat er geen documenten van die spreker
505 meer voorkomen in de test set. De meest karakteristieke woorden uit de resulta-
506 ten van deelvraag 2 suggereren dat woordgebruik van Kamerleden invloed heeft
507 (zie tabel 5). De hypothese is daarom ook dat deze nieuwe classificatie lagere
508 scores vindt.

509 4 Resultaten

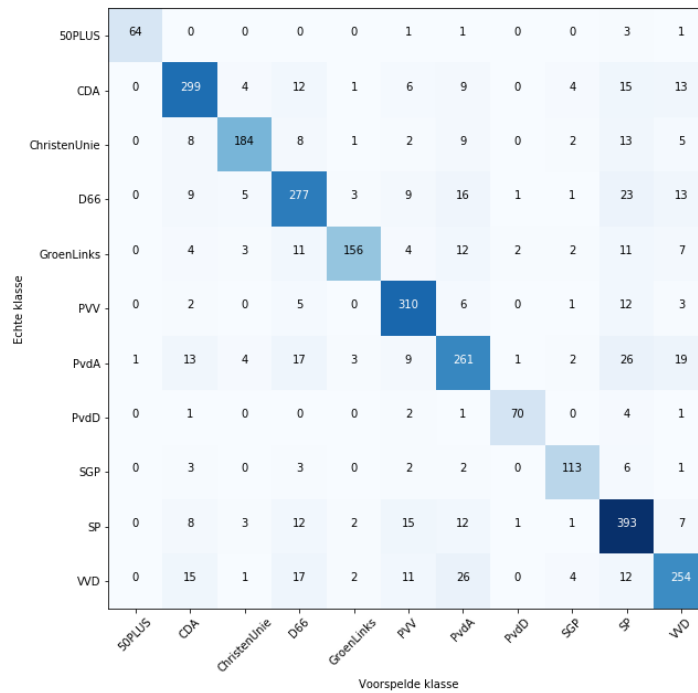
510 4.1 DV1: Beste classificatiemethode

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

517 Tabel 4 laat de scores zien per partij met het aantal documenten in de
518 test set. De *accuracy* voor deze classificatie is 0.80. De F_1 scores per partij
519 liggen tussen de 0.7 en 0.9. De partijen met een sterke focus op één onderwerp,
520 50PLUS, PVV en PvdD, als ook de SGP hebben hoge scores, terwijl de coa-
521 litiepartijen, VVD en PvdA, lagere scores hebben. Figuur 3 laat zien waar de
522 fouten in deze classificatie zitten. De meest karakteristieke features per partij
523 zijn te zien in tabel 5. Met meest karakteristiek worden de n-grams bedoeld die
524 de hoogste coëfficiënt hebben in de classificatie en die dus relatief het meeste
525 belangrijk zijn voor de classificatie van een partij. Hierin is te zien dat vrijwel
526 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

4.2 DV2: Invloed van namen

In tabel 5 was al te zien dat de meest karakteristieke n-grams voornamelijk achternamen van Kamerleden of partijnamen bevatten. In tabel 6 zijn de scores te zien voor een classificatie met alleen achternamen van Kamerleden en partijnamen. De *accuracy* is 0.61. De scores zijn gedaald ten opzichte van de resultaten 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

533 In tabel 7 zijn de F_1 scores te zien van classificatie met achternamen van
534 Kamerleden en partijnamen vervangen. De *accuracy* hiervan is 0.58. De scores
535 zijn aanzienlijk lager dan die uit deelvraag 1 en ook nog lager dan van de clas-
536 sificatie met alleen namen. Wel zijn de scores nog ruim hoger dan de baseline.
537 In tabel 8 is vervolgens te zien welke n-grams het meest karakteristiek zijn per
538 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.

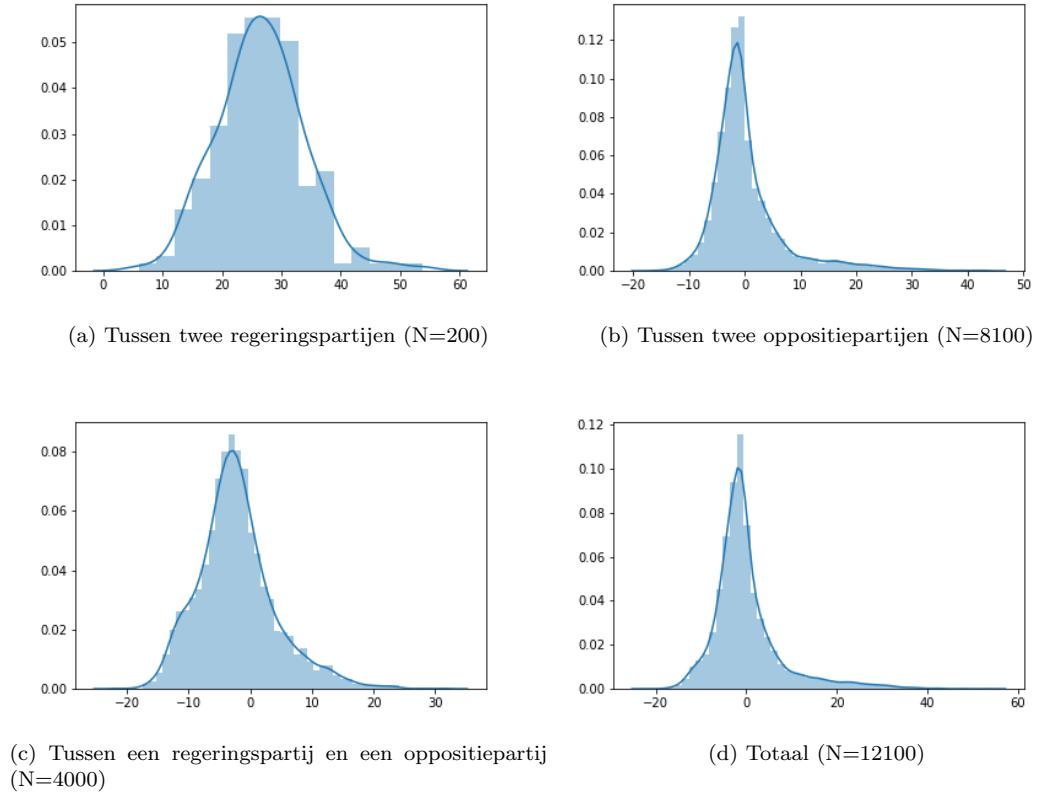
50PLUS	CDA	ChristenUnie	D66	GroenLinks
gepensioneerden	PARTIJ fractie	gezinnen	mijn fractie	zou
ouderen	inwoners	inderdaad	mijn	kamer hierover te
oudere	regering	ik hoop	natuurlijk	schone energie
koopkrachtontwikkeling	PARTIJ	voedselverspilling	fractie	persoonsgebonden
plussers	hier	hoop	buitengewoon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

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	milieu	mevrouw de	mening dat	volgens mij
miljard	kinderen	dierenwelzijn	beantwoording	voorstellen	aruba
natuurlijk	lagere overheden	natuur	punt	van mening dat	speelveld
islam	goede	de bio	voor de beantwoording	bezuinigingen	aangegeven
de islam	circulaire	bio	de beantwoording	de bevolking	volgens
al	tevens	industrie	wel	mensen	en
asielzoekers	circulaire economie	bio industrie	allerlei	huurders	PARTIJ fractie
verzorgingshuizen	mijn partij	klimaatverandering	bewindslieden	van mening	essentieel
brussel	regering tevens	de bio industrie	nadrukkelijk	bevolking	regelgeving

539 4.3 DV3: Oppositie of regering

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



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

Voor alle distributies kan de nulhypothese verworpen worden dat deze normaal verdeeld zijn. In tabel 9 is vervolgens te zien dat er een significant verschil is tussen de distributies binnen regering en 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	7.04×10^{-124}	717042
Tussen twee oppositiepartijen	4.4×10^{-108}	16328471

In tabel 10 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 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	de premier	fractie van PARTIJ	onze
fractie	de fractie	ik hoop	de fractie	niet
wij hebben	fractie van	arbeidsmarkt	de fractie van	deze
KAMERLID	mijn fractie	de arbeidsmarkt	fractie van	immigratie
dank	geweest	hoop	politieke	natuurlijk
PARTIJ fractie heeft	moment	hij	premier	politie
zorgvuldige	termijn	schone energie	deal	de burgers
ons	verschillende	plannen	ik	burger
buitengewoon	beantwoording	kunnen	en	door

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	zegt	PARTIJ
vrouwen	natuur	wel	mensen	onze fractie
belangrijk	bio industrie	beantwoording	problemen	PARTIJ fractie
kinderen	de bio industrie	de bewindslieden	niet	fractie
of	de bio	bewindslieden	vandaar	ondernemers
punt	bio	de voorzitter	is	je
goed	veehouderij	toch	de mensen	want
onderzoek	dierenwelzijn	mijn	bureaucratie	awbz
roc	industrie	een	verdrag	eens
vragen	de natuur	natuurlijk	waarom	antwoorden

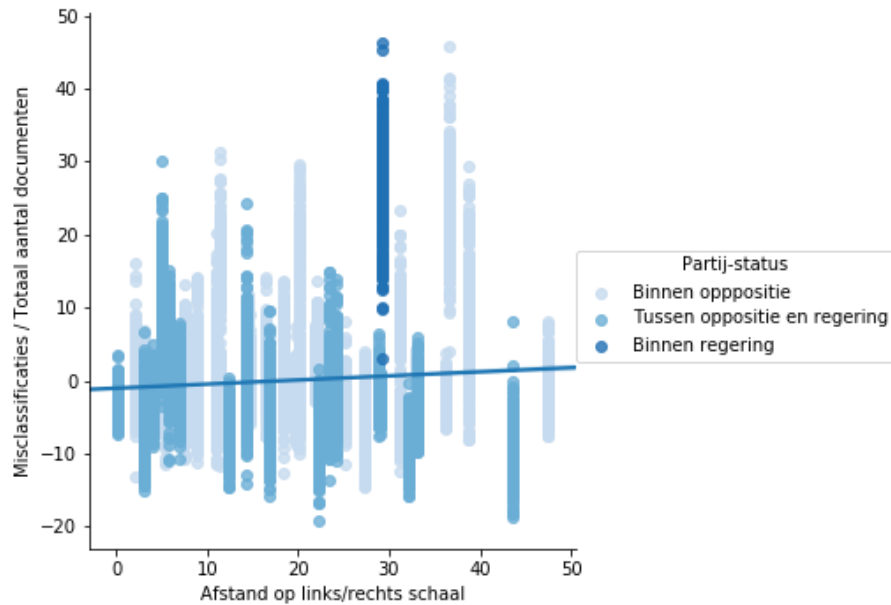
549 In tabel 11 zijn de resultaten van de classificatiescores te zien waarbij de
 550 classificatie getraind is op een zittingsperiode, maar getest op een andere. De
 551 resultaten zijn sterk gedaald, maar nog boven de baseline. De daling verschilt
 552 enorm per partij en zittingsperiode met dalingen van F_1 scores tussen 12 en
 553 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

554 4.4 DV4: Links-rechts as

555 In tabel 5 is de error te zien ten opzichte van de afstand op de links-rechts as.



Figuur 5: 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. 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 5. Daarom wordt er in tabel 12 ook gekeken naar de correlatie tussen afstand op de links-rechts as en error binnen oppositie en tussen regerings- en oppositiepartij. Beide correlaties zijn statistisch significant op het significantieniveau van 0.01, maar opvallend genoeg in tegengestelde richting.

Tabel 12: 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 13 staan de scores van classificatie waarbij de Kamerleden verdeeld zijn over de training en test set. De scores zijn hierbij amper hoger dan de baseline.

Tabel 13: 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	

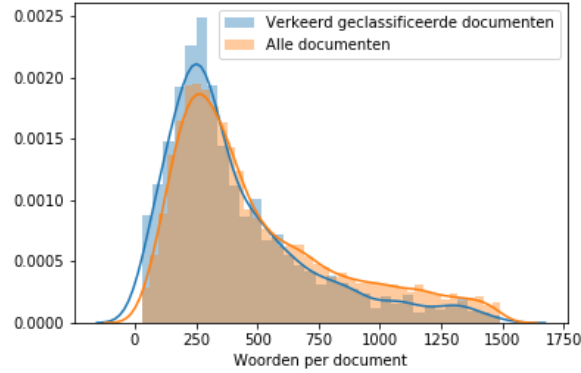
5 Discussie

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%.

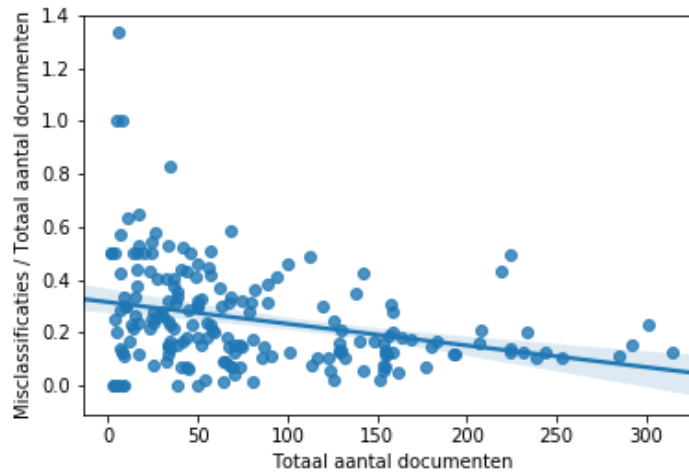
Het onderzoek van Hirst et al. vond dat resultaten afhankelijk kunnen zijn van documentgrootte. Alle documenten in dit onderzoek zijn kleiner dan de grootste documentgrootte uit het onderzoek van Hirst et al. en ook de minimale documentgrootte ligt lager dan de kleinste documentgrootte uit dat onderzoek. 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 6: 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 7: 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} .

598 Dit versterkt het vermoeden dat de classificatie mede plaatsvindt op basis
599 van woordgebruik van individuele sprekers, zoals besproken wordt in deelvraag
600 5.

601 5.2 DV2: Invloed van namen

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

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

616 De classificatiemethode die gebruikt is in deze deelvraag, is gebaseerd op
617 de beste methode voor de dataset uit deelvraag 1. Hierin was gevonden dat een
618 combinatie van uni-, bi- en trigrams het beste resultaat opleverde. In tabel 5
619 is te zien dat trigrams behoren tot de meest karakteristieke n-grams, hoewel de
620 woorden in trigrams vaak overlappen met uni- en bigrams. In tabel 8 daaren-
621 tegen zijn er nog maar een paar trigrams, welke grotendeels procedurele zinnen
622 zijn of toevoeging van een lidwoord op een uni- of bigram. Dit verschil sugge-
623 reert dat trigrams minder belangrijk zijn in de classificatie zonder de namen,
624 dus de classificatiemethode uit deelvraag 1 niet het beste is voor deze classifica-
625 tie. In vervolgonderzoek kan de opzet van deelvraag 1 toegepast worden op de
626 classificatie zonder de namen, om zo te komen tot een classificatiemethode die
627 het beste resultaat oplevert op de classificatie zonder namen.

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

635 5.3 DV3: Oppositie of regering

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

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

641 De verwachting was dat de error normaal verdeeld zou zijn. De verde-
642 lingen uit figuur 4 hebben globaal wel de vorm van een normaal verdeling. In
643 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 aanleiding te 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. En deze verwachte waarde en de daar uit volgende error zijn 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/rechts 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 14: 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>hun</i> <i>collega KAMERLID</i> <i>in</i> <i>aanpak</i> <i>collega</i>	<i>algemeen</i> <i>algemeen overleg</i> <i>toezegging</i> <i>helder</i> <i>overleg</i> <i>aangegeven</i> <i>voor</i> <i>voor PARTIJ</i>
	ChristenUnie	<i>mijn</i> <i>waarop</i> <i>blij</i> <i>collega KAMERLID</i> <i>erg</i>	<i>gaan</i> <i>termijn</i> <i>blij met de</i> <i>volgens</i> <i>volgens mij</i> <i>blij</i> <i>beantwoording</i>
	PvdA		<i>volgens</i> <i>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

668 is het grootst bij VVD, maar valt mee bij de twee andere partijen die gewis-
669 seld zijn van partij-status, ChristenUnie en CDA. Daarnaast is de daling ook
670 heel sterk bij oppositiepartijen GroenLinks en D66, alsook de regeringspartij
671 in beide kabinetten, PvdA. Dat de daling niet consequent groter is bij partijen
672 die gewisseld zijn van partij-status, suggereert dat de invloed van partij-status
673 beperkt is op de classificatie.

674 Dat de experimenten uit Hirst et al. in hun onderzoek wel invloed vin-
675 den, maar in dit onderzoek niet kan komen doordat hun onderzoek zich richt
676 op binaire classificatie, terwijl dit onderzoek meerdere partijen heeft. Zo kan
677 het ontbreken van gemeenschappelijke n-grams komen doordat regeringspartijen
678 zich ook van elkaar moeten onderscheiden in dit onderzoek, waarvoor n-grams
679 die relevant zijn voor partij-status weinig effect hebben, terwijl in het onderzoek
680 van Hirst et al. de regeringspartij alleen onderscheiden hoeft te worden van de
681 oppositiepartij. Daarnaast verklaren zij dat een daling tussen twee zittingsperi-
682 oden met een wisseling van partij-status het gevolg is van deze wisseling, terwijl
683 in dit onderzoek gekeken kan worden naar dit effect voor partijen die wel en
684 niet gewisseld zijn.

685 **5.4 DV4: Links-rechts as**

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

692 Alle correlaties zijn statistisch significant, maar de Pearson correlatie sug-
693 gereert dat het verband zwak is en de effectgrootte klein tot medium. Daarnaast
694 is het ook opvallend dat de twee combinaties van partij-statussen een andere cor-
695 relatierichting hebben. Dit suggereert dat de statistische significantie het gevolg
696 is van de grote steekproef en maar een klein effect [15].

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

703 **5.5 DV5: Woordgebruik van sprekers**

704 De resultaten uit tabel 13 zijn laag, amper hoger dan de baseline. Dit sugge-
705 reert inderdaad dat eerdere classificaties in grote mate toch afhankelijk waren
706 van het woordgebruik van sprekers. Dit is opmerkelijk aangezien vergelijkbare
707 onderzoeken dit effect niet vinden. De meest karakteristieke n-grams van deze
708 classificatie wijken daarnaast grotendeels niet af van die uit tabel 8.

709 Een alternatieve verklaring is dat de classificatie nu mede op basis van
710 woordvoerderschap is. Per onderwerp heeft een partij vaak maar één woord-
711 voerder, met uitzonderingen van wijzigingen in de fractie. Het is aannemelijk
712 dat het taalgebruik afhankelijk is van woordvoerderschap, aangezien er andere
713 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 spreekbeurten 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.

5.6 Algemeen

Het vergelijken van deze resultaten met vergelijkbaar werk is ingewikkeld, aangezien de keuzes en eigenschappen van die onderzoeken het niet een één-op-één vergelijking maken. Voorbeelden hiervan zijn de taal, het parlement, de documentgrootte, baselines, behouden of weglaten van namen, een spreker als document zien en het trainen en testen op dezelfde spreker. Hoewel de resultaten in sommige gevallen lager zijn dan die uit vergelijkbaar werk, is het belangrijk hier rekening mee te houden. Een vervolgonderzoek zou daarom dit onderzoek kunnen reproduceren op een ander parlement om daarmee te kunnen vergelijken.

Dit onderzoek richtte zich hoofdzakelijk op de Handelingen gedurende kabinet-Rutte II. Om te kijken in hoeverre het mogelijk is om deze conclusie door te trekken naar de algemene Handelingen van de Tweede Kamer, kan er in vervolgonderzoek gekeken worden naar meerdere zittingsperiodes. 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). Daarnaast vindt dit onderzoek geen aanwijzingen dat de classificatie afhankelijk is van positie op de links-rechts as. Als rekening wordt gehouden met woordgebruik van individuele Kamerleden, daalt

de *accuracy* verder naar.... Daarmee lijkt de classificatie naar partij-affiliatie in grote mate niet het gevolg van ideologie. Deze conclusie trekt daarmee de bruikbaarheid van tekstclassificatie voor het vinden van een relatie tussen woordgebruik en ideologie in twijfel. Op een aantal punten wijken de bevindingen van dit onderzoek af van vergelijkbare onderzoeken [2, 4]. Voor een vervolgonderzoek kan dit onderzoek uitgebreid worden om de validiteit van dit onderzoek te verhogen.

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] Beia 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 Vellidal. 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, 2014.
- [8] Laura Klompenhouwer. Extra ledenvergadering 50plus om splitsing. *NRC Handelsblad*, June 2014.
- [9] Eric Jones, Travis Oliphant, Pearu Peterson, et al. SciPy: Open source scientific tools for Python, 2001.
- [10] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [11] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. *Introduction to Information Retrieval*. Cambridge University Press, New York, NY, USA, 2008.

- 800 [12] Mahendra Sahare and Hitesh Gupta. A review of multi-class classifica-
801 tion for imbalanced data. *International Journal of Advanced Computer*
802 *Research*, 2(3), 2012.
- 803 [13] *NIST/SEMATECH e-Handbook of Statistical Methods*. NIST/SEMA-
804 TECH, April 2012.
- 805 [14] Andrea Volkens, Pola Lehmann, Theres Matthieß, Nicolas Merz, Sven Re-
806 gel, and Bernhard Weßels. The manifesto data collection. manifesto project
807 (mrg/cmp/marpor). version 2017b, 2017.
- 808 [15] Joseph F. Hair, Jr., Rolph E. Anderson, Ronald L. Tatham, and William C.
809 Black. *Multivariate Data Analysis (6th Ed.)*. Prentice-Hall, Inc., Upper
810 Saddle River, NJ, USA, 2006.