

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



14

Abstract

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

In verschillende onderzoeken zijn parlementaire teksten geclassificeerd naar partij-affiliatie. Dit onderzoek is op zoek gegaan naar de beste classificatiemethode voor de Handelingen van de Tweede Kamer. Daarnaast is gekeken naar in hoeverre dit 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 *accuracy* 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 ontbreken daarentegen. Als laatste daalt de *accuracy* 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	Contents	
32	1 Introductie	3
33	2 Gerelateerd werk	4
34	2.1 Tekstclassificatie van parlementaire teksten	4
35	2.2 Classificatiemethoden	5
36	2.3 Invloed van partijnamen of sprekersnamen	6
37	2.4 Invloed van oppositie of regering	6
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	28
56	5.4 DV4: Links-rechts as	30
57	5.5 DV5: Woordgebruik van sprekers	30
58	5.6 Algemeen	31
59	6 Conclusies	31

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 [1, 2].

69 In diverse landen zijn al onderzoeken gedaan naar het classificeren naar
70 partij-affiliatie op basis van teksten van politici [1, 3, 2, 4, 5]. Met deze tekst-
71 classificatie naar partij-affiliatie proberen onderzoekers uit te vinden in hoeverre
72 ideologie terug te vinden is in teksten van politici. De resultaten van de tekst-
73 classificaties zijn in de meeste gevallen ruim boven de baseline. Diverse onder-
74 zoeken wijzen daarentegen ook naar redenen dat dit niet alleen het gevolg is van
75 ideologie. Zo suggereren de resultaten van Hirst et al. [2] dat de partij-status
76 (oppositie tegenover regering) van invloed is op de classificatie. Daarnaast laat
77 dit onderzoek ook zien dat partijnamen een grote invloed hebben op de classi-
78 ficatie.

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

81 Dit onderzoek richt zich daarom op meerdere classificatiemethoden en
82 daarnaast specifiek op de Nederlandse politiek. De onderzoeksvraag luidt daarom
83 dus ook: "In hoeverre is classificatie naar partij-affiliatie aan de hand van spreek-
84 beurten in de Tweede Kamer het gevolg van ideologie?"

85 Deze vraag wordt beantwoord door de antwoorden te vinden op de vol-
86 gende deelvragen:

- 87 1. Wat is het beste classificatiemodel voor classificatie naar partij-affiliatie
88 in de Tweede Kamer en wat is het resultaat van dit model?
- 89 2. In hoeverre is deze classificatie afhankelijk van achternamen van Kamerleden
90 en partijen?
- 91 3. In hoeverre wordt deze classificatie bepaald door partij-status (oppositie
92 of regering)?
- 93 4. In hoeverre wordt deze classificatie bepaald door positie op de links-rechts
94 as?
- 95 5. In hoeverre wordt deze classificatie bepaald door woordgebruik van sprek-
96 ers?

97 Op basis van vergelijkbare onderzoeken is de hypothese dat in grote mate de
98 classificatie niet het gevolg is van ideologie, maar in beperkte mate wel.

99 Voor de eerste deelvraag zullen Support Vector Machine, Logistische Re-
100 gressie en Naive Bayes met verschillende parameters vergeleken worden aan
101 de hand van *accuracy* en F_1 score. Bij de tweede deelvraag wordt gekeken
102 naar classificatie zonder achternamen van Kamerleden en partijnamen of met
103 alleen achternamen van Kamerleden en partijnamen. De derde vraag bestaat

104 uit meerdere experimenten, waarin gekeken zal worden naar de hoeveelheid mis-
105 classificaties binnen regering of oppositie tegenover tussen regering en oppositie.
106 Daarnaast zal gekeken worden naar overlap in woordgebruik binnen regering en
107 verschil in scores als een partij gewisseld is van partij-status. Bij de vierde
108 vraag zal gekeken worden naar een verband tussen misclassificaties en afstand
109 tussen twee partijen op de links-rechts as. Als laatste zal voor de vijfde vraag
110 de classificatie herhaald worden met Kamerleden verdeeld over training en test
111 set.

112 **Overzicht van scriptie** Sectie 2 bevat vergelijkbare onderzoeken in andere
113 parlementen. Sectie 3 bevat de methodologie van de verschillende deelvragen.
114 Sectie 4 bevat de resultaten. Sectie 5 bevat de evaluatie van zowel de resul-
115 taten als de methodologie. Sectie 6 bevat ten slotte het antwoord op de onder-
116 zoeksvraag.

117 2 Gerelateerd werk

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

122 In dit hoofdstuk zullen verschillende onderzoeken behandeld worden die
123 tekstclassificatie hebben toegepast op parlementaire teksten. Eerst zullen de
124 onderzoeken algemeen besproken worden. Vervolgens is uitgebreider gekeken
125 worden de effecten van verschillende classificatiemethoden. In de latere secties
126 worden de aspecten besproken die in vergelijkbare onderzoeken genoemd worden
127 als van invloed op de classificatie.

128 2.1 Tekstclassificatie van parlementaire teksten

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

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

143 Yu et al. [5] richtten zich vervolgens op zowel het Amerikaanse Huis van
144 Afgevaardigden als de Senaat in 2005. Een document was in dit onderzoek de
145 verzameling van alle speeches van een congreslid en het label de partij. Voor het
146 Huis van Afgevaardigden vonden ze een *accuracy* van 80.1% (baseline van 51.5%)
147 en voor de Senaat 86.0 % (baseline van 55.0%). Ze testten hun classificaties ook

148 op de andere kamer. Van Huis van Afgevaardigden naar Senaat leverde dit een
149 *accuracy* op van 88.0% (baseline van 55.0%) en andersom 67.6% (baseline van
150 51.5%). Hun verklaring voor dit verschil was dat het Huis van Afgevaardigden
151 sterker verdeeld is langs partijlijnen.

152 Vervolgens herhaalden ze de classificaties op het Huis van Afgevaardigden
153 uit 2015, maar testten ditmaal op de Senaat elk jaar tussen 1989 en 2006 afzon-
154 derlijk. Hierin was een stijging in *accuracy* van 60.0% (baseline van 55.0%) in
155 1989 naar 87.0% (baseline van 55.0%) in 2006 te zien, maar met twee duidelijke
156 dalen. Ze presenteren twee mogelijke verklaringen voor de trend; het veranderen
157 van de onderwerpen en het sterker verdeeld worden van het Congres.

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

163 Het onderzoek bevat ook een classificatie van het Europees Parlement.
164 Hierbij voegden ze alle teksten van een parlements lid bij elkaar en deelden die
165 op in documenten van gelijke grootte. Voor documentgrootte van 267 woorden
166 werd een *accuracy* van 44.0% gevonden oplopend tot 61.8% (baseline van 38-
167 39%) voor documentgrootte van 6666.

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

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

177 Høyland et al. trainden een classificatie voor partij-affiliatie op basis van
178 teksten van het vijfde Europese Parlement (1999-2004) en testten vervolgens
179 op het zesde Europese Parlement (2004-2009) [7]. Alle teksten van een spreker
180 waren samengevoegd tot één document. 40% van de sprekers in de test set zaten
181 ook in de training set. Hier werd een *macro* F_1 score van 0.464 (baseline van
182 0.097) en *accuracy* van 0.551 (baseline van 0.410) verkregen. De baseline is in
183 dit onderzoek op basis van altijd classificeren als grootste partij, terwijl voor F_1
184 score de baseline hoger ligt als hiervoor gekozen wordt voor gokken gewogen bij
185 grootte van een klasse.

186 2.2 Classificatiemethoden

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

190 Yu et al. [5] maakten gebruik van Support Vector Machines en Naive
191 Bayes, waarvan de varianten multinomial en Bernoulli. De features waren uni-
192 grams, met minimale woordfrequentie van drie en de top 50 meest voorkomende
193 woorden weggelaten. Voor de wegingen van de features bij Support Vector
194 Machines werd geëxperimenteerd met *boolean*, *tf-norm* en *tf-idf*. De beste clas-
195 sificatiemethode was afhankelijk van de dataset. Voor het Huis van Afgevaardig-

den was het Support Vector Machines met als weging *tf-idf* en voor de Senaat Bernouilli 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 variaties maakten in hun onderzoek geen grote verschillen en uiteindelijk is gekozen voor het niet stemmen, het weglaten van woorden die in minder dan vijf documenten voorkomen en resultaten van zowel met als zonder de top 500 meest frequente woorden. Daarnaast werd geëxperimenteerd met vier wegenen voor woorden: *boolean*, *tf*, *tf-norm* en *tf-idf*, waarvan *tf-idf* het beste resultaat opleverde.

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

Ferreira maakten gebruik van twee classificatiemethoden: Logistische regressie en *margin-infused relaxed algorithm* (MIRA) [3]. Logistische regressie werd aangevuld met *group Lasso* regularisatie, wat het beste resultaat opleverde. Voor wegenen van woorden werd geëxperimenteerd met *tf*, *tf-idf*, Δ -*tf-idf* en Δ -*BM-25*. Daarnaast wordt er gebruik gemaakt van woordclustering, *Concise Semantic Analysis* en stylometrische eigenschappen. Op *Part-Of-Speech tagging* na hadden stylometrische eigenschappen een duidelijke negatieve invloed op de classificatie.

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

2.3 Invloed van partijnamen of sprekersnamen

Diermeier et al. [4] lieten de namen van de sprekers en verwijzingen naar staten die de senatoren representeren weg, omdat deze volgens hen de classificatie te makkelijk zouden maken. Hirst et al. [2] vonden inderdaad dat partijnamen - en het weglaten daarvan - bij het Europees Parlement een grote invloed hebben op de classificatie. Bij het Europees Parlement was te zien dat een spreker de eigennaam 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. [2] 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

242 bij het 39e parlement bij de conservatieven (regering) te vinden waren. An-
243 dersom gebeurde hetzelfde met één van de tien woorden van de conservatieven
244 (oppositie) in het 36e parlement naar liberalen (oppositie) in het 39e parlement.

245 In hetzelfde onderzoek trainden ze ook hun classificaties op het ene par-
246 lement en testten deze op het andere parlement. Hierbij vonden zij in beide
247 gevallen een *accuracy* ver onder de baseline.

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

250 3 Methodologie

251 3.1 De data

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

262 Deze dataset bestaat uit een aantal soorten spreekbeurten; debatbijdra-
263 gen, interrupties en antwoorden. Een debatbijdrage is de eerste onafgebroken
264 spreekbeurt die een spreker geeft achter een spreekgestoelte, aangeduid in de
265 xml-file met het attribuut *nieuw="ja"*. Dit kan een bijdrage in een debat zijn
266 of een vraag tijdens een vragenuur. Interrupties zijn de vragen die andere politici
267 stellen vanachter de interruptiemicrofoon aan een spreker. De antwoorden zijn
268 vervolgens de reactie van een spreker achter het spreekgestoelte op een inter-
269 ruptie. Aangezien een debatbijdrage geïnterrumped kan worden, kan deze
270 inhoudelijk doorlopen in een antwoord van een spreker. Vergelijkbare onder-
271 zoeken voegen vaak alle teksten van een spreker samen tot één document. Dit
272 was alleen niet mogelijk voor dit onderzoek met de hoeveelheid kleine partijen
273 in de Tweede Kamer, die dan niet altijd in een training of test set zijn vertegen-
274 woordigd. Daarom was in dit onderzoek ervoor gekozen om een debatbijdrage
275 samengevoegd met alle bijbehorende antwoorden te beschouwen als één docu-
276 ment.

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

282 Deze dataset bevatte vervolgens naast de verkozen partijen na de Tweede
283 Kamerverkiezingen van 2012 ook afsplitsingen van die partijen (tien in totaal) en
284 bezoeken van vertegenwoordigingen van Nederlandse partijen uit het Europees
285 Parlement (tien in totaal). Omdat van beide categorieën relatief weinig data
286 was en er overlap zat met hun oorspronkelijke of gelieerde partij, waren deze er
287 uit gehaald. 50PLUS is in 2014 [8] uiteengevallen in twee fracties die aanspraak

288 maakten op de partij-affiliatie 50PLUS. Vanaf dit moment zijn deze documenten
 289 niet meer meegenomen om onduidelijkheid te voorkomen.

290 De documenten verschilden in grootte (aantal woorden). De distributie
 291 van documentgrootte lijkt op een lognormale verdeling, maar met een Kolmogorov-
 292 Smirnov test was hier geen bewijs voor gevonden [9].

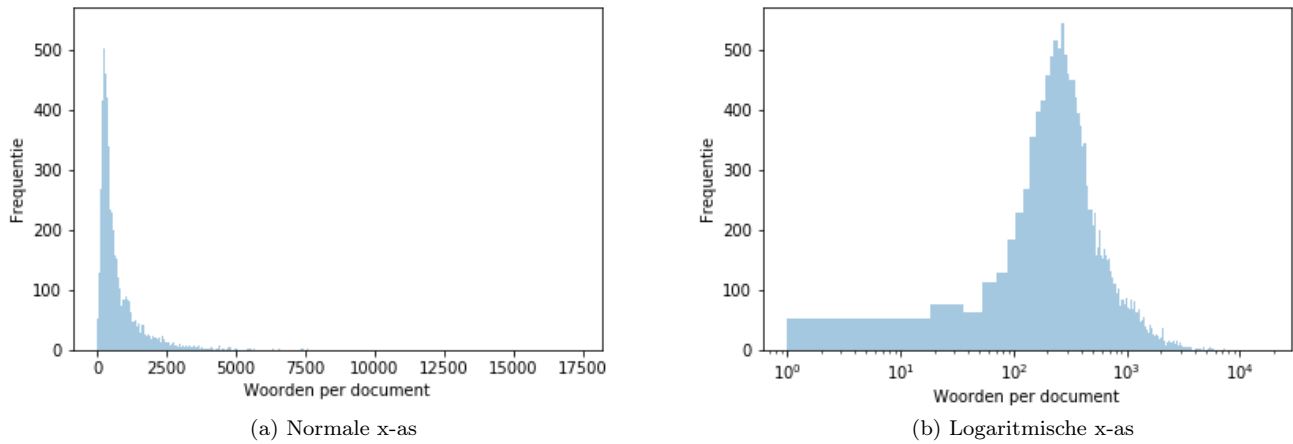


Figure 1: Aantal woorden per document

293 Om toch de uitschieters er uit te halen, was aangenomen dat de distributie
 294 wel lognormaal verdeeld is en waren daarmee de documenten buiten het be-
 295 trouwbaarheidsinterval van 95% eruit gehaald. De documenten met een lengte
 296 van minimaal 28 en maximaal 1492 woorden bleven daarmee over. De gemid-
 297 delde documentlengte is daarna 498 woorden en de mediaan is 386 woorden.

Table 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

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

303 3.2 Methoden

304 3.2.1 DV1: Beste classificatiemethode

305 Om deze deelvraag te beantwoorden zijn een aantal classificatiemethoden vergeleken.
306 Aangezien het niet mogelijk was om alle classificatiemethoden te vergelijken,
307 beperkte dit onderzoek zich tot classificatiemethoden die gebruikt zijn in vergeli-
308 jkbare onderzoeken, zoals besproken in sectie 2.2. Er was ervoor gekozen om
309 alleen gebruik te maken van methoden waarvan reeds implementaties beschik-
310 baar waren in scikit-learn. Voor alle methoden werd gezocht naar de beste
311 parameters, ook wel bekend als een grid search. Deze grid search werd gedaan
312 door 5-fold cross-validation, waarbij de training set steeds 80% was en de test set
313 20% van de totale dataset. Een totaal aantal van 6480 combinaties van metho-
314 den en parameters zijn getest. De verwachting was dat de scores lager zijn dan
315 die gevonden in het gerelateerd werk, omdat de documentgrootte kleiner is en
316 de baseline scores lager zijn.

317 **Pre-processing** Voor pre-processing is gebruik gemaakt van tokenisation en
318 lowercasing. Voor tokenisation is de reguliere expressie
319 $w+$ gebruikt, waardoor alles behalve letters en cijfers weggehaald wordt. Ook is
320 er gevarieerd tussen wel of geen gebruik maken van stemming. In het geval van
321 stemming is gebruik gemaakt van de Snowball Stemmer van de Python NLTK
322 module.

323 **Bag-of-words model** Bag-of-words model is de meest gebruikte representatie
324 van data in vergelijkbare onderzoeken. Deze is ook gebruikt in dit onderzoek.
325 Bij het bag-of-words model wordt elk document gerepresenteerd als een vector,
326 waarbij elke kolom een woord is met een bijbehorende waarde. Voornaamste
327 beperking van dit model is dat het geen rekening houdt met de volgorde van
328 woorden, wat een groot effect kan hebben op de betekenis van een document.

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

340 **Support Vector Machine en Logistische Regressie** De meest voorkomende
341 techniek in vergelijkbaar onderzoek is Support Vector Machine (SVM). Een an-

342 dere techniek die gebruikt wordt, is logistische regressie. Beide hebben een eigen
 343 implementatie in scikit-learn, maar deze implementaties zijn niet efficiënt met
 344 grote datasets. Om deze reden is er in beide gevallen voor gekozen om gebruik
 345 te maken van de functie `SGDClassifier`, die beide technieken leert met *stochastic*
 346 *gradient descent learning*. Voor regularisatie was hier geëxperimenteerd met L1
 347 en L2 regularisatie en een combinatie van beide genaamd Elasticnet. De andere
 348 parameters zijn gelaten op de standaardwaarden van scikit-learn [10]. Een be-
 349 langrijke onaangepaste waarde was die van maximaal aantal iteraties, waarvoor
 350 de scikit-learn standaard 5 is. Volgens scikit-learn convergeert de `SGDClassifier`
 351 rond de $10^6/n$ iteraties waar n het aantal documenten in de training set is. In
 352 het geval van deze dataset zou dat 84 iteraties zijn. Vanwege de grootte van
 353 de grid search was het voor dit onderzoek niet mogelijk het maximaal aantal
 354 iteraties te verhogen tijdens de grid search. De resultaten buiten de grid search
 355 zullen gebaseerd zijn op een maximaal aantal iteraties van 100.

356 **Naive Bayes** Een andere techniek die gebruikt wordt voor politieke tekstclas-
 357 sificatie is Naive Bayes. Dit algoritme neemt aan dat elke *feature* onafhankelijk
 358 is ten op zichte van de rest. Dit is bij tekstclassificatie vaak niet het geval
 359 omdat het gebruik van sommige woorden gepaard kan gaan met het gebruik
 360 van andere woorden. Daarnaast is het gebruik van meerdere n-grams in een
 361 classificatie schending van de aanname, want als bijvoorbeeld een bigram er in
 362 voorkomt dan komen ook beide unigrams er in voor. Desalniettemin blijkt Naive
 363 Bayes effectief te zijn voor tekstclassificatie [6, 10]. Hiervoor zijn de functies van
 364 scikit-learn `MultinomialNB` en `BernoulliNB` gebruikt [6, 10].

365 **Beoordelen van kwaliteit** De meest gebruikte methoden om kwaliteit van
 366 politieke tekstclassificatie te beoordelen zijn *accuracy* en F_1 score, die opge-
 367 bouwd is uit *recall* en *precision*. Deze scores worden berekend op basis van vier
 368 variabelen. Deze variabelen geven weer hoeveel documenten wel of niet bij een
 369 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>

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

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

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

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

373 *Accuracy* is het percentage van documenten dat correct geclassificeerd is. *Ac-*
 374 *curacy* wordt voor de hele classificatie gedaan en niet per klasse. *Precision* is
 375 het percentage van documenten geclassificeerd als een partij, dat ook bij die
 376 partij hoort. *Recall* is het percentage documenten van documenten behorende

377 tot een partij, dat ook als die partij geclassificeerd is. F_1 is het harmonisch
378 gemiddelde van *recall* en *precision*. *Precision*, *recall* en daarmee F_1 worden per
379 partij berekend. Er zijn drie varianten om deze scores voor de hele classificatie
380 te berekenen.

381 Allereerst is er *micro*, waarbij alle variabelen bij elkaar opgeteld worden en
382 vervolgens de scores berekend. Dit leidt ertoe dat resultaten van partijen met
383 veel documenten belangrijker zijn. Als een classificatie kleine partijen groten-
384 deels fout classificeert, kan deze score alsnog hoog zijn. In het geval van meer
385 dan twee partijen is dit hetzelfde als *accuracy*.

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

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

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

398 3.2.2 DV2: Invloed van namen

399 In Diermeier et al. [4] werd aangenomen dat namen een groot effect hebben
400 op de classificatie. Hirst et al. [2] bevestigden dit voor het Europees Par-
401 lement. Aangezien hier bij deelvraag 1 niet voor was gekozen, werd bij deze
402 deelvraag gekeken hoe groot het effect hiervan is, specifiek gericht op partijna-
403 men en achternamen van Kamerleden. Op basis van vergelijkbaar onderzoek is
404 de hypothese dat de achternamen van Kamerleden en partijnamen van invloed
405 zijn.

406 Voor deze deelvraag werd wederom een classificatie gedaan met de classi-
407 ficatiemethode die resulteerde uit deelvraag 1. In deze classificatie werden alle
408 partijnamen vervangen door *PARTIJNAAM* en alle achternamen van Kamerleden
409 vervangen door *KAMERLIDNAAM*. Deze namen waren uit de Handelingen
410 gehaald. Voor partijnamen waren ook lidwoorden toegevoegd en voor achter-
411 namen van Kamerleden zijn ook verkortingen meegenomen. Dit laatste omdat
412 bijvoorbeeld *Van Haersma Buma* vaak aangesproken wordt als *Buma*. Voorna-
413 men van Kamerleden worden zelden tot nooit gebruikt, dus die waren er niet uit-
414 gehaald. Een nadeel van deze aanpak is dat ook namen van niet-Kamerleden of
415 andere woorden weggehaald kunnen worden als deze hetzelfde zijn als naam van
416 een Kamerlid. Door gebruik van gevoeligheid voor hoofdletters was geprobeerd
417 dit te voorkomen. Een opvallend voorbeeld hiervan is de naam Rutte, die zowel
418 behoort tot het Kamerlid Arno Rutte als de premier Mark Rutte. Steekproef-
419 gewijs was gekeken of er nog namen achter zijn gebleven, maar die waren niet
420 gevonden.

421 Ook werd gekeken naar classificatie met alleen partijnamen en achterna-
422 men van Kamerleden. Alle andere woorden worden weggehaald. Namen van
423 Kamerleden en partijen die niet aan elkaar geschreven worden, zoals *Partij van*
424 *de Arbeid*, zijn aan elkaar geschreven zodat het één feature is. Doordat alle

425 andere woorden weggehaald zijn, waren de bi- en trigrams combinaties van na-
426 men die zinnen uit elkaar kunnen staan, dus die niet meer informatie geven dan
427 unigrams. Daarom werd er gebruikt van de classificatiemethode uit deelvraag
428 1, maar met alleen unigrams.

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

432 3.2.3 DV3: Oppositie of regering

433 Om deze deelvraag te beantwoorden zijn drie experimenten uitgevoerd. Twee
434 daarvan zijn gebaseerd op experimenten uit Hirst et al. [2] voor dezelfde vraag.
435 De derde is ontwikkeld voor dit onderzoek. Met deze laatste wordt begonnen.
436 Bij deze deelvraag is de classificatiemethode uit deelvraag 1 zonder achternamen
437 van Kamerleden en partijnamen gebruikt. De hypothese is op basis van de
438 bevindingen van Hirst et al. dat de classificatie inderdaad afhankelijk is van
439 partij-status.

440 Als er een afhankelijkheid is van partij-status, dan is het te wachten dat het
441 aantal misclassificaties minus verwachte waarde binnen regeringspartijen en bin-
442 nen oppositiepartijen hoger ligt dan tussen een oppositiepartij en regeringspartij.
443 De verwachte waarde is afhankelijk van het aantal documenten van een partij in
444 de training set [12]. Aangezien de test set uit dezelfde set als de training werd
445 gehaald, is de verwachte waarde ook afhankelijk van het aantal documenten van
446 een partij in de test set. Uit de voorverkenning (op basis van resultaten uit
447 deelvraag 1 en 2) bleek deze correlatie tussen het aantal *false positives* van een
448 partij en het aantal documenten behorend tot die partij.

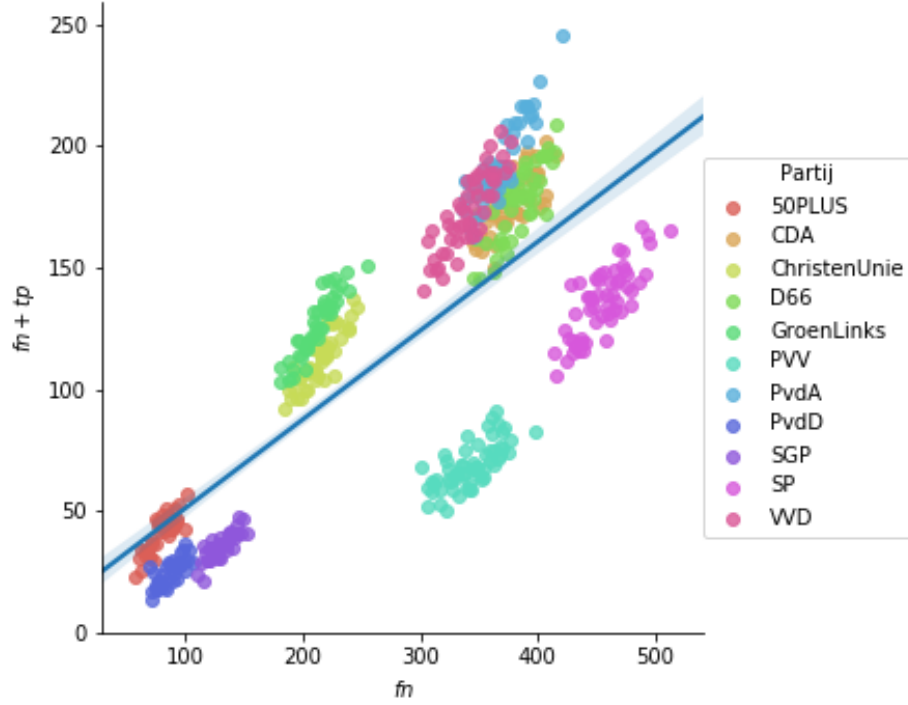


Figure 2: Het aantal *false positives* ten opzichte van het aantal documenten behorend tot die partij (*false negatives* en *true positives*). Dit is op basis van 100 classificaties met verschillende train en test set. De Pearson correlatie is 0.77 en de p-waarde 5.40×10^{-101} .

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

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

waar $i \neq j$. 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 [13]. Om te kijken of er een bias is, werden de distributies binnen regeringspartijen en binnen oppositiepartijen vergeleken met de distributie tussen beide groepen. Om de invloed van variantie door de willekeurige splitsing documenten voor trainen en testen te beperken, werd de classificatie 100 keer gedaan. Als de distributies normaal verdeeld zijn, vond de statistische

464 test plaats op basis van een eenzijdige t-toets. Als de distributies niet normaal
465 verdeeld zijn, vond dit plaats door een Mann-whitneytoets. Het gekozen signif-
466 icantieniveau (α) is 0.01. De nulhypothese is dat er geen verschil is tussen de
467 verdelingen. De alternatieve hypothese is dan dat de distributie van binnen op-
468 positie of regering groter is dan die tussen een regerings- en oppositiepartij. Op
469 basis van de bevindingen van Hirst et al. was de hypothese dat de nulhypothese
470 verworpen kan worden.

471 In het eerste experiment gebaseerd op Hirst et al. zijn de meest karakter-
472 istieke woorden per partij van de ene zittingsperiode vergeleken met de meest
473 karakteristieke woorden per partij van de andere zittingsperiode, waar een kabi-
474 net uit andere partijen bestond. De verwachting is dat als de classificatie niet het
475 gevolg is van partij-status dat de woorden bij een partij blijven en niet gekop-
476 peld zijn aan in oppositie of regering zitten. Aansluitend bij de hypothese is dus
477 de verwachting dat woorden wel wisselen van partij wanneer ze van partij-status
478 gewisseld zijn.

479 In het tweede experiment gebaseerd op Hirst et al. zijn de classificaties
480 getraind op een zittingsperiode en getest op een andere zittingsperiode, waar
481 wederom een kabinet uit andere partijen bestond. Als de classificatie afhankelijk
482 was van partij-status was de verwachting dat de scores van partijen die gewisseld
483 zijn van partij-status sterker gedaald zijn dan partijen die niet van partij-status
484 zijn veranderd. Op basis van de hypothese is dan ook de verwachting dat bij de
485 partijen die gewisseld zijn partij-status een sterkere daling te zien is.

486 Als vergelijkingsmateriaal was voor deze experimenten een tweede dataset
487 nodig uit een ander kabinet. Hiervoor is het wenselijk dat dit kabinet bestaat uit
488 andere partijen dan kabinet-Rutte II. Daarnaast is het ook wenselijk als het niet
489 te ver terug is, zodat onderwerpen en taalgebruik enigszins overeenkomstig zijn.
490 Omdat kabinet-Rutte I een minderheidskabinet was met een bijzondere partij-
491 status voor de PVV, is ervoor gekozen om de Handelingen van de Tweede Kamer
492 tijdens het missionaire kabinet-Balkenende IV (22 februari 2007 tot 20 februari
493 2010) te gebruiken. Dit kabinet bestond uit CDA, PvdA en ChristenUnie.

494 De partij 50PLUS bestond nog niet gedurende kabinet-Balkenende IV, dus
495 documenten van deze partij zijn weggelaten. Verder heeft dezelfde verwerking
496 van data plaatsgevonden, zoals beschreven in 3.1. Alleen de minimum- en max-
497 imumlengte is overgenomen van de dataset van kabinet-Rutte II.

Table 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 werd wederom formule 5 gebruikt als verwachte waarde en formule 6 als error. De hypothese is dat de classificatie deels afhankelijk is van positie op de links-rechts as.

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 de links-rechts as, op basis van het verkiezingsprogramma van dat jaar. Voor de dataset van kabinet-Rutte II is gebruikt gemaakt van de scores op basis van de verkiezingsprogramma's voor de verkiezingen van 2012.

Table 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

512 of er een correlatie is tussen de error van twee partijen en de afstand op de
513 links-rechts as van die partij. Het significantieniveau (α) hiervoor is opnieuw
514 0.01. De nulhypothese is dat er geen negatieve correlatie is tussen de error en
515 de afstand op de links-rechts as. De alternatieve hypothese is dat er wel een
516 negatieve correlatie is tussen de error en de afstand op de links-rechts as.

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

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

523 3.2.5 DV5: Woordgebruik van sprekers

524 De vorige classificaties trainden op documenten en werden getest op andere
525 documenten, maar wel van dezelfde sprekers als uit de training set. Naast
526 de ideologie kan de classificatie daarom ook getraind zijn op het taalgebruik
527 van sprekers. Als een Kamerlid bijvoorbeeld een woord regelmatig in speeches
528 gebruikt, maar niet wordt gebruikt door zijn partijgenoten, wordt dit wel gezien
529 als een belangrijk woord voor de classificatie naar partij-affiliatie. Hirst et al.
530 [2] plaatsten al een soortgelijke kanttekening bij de resultaten van Diermeier et
531 al. [4]. De hypothese is dat de classificatie afhankelijk is van woordgebruik van
532 sprekers

533 Om te kijken of dit effect er is, wordt er opnieuw een classificatie gedaan
534 met de classificatiemethode uit deelvraag 1 zonder achternamen van Kamerleden
535 en partijnamen. Ditmaal worden alleen niet de individuele documenten verdeeld
536 over de training en test set, maar worden de Kamerleden, met bijbehorende
537 documenten, verdeeld over de training en test set. Als taalgebruik van een
538 spreker in de training set voorheen invloed had op de classificatie, zal dat nu geen
539 effect meer hebben omdat er geen documenten van die spreker meer voorkomen
540 in de test set. De verwachting is daarom ook dat deze classificatie lagere scores
541 vindt dan die van deelvraag 2.

542 4 Resultaten

543 4.1 DV1: Beste classificatiemethode

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

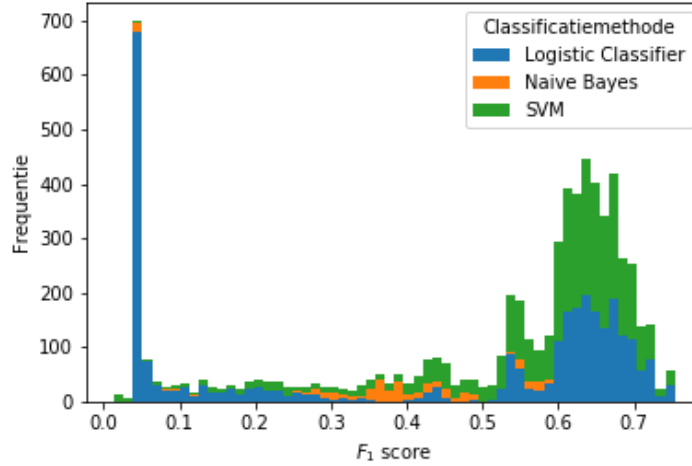


Figure 3: Histogram van de grid search met de F_1 scores van de classificatiemethoden

547 Het beste resultaat werd bereikt met Support Vector Machines gebruik-
 548 makend van *stochastic gradient descent learning* en L2 regularisatie. In de grid
 549 search behaalde deze methode een F_1 score en *accuracy* van 0.75. Voor beide
 550 scores was dit het hoogste van de grid search. De woorden waren hierbij gestemd.
 551 De features waren zowel unigrams, bigrams als trigrams. Geen features zijn
 552 weggelaten door minimale of maximale documentfrequenties. De waarden van
 553 deze features waren *tf-idf* scores. Het maximum aantal iteraties was 5 voor de
 554 grid search, maar de rest van resultaten zijn op basis van 100 iteraties.

555 Tabel 4 laat de scores zien per partij met het aantal documenten in de
 556 test set. De *accuracy* voor deze classificatie is 0.80. De F_1 scores per partij
 557 liggen tussen de 0.7 en 0.9. De partijen met een sterke focus op één onderwerp,
 558 50PLUS, PVV en PvdD, als ook de SGP hebben hoge scores. De coalitiepartijen,
 559 VVD en PvdA, daarentegen hebben lagere scores. Figuur 4 laat zien waar de
 560 fouten in deze classificatie zitten. De meest karakteristieke n-grams per partij
 561 zijn te zien in tabel 5. Met meest karakteristiek worden de n-grams bedoeld die
 562 de hoogste coëfficiënt hebben in de classificatie en die dus relatief het meeste
 563 belangrijk zijn voor de classificatie van een partij. Hierin is te zien dat vrijwel
 564 alle n-grams achternamen van Kamerleden of partijnamen bevatten.

Table 4: Classificatie scores per partij van de beste classificatiemethode (SVM). Gemiddelde van vijf splitsingen van training en test set.

	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

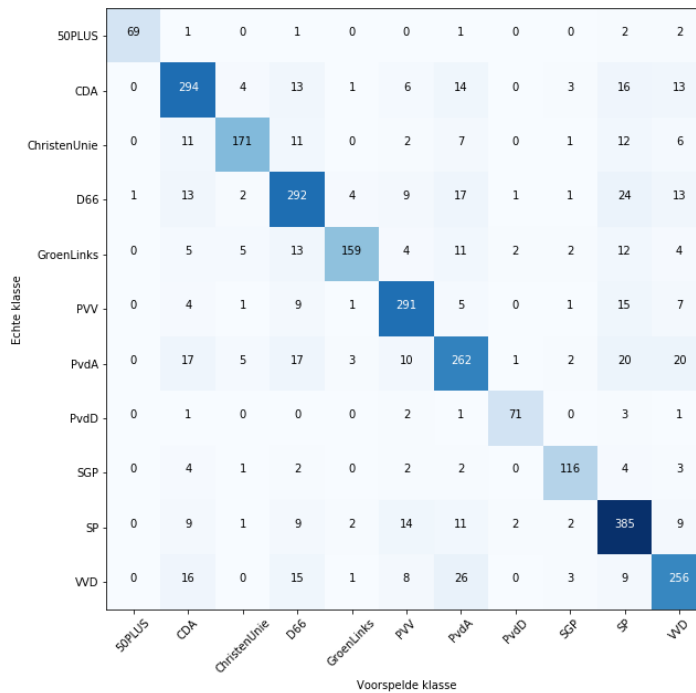


Figure 4: Confusion matrix van de beste classificatiemethode (SVM). Gemiddelde van vijf splitsingen van training en test set.

Table 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

Table 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

565 4.2 DV2: Invloed van namen

566 In tabel 5 was al te zien dat de meest karakteristieke n-grams voornamelijk
567 achternamen van Kamerleden of partijnamen bevatten. In tabel 6 zijn de scores
568 te zien voor een classificatie met alleen achternamen van Kamerleden en par-
569 tijnamen. De *accuracy* is 0.61. De scores zijn gedaald ten opzichte van de
570 resultaten van deelvraag 1, maar hoger dan de baseline scores.

Table 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

571 In tabel 7 zijn de F_1 scores te zien van classificatie met achternamen van
572 Kamerleden en partijnamen vervangen. De *accuracy* hiervan is 0.58. De scores
573 zijn lager dan die uit deelvraag 1 en lager dan van de classificatie met alleen
574 namen. Wel zijn de scores nog steeds hoger dan de baseline. In tabel 8 is
575 vervolgens te zien welke n-grams het meest karakteristiek zijn per partij voor
576 deze classificatie.

Table 7: 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 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

Table 8: 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

Table 8: 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	bewindsliden	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

577 4.3 DV3: Oppositie of regering

578 In figuur 5 zijn de distributies van de errors, zoals gedefinieerd in formule 6 te
579 zien van combinaties van regerings- en oppositiepartijen. Bijgevoegd zijn het
580 aantal combinaties (N), de mediaan en de interkwartielafstand (IKA).

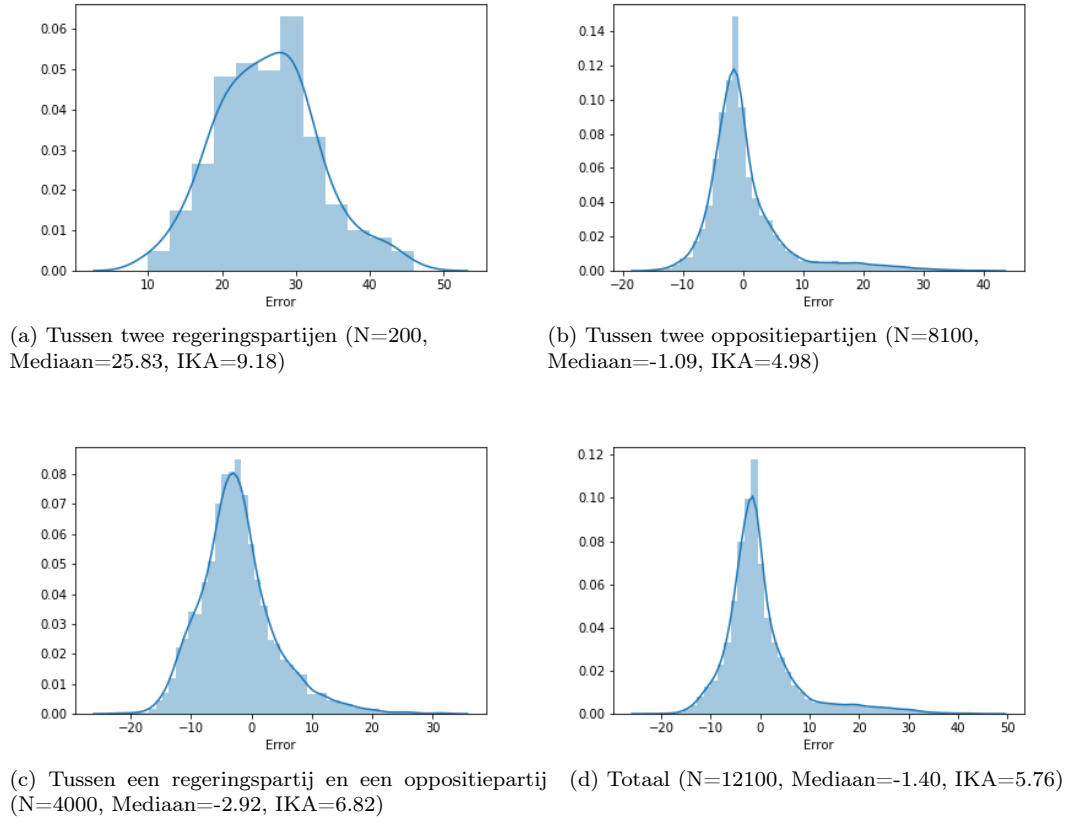


Figure 5: Genormaliseerde distributie van de error uit formule 6 voor de verschillende combinaties.

581 Voor alle distributies was de nulhypothese verworpen worden dat deze
582 normaal verdeeld zijn ($p < 0.01$). In tabel 9 is vervolgens te zien dat er een sig-
583 nificant verschil is tussen de distributies binnen regering en oppositie tegenover
584 de distributie tussen regering en oppositiepartij. Tussen regeringspartijen zijn er
585 gemiddeld 26.11 misclassificaties meer dan verwacht en tussen oppositiepartijen
586 gemiddeld 0.43.

Table 9: 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

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

Table 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	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

Table 10: Meest karakteristieke n-grams per partij op basis van classificatie uit deelvraag 2 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

590 In tabel 11 zijn de resultaten van de classificatiescores te zien waarbij de
 591 classificatie getraind is op een zittingsperiode, maar getest op een andere. De
 592 resultaten zijn sterk gedaald, maar nog boven de baseline. De daling verschilt
 593 enorm per partij en zittingsperiode met dalingen van F_1 scores tussen 12 en
 594 92%.

Table 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

595 4.4 DV4: Links-rechts as

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

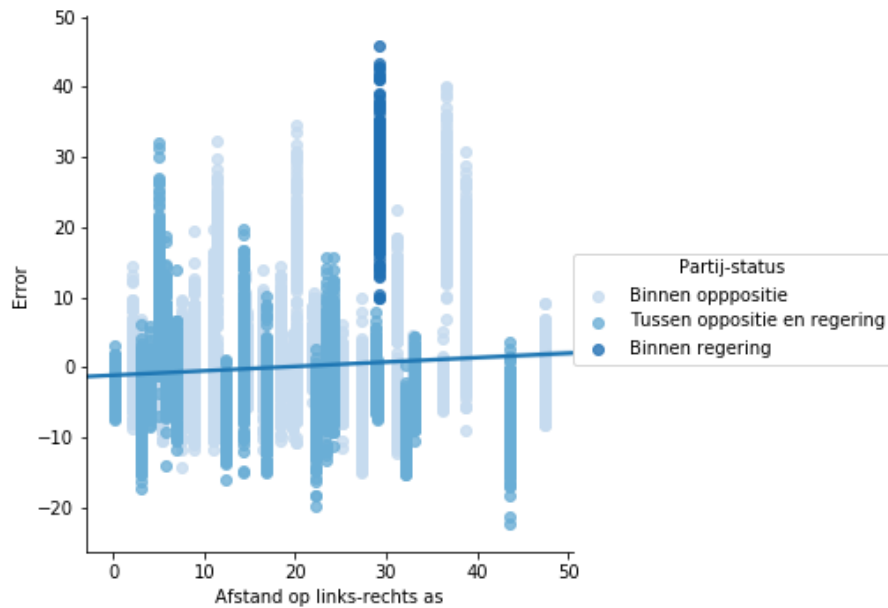


Figure 6: Error ten opzichte van de afstand op de links-rechts as van twee partijen. Gebaseerd op 100 classificaties met verschillende test en train set. De Pearson correlatie is 0.09 en de p -waarde 2.39×10^{-20} .

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. 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 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.

Table 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.

Table 13: Classificatierapport van beste classificatie met de Kamerleden verdeeld over training en test set. Gemiddelde van tien 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 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 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.

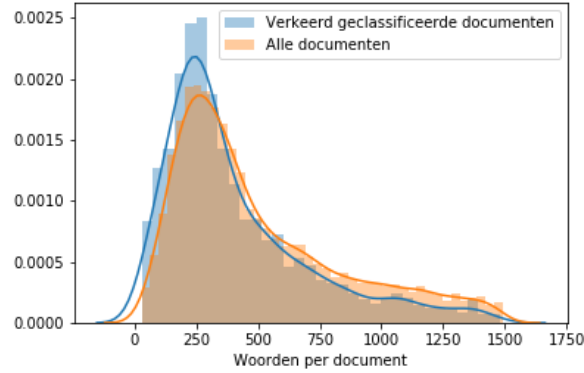


Figure 7: Genormaliseerde distributie van documentlengtes van foutief geclas-
sificeerde documenten en alle documenten. Totaal van 5-fold cross-validation,
waardoor documenten vaker voor kunnen komen. Mediaan documentlengte van
foutief geclassificeerde documenten is 321 en voor alle documenten 386.

632 Voor een vervolgonderzoek kan uitgebreider gekeken worden naar dit effect
633 en wat dit betekent voor de resultaten. Het percentage documenten van een
634 vragenuur is tweemaal zo hoog bij foutief geclassificeerde documenten, maar dit
635 lijkt te komen doordat deze documenten vaak kleiner zijn (mediaan is 286).
636 Er is verder nog gekeken naar andere verbanden tussen documenten die
637 verkeerd zijn geclassificeerd. Daarbij is nog te zien dat sprekers met weinig doc-
638 umenten relatief iets meer voorkomen in verkeerd geclassificeerde documenten.

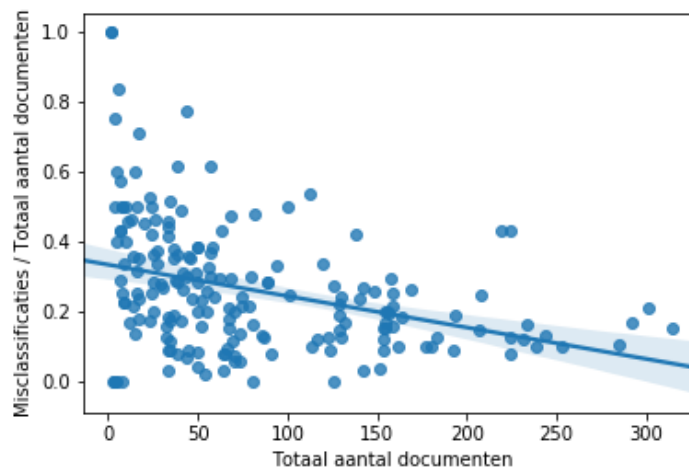


Figure 8: Aantal misclassificaties gedeeld door totaal aantal documenten per
spreker tegenover totaal aantal documenten van een spreker. Misclassificaties
zijn totaal van 5-fold cross-validation, waardoor documenten vaker mee kunnen
tellen. De pearson correlatie is -0.28 en de p -waarde 1.07×10^{-4} .

639 Dit versterkt het vermoeden dat de classificatie mede plaatsvindt op basis
640 van woordgebruik van individuele sprekers, zoals besproken wordt in deelvraag
641 5.

642 5.2 DV2: Invloed van namen

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

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

657 De classificatiemethode die gebruikt is in deze deelvraag, is gebaseerd op
658 de beste methode voor de dataset uit deelvraag 1. Hierin was gevonden dat een
659 combinatie van uni-, bi- en trigrams het beste resultaat opleverde. In tabel 5
660 is te zien dat trigrams behoren tot de meest karakteristieke n-grams, hoewel de
661 woorden in trigrams vaak overlappen met uni- en bigrams. In tabel 8 daarente-
662 gen zijn er nog maar een paar trigrams, welke grotendeels procedurele zinnen
663 zijn of toevoeging van een lidwoord op een uni- of bigram. Dit verschil sug-
664 gereert dat trigrams minder belangrijk zijn in de classificatie zonder de namen,
665 dus de classificatiemethode uit deelvraag 1 niet het beste is voor deze classifi-
666 catie. In vervolgonderzoek kan de opzet van deelvraag 1 toegepast worden op
667 de classificatie zonder de namen, om zo te komen tot een classificatiemethode
668 die het beste resultaat oplevert op de classificatie zonder namen.

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

676 5.3 DV3: Oppositie of regering

677 In tabel 4 is het opvallend dat de coalitiepartijen lage scores krijgen. Daarnaast
678 laat figuur 4 zien dat er een hoge overlap zit tussen deze twee partijen. De
679 resultaten van het eerste experiment, ontwikkeld voor dit onderzoek, vinden
680 een afhankelijkheid van partij-status. De twee andere experimenten versterken
681 deze bevindingen niet, zoals wel het geval was bij Hirst et al. [2]. Hieronder
682 worden de resultaten nader besproken.

683 De statistische toetsresultaten in tabel 9 laten zien dat inderdaad de error
684 groter is binnen oppositie of regering dan tussen een regerings- en oppositiepar-

685 tij. Met name regeringspartijen lijken lastiger uit elkaar te halen. Dit suggereert
 686 dat inderdaad partij-status invloed heeft op de classificatie.

687 De verwachting was dat de error normaal verdeeld zou zijn. De verdelin-
 688 gen uit figuur 5 hebben globaal wel de vorm van een normaal verdeling. In
 689 figuur 2 is het daarnaast opvallend dat partijen zoals SP en PVV ruim onder de
 690 regressielijn zitten, terwijl andere partijen er een stuk boven zitten. Dit geeft
 691 aanleiding te vermoeden dat er naast het aantal documenten van een partij
 692 nog meer factoren van invloed zijn op het aantal misclassificaties en daarmee
 693 de verwachte waarde. En deze verwachte waarde en de daar uit volgende er-
 694 ror zijn een belangrijke aanname van deze methode. Voor deze methode is het
 695 dus belangrijk uit te vinden of dit een goede benadering is van de verwachte
 696 waarde. In deelvraag 4 wordt gekeken of links/rechts positie hier nog invloed
 697 heeft. Voor een vervolgonderzoek kan nog verder gekeken worden naar invloeden
 698 op verwachte waarde of andere confounding biases.

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

Table 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>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>
			<i>voor</i>
	ChristenUnie		<i>voor PARTIJ</i>
			<i>gaan</i>
		<i>mijn</i>	<i>termijn</i>
		<i>waarop</i>	<i>blij met de</i>
		<i>blij</i>	<i>volgens</i>
PvdA	<i>collega KAMERLID</i>	<i>volgens mij</i>	
	<i>erg</i>	<i>blij</i>	
		<i>beantwoording</i>	
		<i>volgens</i>	
		<i>volgens mij</i>	

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

709 zeggen over een regeringspartij.

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

720 Dat de experimenten uit Hirst et al. in hun onderzoek wel invloed vin-
721 den, maar in dit onderzoek niet kan komen doordat hun onderzoek zich richt
722 op binaire classificatie, terwijl dit onderzoek meerdere partijen heeft. Zo kan
723 het ontbreken van gemeenschappelijke n-grams komen doordat regeringspartijen
724 zich ook van elkaar moeten onderscheiden in dit onderzoek, waarvoor n-grams
725 die relevant zijn voor partij-status weinig effect hebben, terwijl in het onderzoek
726 van Hirst et al. de regeringspartij alleen onderscheiden hoeft te worden van de
727 oppositiepartij. Daarnaast verklaarden zij dat de daling tussen twee zittingspe-
728 riodes het gevolg was van de wisseling van partij-status. In dit onderzoek kon
729 daarentegen gekeken worden naar effecten op partijen niet die niet van partij-
730 status zijn gewisseld. Hierin was te zien dat de daling ook aanwezig was bij
731 partijen die niet gewisseld zijn van partij-status.

732 5.4 DV4: Links-rechts as

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

739 Alle correlaties zijn statistisch significant, maar de Pearson correlatie en
740 daarmee effectgrootte is klein. Daarnaast is het ook opvallend dat de twee
741 combinaties van partij-statussen een andere correlatierichting hebben. Dit sug-
742 gereert dat de statistische significantie het gevolg is van de grote steekproef en
743 maar een klein effect [15].

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

750 5.5 DV5: Woordgebruik van sprekers

751 De resultaten uit tabel 13 zijn laag, amper hoger dan de baseline. Dit sug-
752 gereert inderdaad dat eerdere classificaties in grote mate toch afhankelijk waren
753 van het woordgebruik van sprekers. Dit is opmerkelijk aangezien vergelijkbare

754 onderzoeken dit effect niet vinden. De meest karakteristieke n-grams van deze
755 classificatie wijken daarnaast grotendeels niet af van die uit tabel 8.

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

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

774 5.6 Algemeen

775 Het vergelijken van deze resultaten met vergelijkbaar werk is ingewikkeld, aangezien
776 de keuzes en eigenschappen van die onderzoeken het niet een één-op-één vergeli-
777 jking maken. Voorbeelden hiervan zijn de taal, het parlement, de documentg-
778 rootte, baselines, behouden of weglaten van namen, een spreker als document
779 zien en het trainen en testen op dezelfde spreker. Hoewel de resultaten in som-
780 mige gevallen lager zijn dan die uit vergelijkbaar werk, is het belangrijk hier
781 rekening mee te houden. Een vervolgonderzoek zou daarom dit onderzoek kun-
782 nen reproduceren op een ander parlement om daarmee te kunnen vergelijken.

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

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

796 6 Conclusies

797 Dit onderzoek vindt een *accuracy* en F_1 score van 0.80 voor het classificeren
798 van spreekbeurten in de Tweede Kamer naar partij-affiliatie. De beste classifi-

799 catiemethode maakt gebruik van Support-Vector Machines. De baseline scores
800 zijn respectievelijk 0.11 en 0.15. Als rekening wordt gehouden met partijnamen
801 en achternamen Kamerleden daalt de *accuracy* naar 0.58 en de F_1 score naar
802 0.57. Dit onderzoek vindt aanwijzingen dat deze classificatie afhankelijk is van
803 de partij-status (oppositie of regering). Daarnaast vindt dit onderzoek geen
804 aanwijzingen dat de classificatie afhankelijk is van positie op de links-rechts as.
805 Als rekening wordt gehouden met woordgebruik van individuele Kamerleden,
806 dalen de *accuracy* en F_1 verder naar 0.27. Daarmee lijkt de classificatie naar
807 partij-affiliatie in grote mate niet het gevolg van ideologie. Deze conclusie trekt
808 daarmee de bruikbaarheid van tekstclassificatie voor het vinden van een relatie
809 tussen woordgebruik en ideologie in twijfel. Op een aantal punten wijken de
810 bevindingen van dit onderzoek af van vergelijkbare onderzoeken. Voor een ver-
811 volgonderzoek kan dit onderzoek uitgebreid worden met een aantal aanbevelin-
812 gen.

813 References

- 814 [1] Felix Bießmann. Automating political bias prediction. *CoRR*,
815 abs/1608.02195, 2016.
- 816 [2] Graeme Hirst, Yaroslav Riabinin, Jory Graham, and Magali Boizot-roche.
817 Text to ideology or text to party status? In Bertie Kaal, Isa Maks, and An-
818 nemarie van Elfrinkhof, editors, *From Text to Political Positions*, chapter 5,
819 pages 93–115. John Benjamins Publishing Company, Amsterdam, 2014.
- 820 [3] Vasco Ferreira. Using textual transcripts of parliamentary interventions for
821 profiling portuguese politicians. 2016.
- 822 [4] Daniel Diermeier, Jean-François Godbout, Bei Yu, and Stefan Kaufmann.
823 Language and ideology in congress. *British Journal of Political Science*,
824 42(1):31–55, 2012.
- 825 [5] Beia Yu, Stefan Kaufmann, and Daniel Diermeier. Classifying party affili-
826 ation from political speech. *Journal of Information Technology & Politics*,
827 5(1):33–48, 2008.
- 828 [6] Maneesh Bhand, Dan Robinson, and Conal Sathi. Text classifiers for po-
829 litical ideologies, 2009.
- 830 [7] Bjørn Høyland, Jean-François Godbout, Emanuele Lapponi, and Erik Vell-
831 dal. Predicting party affiliations from european parliament debates. In
832 *Proceedings of the ACL 2014 Workshop on Language Technologies and*
833 *Computational Social Science*, pages 56–60. Association for Computational
834 Linguistics, 2014.
- 835 [8] Laura Klompenhouwer. Extra ledenvergadering 50plus om splitsing. *NRC*
836 *Handelsblad*, June 2014.
- 837 [9] Eric Jones, Travis Oliphant, Pearu Peterson, et al. SciPy: Open source
838 scientific tools for Python, 2001.

- 839 [10] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel,
840 M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Pas-
841 sos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-
842 learn: Machine learning in Python. *Journal of Machine Learning Research*,
843 12:2825–2830, 2011.
- 844 [11] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. *Intro-*
845 *duction to Information Retrieval*. Cambridge University Press, New York,
846 NY, USA, 2008.
- 847 [12] Mahendra Sahare and Hitesh Gupta. A review of multi-class classifica-
848 tion for imbalanced data. *International Journal of Advanced Computer*
849 *Research*, 2(3), 2012.
- 850 [13] *NIST/SEMATECH e-Handbook of Statistical Methods*. NIST/SEMAT-
851 ECH, April 2012.
- 852 [14] Andrea Volkens, Pola Lehmann, Theres Matthieß, Nicolas Merz, Sven
853 Regel, and Bernhard Weßels. The manifesto data collection. manifesto
854 project (mrg/cmp/marpor). version 2017b, 2017.
- 855 [15] Joseph F. Hair, Jr., Rolph E. Anderson, Ronald L. Tatham, and William C.
856 Black. *Multivariate Data Analysis (6th Ed.)*. Prentice-Hall, Inc., Upper
857 Saddle River, NJ, USA, 2006.