

Measuring Rhetorical Similarity with Machine Learning

Nicolai Berk

University of Amsterdam

July 9, 2020

Relevance

Group similarity of interest for many areas in political science:

- Representation
 - Do certain (e.g. female/working class/PoC) MPs communicate differently?
 - How similar is legislation to interest group's demands?
 - Do MP's become more similar over time?
 - Are directly elected MPs communicating differently?
- Government formation
 - Which parties communicate more similar and are more likely to form a coalition?
- Party Competition
 - Did parties become more or less similar to each other?
 - Did parties fill political opportunity spaces opened by other parties?
 - Which MPs are likely to leave their party?
- Polarisation, ...

Current Methods

- Simple similarity measures: cosine, jaccard
→ no weighting of terms
 - Scaling methods: WORDSCORE, WORDFISH (Laver, Benoit, and Garry 2003; Slapin and Proksch 2008)
→ no similarity measures, cannot measure similarity of a document to a collection of texts
 - Peterson and Spirling 2018: Machine learning accuracy to measure polarisation → no estimates for single texts
- **Use ML to scale texts on 'groupness'**

A supervised ML measure of group similarity

- ① Pre-process texts
- ② Select well-performing classifier
- ③ Balance data, using SMOTE (Chawla et al. 2002)
- ④ BoW-classifier trained on full set
- ⑤ Predicted probabilities assigned to all texts
- ⑥ Train one classifier/time-period to 'control' for changing language

Application

- How similar are speakers and parties to the radical right?
- Parliamentary speeches from lower Chambers in AT, DE, NL (Rauh and Schwalbach 2020)
- estimates of similarity to radical right party obtained using Python's `scikit-learn`
- Logistic regression on unstemmed text, tfidf-weighted

Validation

Three approaches

- Content of RR language
- Government formations with RR participation (AT & NL)
- Speaker estimates of party exits (DE & NL)

Best predictor Words

+0.796	SPÖ	+5.317	Bürger	+10.354	immigratie
+0.688	geehrten	+5.209	Merkel	+9.579	enzovoorts
+0.677	einmal	+5.037	deutschen	+9.199	islam
+0.616	eben	+4.758	AfD	+9.081	PVV
+0.610	Patienten	+4.564	Deutschland	+8.541	islamitische
+0.604	Freiheitliche	+4.491	Altparteien	+7.271	Brussel
+0.590	Hohes	+4.149	hier	+7.216	Mogelijk
+0.550	Kollegin	+4.066	Regierung	+6.936	immigratiebeleid
+0.535	Jahren	+3.879	Vielen	+6.902	illegalen
+0.533	Bevölkerung	+3.843	bedanke	+6.882	137d
+0.477	Auch	+3.575	Damen	+6.562	tekenen
+0.460	Bundesminister	+3.501	Danke	+6.523	dierenpolitie
+0.458	Rednerpult	+3.468	Herrn	+6.414	huishoudelijke
+0.448	Zur	+3.447	Arbeitnehmer	+6.357	gigantisch
+0.447	Wien	... 11658 more positive ...		+6.012	kortingen
+0.443	geehrte	... 15143 more negative ...		+5.965	bene
+0.443	Sozialdemokratie	-3.385	Soldatinnen	+5.863	statushouders
... 6818 more positive ...		-3.425	Staatssekretär	+5.856	hemelsnaam
... 7981 more negative ...		-3.448	Herzlichen	... 13823 more positive ...	
-0.445	Kurz	-3.463	wichtig	... 25807 more negative ...	
-0.448	halte	-3.792	Koalitionsvertrag	-5.966	Kops
-0.456	Finanzminister	-3.893	letzten	-6.149	afweging
-0.517	Entschließungsantrag	-3.943	Deswegen	-6.164	Tegelijkertijd
-0.529	Vielen	-4.160	Linke	-6.209	afspraken
-0.540	FPÖ	-4.327	Kollegen	-6.311	zorgvuldig
-0.552	Bürgerinnen	-4.486	Deshalb	-6.387	ontraad
-0.589	Menschen	-4.504	brauchen	-6.962	heer
-0.619	Liebe	-4.567	Bürgerinnen	-6.977	kijken
-0.624	Ministerin	-4.604	Demokraten	-7.071	Agema
-0.727	Bundesministerin	-4.618	sagen	-7.877	helder
-0.758	Kolleginnen	-4.884	finde	-7.888	overwegende
-0.841	ÖVP	-9.485	Kolleginnen	-8.770	Baudet

(a) Austria

(b) Germany

(c) Netherlands

Best predictor Words I

What distinguishes the radical right in **Austria**?

- talking about **opposition**
- talking about Vienna (SD-governed)
- talking about health reform (RR-led ministry)
- less likely to address house as colleagues
- less likely to use gender-inclusive language

Best predictor Words II

What distinguishes the radical right in **Germany**?

- talking about **government**
- talking about Germany
- less likely to address house as colleagues
- less likely to use gender-inclusive language
- less causal language

Best predictor Words III

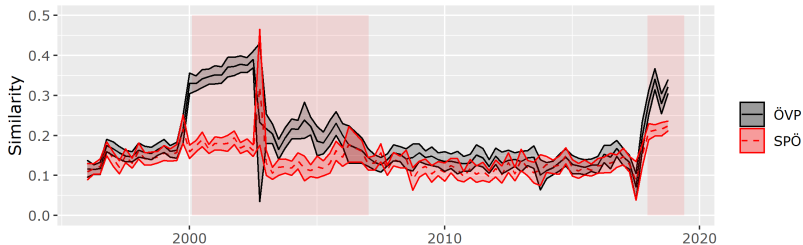
What distinguishes the radical right in **the Netherlands**?

- strong immigration focus
- more informal language
- less nuancing language

Radical-right government participation I

Quarterly similarity estimates of radical-right FPÖ to...

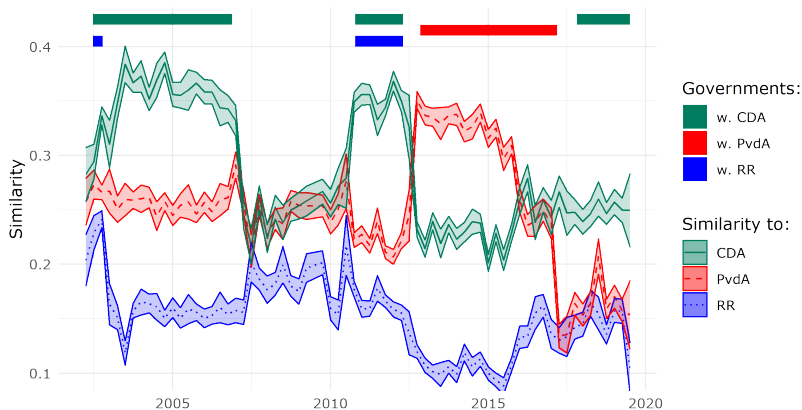
Red areas indicate ÖVP-FPÖ coalition governments



Radical-right government participation II

Quarterly similarity estimates of liberal VVD

Top bars indicate governing status



Speaker distribution I

Leaving the AfD

- One day after the national election in 2017, Frauke Petry and Mario Mieruch left the party
- They became independent members of the Bundestag

Leaving the VVD

- Geert Wilders left the VVD in September 2004
- this followed a rift with his party over the immigration issue
- as well as a series of escalations (position paper in July, vote on Turkey end of August)

Speaker distribution II

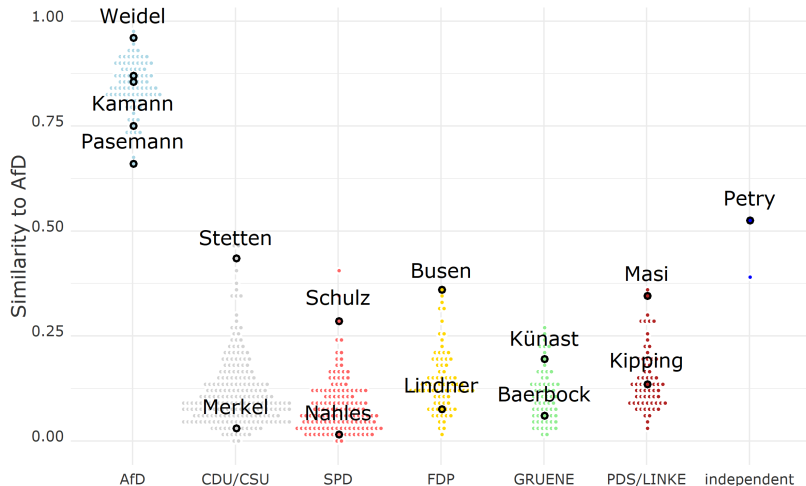
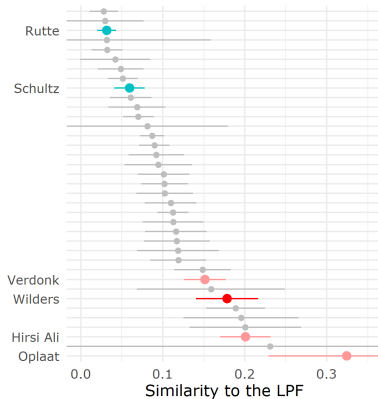
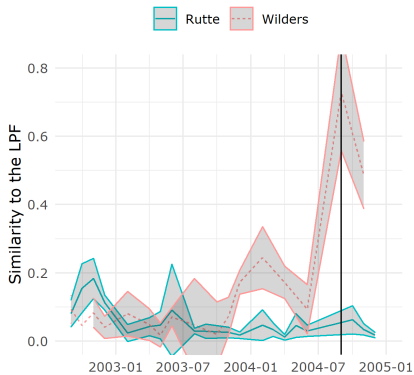


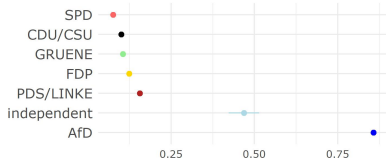
Figure: Speakers in the current German Bundestag.

Speaker distribution III

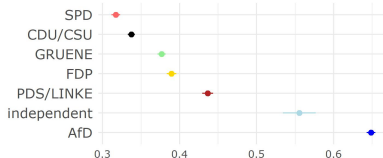


Comparable Measures

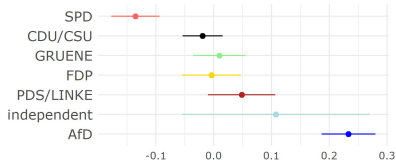
Supervised estimate



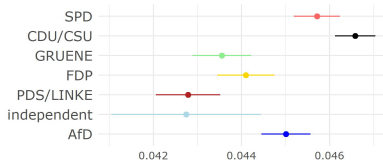
Supervised estimate (subset)



Wordfish



Cosine similarity



Summary

- SML estimate gives precise measure of rhetorical similarity
- the meaning of these differences can be assessed in detail
- outperforms other methods like WORDFISH and cosine similarity

Reviewer comments

- Scaling speeches on 'partyness' → not a novel method
- Doesn't build on Peterson and Spirling 2018 (scaling, not aggregate accuracy)
- Cosine similarity not a measure of similarity (document to corpus)

Possible routes forward

Substantive piece

- Testing classical theories of RR accommodation
- Improving predictors of coalition formation
- Historical: Showing accommodation of NSDAP by conservatives (Levitsky & Ziblatt)

Methods contribution

- Arguing for extended use of SML for better content validity
- Showing a possible application

Thank you!

Questions?

Resources I



Chawla, Nitesh V. et al. (2002). “SMOTE: Synthetic Minority Over-sampling Technique”. In: *Journal of Artificial Intelligence Research* 16, pp. 321–357. ISSN: 10769757. DOI: 10.1613/jair.953. arXiv: 1106.1813. URL: <https://arxiv.org/pdf/1106.1813.pdf%7B%5C%%7D0Ahttp://www.snopes.com/horrors/insects/telamonia.asp>.



Laver, Michael, Kenneth Benoit, and John Garry (2003). “Extracting policy positions from political texts using words as data”. In: *American Political Science Review* 97.2, pp. 311–331. ISSN: 00030554. DOI: 10.1017/S0003055403000698.



Peterson, Andrew and Arthur Spirling (2018). “Classification Accuracy as a Substantive Quantity of Interest: Measuring Polarization in Westminster Systems”. In: *Political Analysis* 26.1, pp. 120–128. ISSN: 14764989. DOI: 10.1017/pan.2017.39.

Resources II



Rauh, Christian and Jan Schwalbach (2020). “The ParlSpeech V2 data set: Full-text corpora of 6.3 million parliamentary speeches in the key legislative chambers of nine representative democracies”. In: *Harvard Dataverse, V1*, pp. 1–14. DOI: <https://doi.org/10.7910/DVN/L40AKN>, . URL: <https://dataverse.harvard.edu/dataverse/ParlSpeech>.



Slapin, Jonathan B. and Sven Oliver Proksch (2008). “A scaling model for estimating time-series party positions from Texts”. In: *American Journal of Political Science* 52.3, pp. 705–722. ISSN: 00925853. DOI: 10.1111/j.1540-5907.2008.00338.x.