

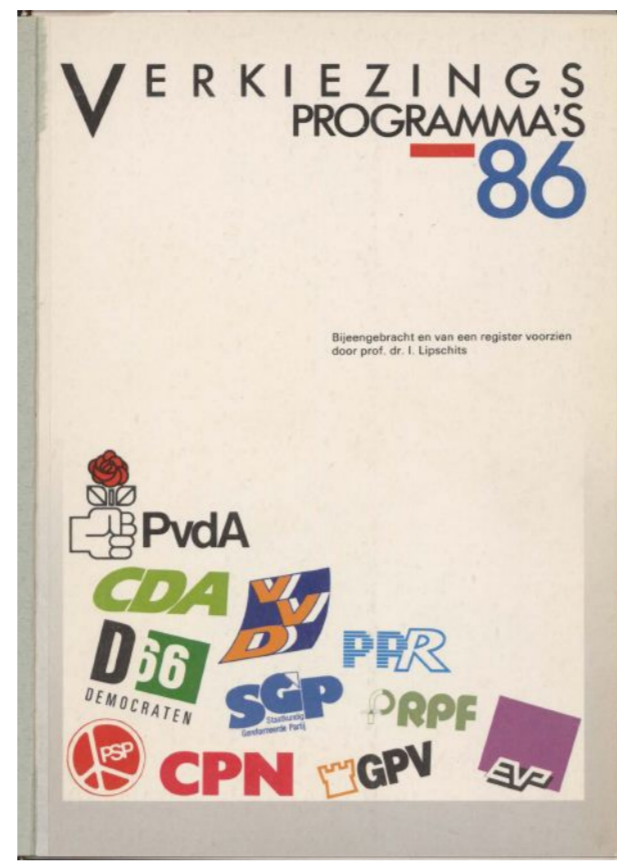
Automatic thematic classification of election manifestos

Suzan Verberne¹, Eva D'hondt¹, Antal van den Bosch¹, Maarten Marx²

1. Radboud University Nijmegen
2. University of Amsterdam

Introduction

- ▶ Isaac Lipschits, Dutch political scientist, annotated party programmes for the Dutch elections (1977-1998) with themes.
- ▶ Goals in the Political Mashup project:
 - ▷ Digitize the 1977-1998 Lipschits collections
 - ▷ Build an automatic classifier for more recent, unclassified editions

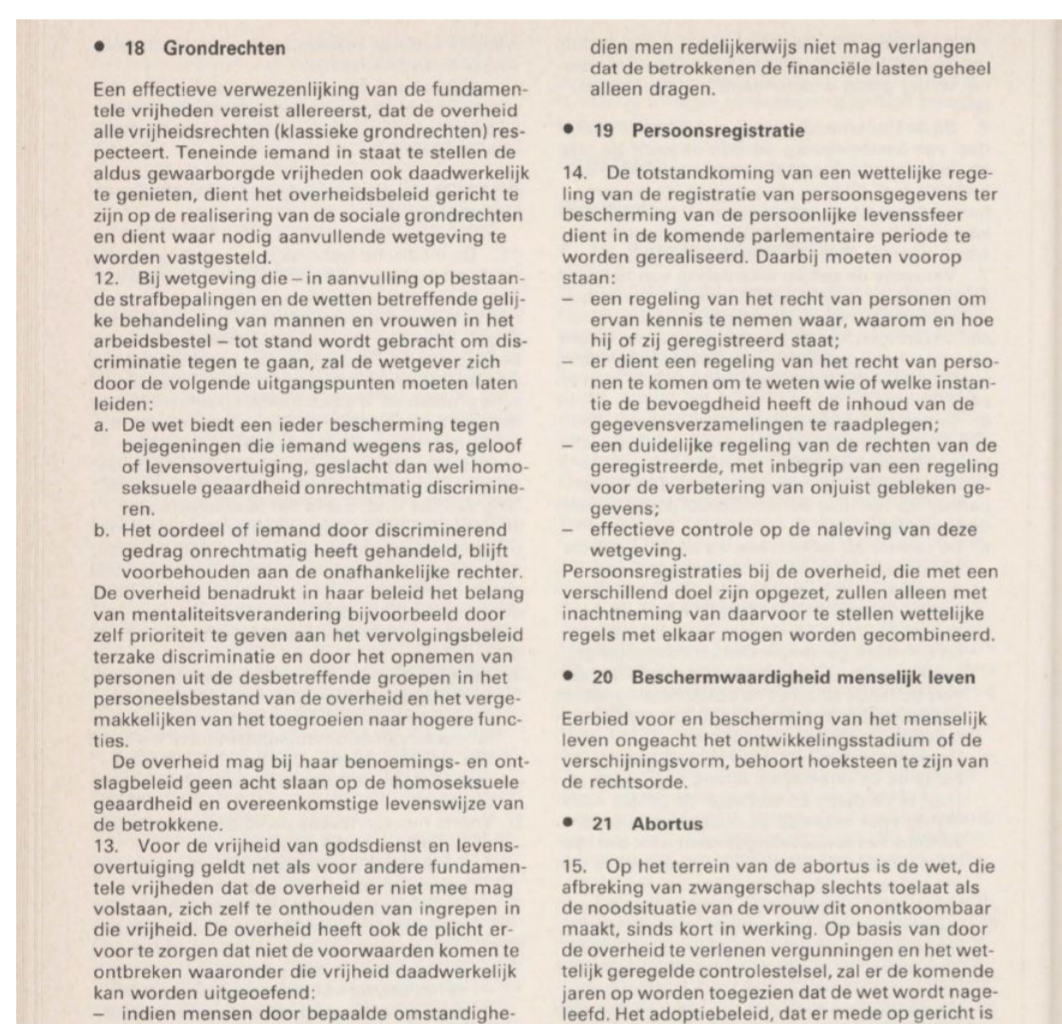


Approach

1. Convert the scanned PDFs of 1986, 1994 and 1998 to enriched publications
2. Use these data to train a classifier
3. Run the classifier on 2006-2012 data
4. Ask an expert to evaluate a sample of the labels for the 2006-2012 data

Converting the PDF books to enriched publications

- ▶ The body of Lipschits data are the manifesto texts, labelled with numbers at the places in the text where the topic changes



- ▶ At the back of the book is a register of alphabetically sorted themes with the associated texts per party by topic number.
- ▶ Lipschits classification is a multi-label classification task

theme	parties
18 Grondrechten	CDA 14, 19, 23, 24, 27, 29, 31, 32, 46, 73; CPN 1-3, 7, 9, 28, 30, 37, 40, 43, 44, 46, 48, 49; D'66 3, 14, 18, 24, 26, 33, 35, 36; EVP 20, 40, 59, 60, 64, 65; GPV 9, 20, 25, 55, 111, 114, 120; PPR 6, 16, 17, 19, 21, 23, 26, 40, 41; PSP 5-10, 18-23, 42-44, 46; PvdA 3, 5, 6, 30, 35, 39, 42; RPF 3, 5, 10, 24-26, 28, 29, 38, 39, 55, 56, 65, 75, 104, 112; SGP 1, 2, 9, 15, 16, 18, 26, 28, 31, 62, 65; VVD 5, 7, 8, 10, 13, 15, 43, 49, 84
19 Persoonsregistratie	CDA 53, 82; CPN 34; D'66 24; EVP 54; PSP 38, 39; PvdA 35; RPF 115; SGP 37; VVD 21
20 Beschermwaardigheid menselijk leven	CDA 2, 15, 36, 41, 45, 75, 77; CPN 9, 17, 27, 39; D'66 9, 11, 15, 40; EVP 62; GPV 21, 22, 49, 135; PPR 5, 6, 46; RPF 3, 29, 40, 41; SGP 13, 64, 76; VVD 2, 3, 9, 10, 17
21 Abortus	CDA 53, 82; CPN 34; D'66 24; EVP 54; PSP 38, 39; PvdA 35; RPF 115; SGP 37; VVD 21

Table : Statistics on the classification data

year	# words	vocabulary size	# texts	# themes	# words/text	# themes/text
1986	290,942	32,836	797	214	365	10.5
1994	269,270	32,499	951	210	283	6.9
1998	244,697	31,162	826	218	296	8.3

Example of data visualization

Dispersion matrix of the 1998 manifesto of the 'pensioners party'

#THEME	DISPERSION VALUE	DISPERSION MATRIX
ouderenbeleid	0.409	
sociale_voorzieningen	0.344	
gezondheidszorg	0.276	
minimuminkomen	0.276	
ADW_AOWers	0.272	
werkgelegenheid	0.272	
loon_en_inkomensbeleid	0.258	
ziekenverzorging	0.258	
economische_groei	0.254	
armoedebestrijding	0.224	
gehandicapten	0.224	
belastingen_algemeen	0.219	
Europese_samenwerking	0.219	
individuele_huursubsidie	0.219	
pensioenen	0.219	
pensioenfondsen	0.219	
thuiszorg	0.219	
WAO_WAOers	0.219	
wetshandhaving	0.219	
Ziekenfonds	0.219	

Classification experiments

- ▶ We aim to develop a classifier which assigns themes to unseen Dutch election manifestos written after Lipschits' work
- ▶ We have to rely on the older data from the eighties and nineties for training and optimization of the classifier
- ▶ System was tuned by testing on 1998 data, while using older data as training material
- ▶ Balanced Winnow, implementation in the Linguistic Classification System (LCS)

Results on 1998 data

Table : Results for the automatic classification of the 826 texts from the 1998 data training data

	# themes	Precision	Recall
1986 and 1994 data	320	68.8%	37.3%
1986 and 1994 data, only 1994 themes	211	72.2%	37.1%

Adding bigrams, removing stopwords and/or lemmatization did not improve these results

Labelling new, unseen election manifestos

Table : Statistics on the unseen data

year	# words	vocabulary size	# texts	# words/text
2006	235,949	35,547	4,771	49.5
2010	756,254	35,547	21,329	35.5
2012	246,004	30,477	5,880	41.8

Evaluation: A sample of the automatically labelled data was manually evaluated by expert (political journalist)

Text to be assessed

We are fully committed to employment by making European funds available for investment in education, research, innovation, energy efficiency, infrastructure and digitization. These investments are financed or guaranteed by the relevant EU funding programmes, the European Investment Bank, European project bonds, making better use of Structural Funds and the European Social Fund, and alternative exploitation of agricultural subsidies.

This text cannot be assessed thematically (keep form empty)

Assigned themes

- **European cooperation**
 Relevant Not relevant I don't know
Comment (optional):
- **employment**
 Relevant Not relevant I don't know
Comment (optional):
- **infrastructure**
 Relevant Not relevant I don't know
Comment (optional):

Missing theme (optional)

-

Evaluation

- ▶ 193 texts from 2006-2012 data were assessed by the expert
- ▶ In addition, she judged 50 text fragments from the 1998 data (manually labelled by Lipschits, but the expert was not aware of that)

Table : Results for the automatic classification of 193 text fragments from 2006 onwards into the 218 themes defined by Lipschits in 1998. The themes were manually evaluated by an expert.

Year	# texts	# themes/text	Precision	Recall
2006 (trained on 86+94+98)	39	1.7	74.2%	45.8%
2010 (trained on 86+94+98)	122	1.9	77.4%	55.0%
2012 (trained on 86+94+98)	32	2.0	84.6%	56.1%
1998 (manually by Lipschits)	50	7.7	71.5%	89.0%

Conclusions

1. In the thematic classification of political texts, a high level of detail seems to be preferred by domain experts
2. Change of themes over the years affects recall of the learned classifier,
3. but precision is comparable to the precision of annotations from a human expert
4. Thus when using old political texts to classify new texts, work is needed to link and expand the set of themes to newer topics.