**WZB**

**Text-mining international politics**
## Session 5
## Text scaling approaches



Dr. Christian Rauh – Block seminar Charles University Prague     1

1

---

**WZB**

# Automated scaling of texts

– Scaling techniques …

   … automatically distribute documents across a latent (underlying) scale (dimension)

   … are used to infer the position of a document's author

   … were mainly developed in studying the ideological positions that drive party manifestos or political speeches (left–right dimension)

   … are increasingly applied to other questions such as lobbying success, e.g.
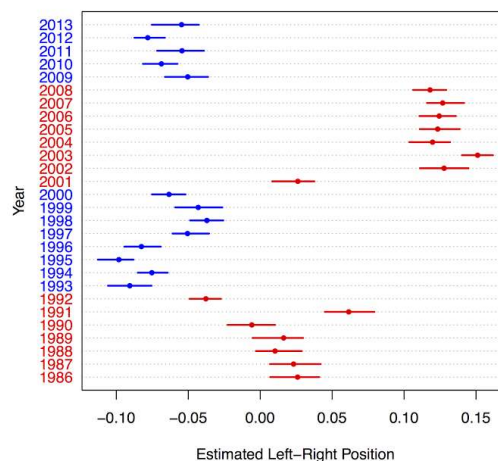
– Basic idea

   Estimate text positions by focussing on language that discriminates most strongly among the texts (i.e. give strong weight to terms that occur very frequently in some texts but only very infrequently in others)

2

2

**WZB**

# Unsupervised scaling

- *Wordfish* (Slapin and Proksch 2008)
  - Assumes that there is only exactly one latent dimension structuring the text corpus!
  - Algorithm weights term frequencies so that that there is a maximum distance between the texts in the corpus
  - Rare terms influence the results strongly
  - Resulting positions can only be interpreted relative to each other
  - Content of the scale has to be interpreted ex-post

- *Wordshoal* (Lauderdale and Herzog 2016)
  - Two-stage approach:
    1. Scale variation in word usage with Wordfish for specific 'debates'
    2. Use factor analysis to construct a common scale across debates
  - Relaxes assumption of uni-dimensionality
    Discriminating power of individual words may vary across debates
  - Geared more towards scaling latent *actor* rather than *text* positions
  - Still unsupervised, interpretation only *ex–post*

3

3

**WZB**

# Wordfish example

**US State of the Union Address Positions**



Source: The Monkey Cage / Benjamin Lauderdale          4
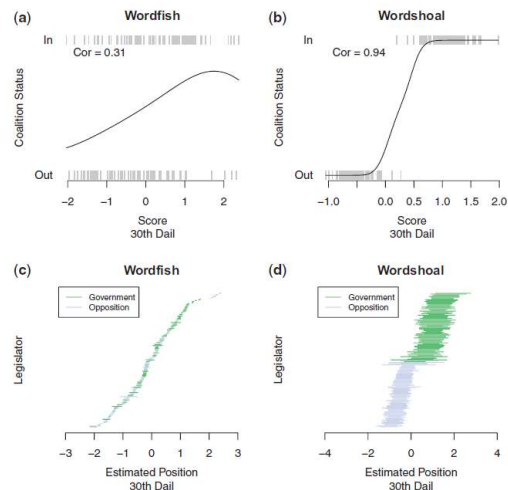
4

**Wordshoal validation example**

Fig. 2 The association between the estimated positions of each legislator and their status as members of the coalition versus opposition, with correlation and local linear smooth, under Wordfish (left) and our approach (right), for the 30th Dáil. In the bottom row, we show the 95% intervals associated with the estimates for each legislator under Wordfish (left) and Wordshoal (right).
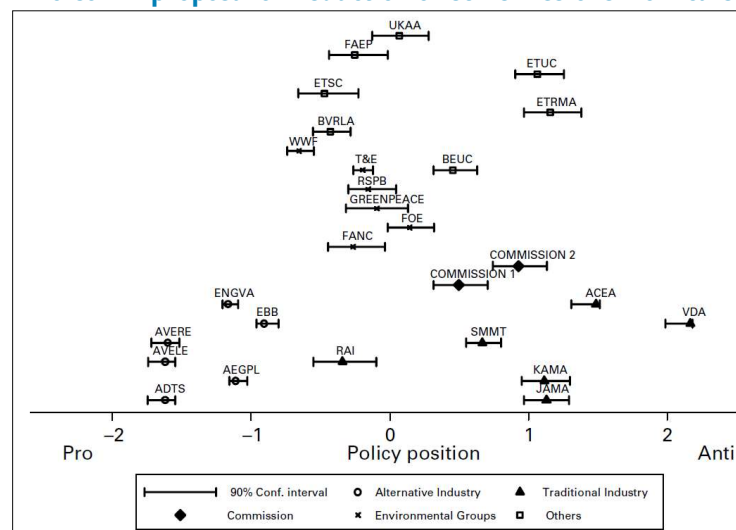
Source: Lauderdale / Herzog (2016, *Political Analysis 24*: 382)

5



**Wordfish example**
**EU Comm proposal on reduction of CO2 emissions from cars**

Source: Klüver (2009, European Union Politics 10(4): 543)

6

**WZB**

## Supervised scaling

- *Wordscores*   (Laver, Benoit and Garry 2003)
  - o Researcher supplies reference texts with 'known' values across the latent scale(s)
  - o Algorithm retrieves and weights the relative term frequencies in these texts
  - o Virgin texts are then positioned on the latent dimension along the weights of the terms they contain
  - o Cf. machine learning in Session 4

7

7

---

**WZB**

## Applying Wordfish
## to our running example

- What differentiates national delegates in the United Nations General Assembly according to the relative frequency of words they use when speaking about climate change?
- ➢ And: Does this meaningfully capture expressed political *positions* on climate change issues?

- Approach
- ➢ Apply the Wordfish algorithm (as implemented in *quanteda*) to the corpus of 100-term window around climate change references aggregated to country (!pre-processing!)
- ➢ Scrutinize term weights ('betas') and document positions ('thetas')

8

8

Wordfish: Estimated word stem positions in UNGA climate change talk

9



**Wordfish positions of climate change related speeches in UN General Assembly**

Wordfish theta estimates of 100-term windows around 'climate change' / 'global warming' pooled by country

10

**WZB**

**How do states position themselves on climate change?**
Explaining the estimated Wordfish position around climate change references by some crude national-level variables



Standardized regression coefficient
with 95 and 99% confidence intervals

Linear Model, n = 191 countries, Adj. R2: .46.

11

---

**WZB**

## Applying Wordscores
## to our running example

- In how far do speeches of national delegates in the UNGA use language of climate sceptics or climate activists?
➢ And: Does this meaningfully capture expressed political *positions* on climate change issues?

- **Approach**
o Corpus of 3000+ reference texts: scrape c*limate–change related news* (!) from websites of *The Heartland Institute* (climate change sceptics or deniers; reference score: –1) and *The Ecologist* (climate activists; +1)
o Train a Wordscores model via *quanteda* on this corpus and analyze the resulting term weights
o Scale UNGA speeches (pooled by country) along this model and see whether we find something meaningful

12

**Example of a Wordscores model**
**Identify language that separates climate change skeptics and activists**

<- Skeptics | Activists ->

carbon-dioxide
american
causing
attorney
lehr
computer
state
scientists
regulatory
taylor
sources
predictions
curry
rule
assertions
endangerment
mann
nipcc
hurricanes
americans
scientific
tax
panel
temperatures
atmospheric
concerning
claims
agency
skeptics
temperature
alarmism
mandates
restrictions
hurricane
ipcc
science
gore
source
peer-reviewed
data
noaa
regulations
watts
dioxide
cpp
models
warming
administration
alarmist
alarmists

oxymoron
entitlements
primed
dips
first-ever
trope
relics
ky
commendable
aluminum
six-fold
mpg
oftentimes
tears
latched
ideologically-driven
sharpen
exhorting
financially
slices
evermore
cassava
regularity
impulse
males
providence
bathtub
air-conditioners
shaving
uniformity
airtight
programr
instinctively
sped
honeymoon
rift
four-year-old
creed
quake
culpability
compressor
exerts
transitioned
immediacy
demolishing
motorized
zero-emission
resource-rich
acquisitions
woeful

transport
copenhagen
communities
sustainable
tonnes
indigenous
aviation
trading
britain'
deforestation
towards
tar
campaigners
shell
forest
local
community
camp
tackling
amazon
labour
tackle
meat
food
protest
need
airport
biodiversity
ecological
heathrow
farming
footprint
investment
around
mc
corporate
waste
us
conservation
rainforest
consumption
organic
bank
police
eu
destruction
british
finance
extraction

Wordscores
-0.4    0.0    0.4

Reference texts to train the WS algorithm are climate-change related news from two outlets:
1) 'Heartland Institute' (Skeptics, reference score = -1, n = 1,322)
2) 'The Ecologist' (Activists, reference score = 1, n= 1,851)

Plot shows the terms that are closest to the minimum, mean, and maximum of the estimated score distribution
(61,675 terms scored in total)

13

---

**Wordscores positions of climate change related speeches in UN General Assembly**
Wordscores estimates of 100-term windows around 'climate change' / 'global warming' references pooled by country

*Closest to Climate Activist Language:*
Iraq
Syria
Ireland
Lebanon
Colombia
Afghanistan
Dominican Republic
United States
Romania
Israel
Djibouti
Côte d'Ivoire
New Zealand
Bolivia
United Kingdom
Senegal
Malaysia
Liechtenstein
France
Thailand
Malawi
Finland
Iran
Indonesia
Brazil

*Closest to Climate Sceptic Language:*
Sierra Leone
Austria
Sudan
Montenegro
Albania
St. Vincent & Grenadines
Mauritania
San Marino
Japan
Pakistan
Kyrgyzstan
Egypt
Equatorial Guinea
Mongolia
Trinidad & Tobago
Cuba
United Arab Emirates
Suriname
Morocco
El Salvador
Kuwait
Turkmenistan
Moldova
Azerbaijan
Russia

0.24    0.28    0.32

14

**WZB**

**How do states position themselves on climate change?**
Explaining the estimated Wordscores position (skeptics/activists) around climate change references by some crude nation
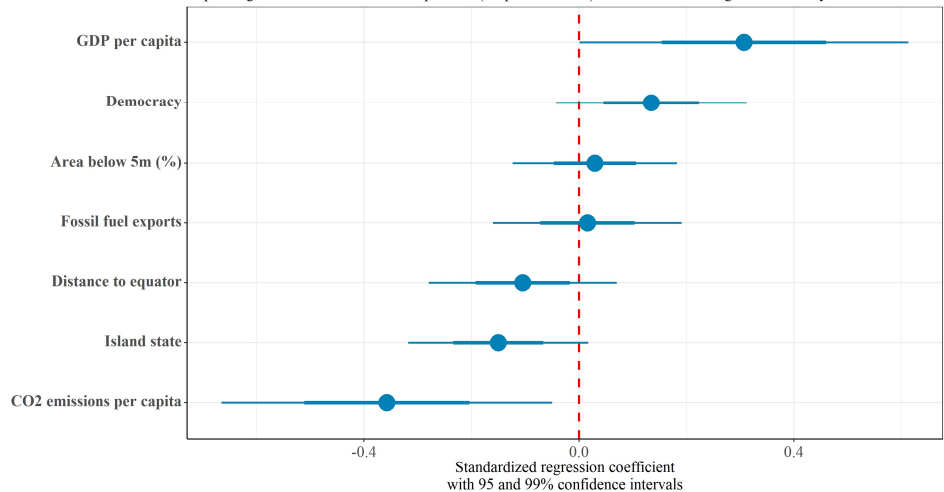


GDP per capita
Democracy
Area below 5m (%)
Fossil fuel exports
Distance to equator
Island state
CO2 emissions per capita

-0.4    0.0    0.4

Standardized regression coefficient
with 95 and 99% confidence intervals

Linear Model, n = 191 countries, Adj. R2: .02.

---

**WZB**

## Pitfalls of automated scaling

– Scaling works only:
  o with documents that are very focussed on the theorized dimension (cf. party manifestos vs. newspaper articles)
  o if documents come from the same context in which the language is used identically (political speeches vs. news outlets?)

$\Rightarrow$ Scaling procedures make strong assumptions!
$\Rightarrow$ Scaling procedures require particularly careful validation!