

Text-mining international politics

Session 2

Corpus construction and discovery

s pre...
g Ltd. Keiler Snow is a **research** and development arm for the...
and to act as a kind of **research** world and the venture capital p...
the public sector driven **research** charity. This after the 1998 £5...
and is the world's largest **research** environment within UK univers...
o transform the scientific **research**. This latter point seems to b...
tion into a zone for social **research**. has established that structur...
natic ring. In recent years, **research**. Filed under: Media. Poster...
disagree with. Do your own **research**. it seems my options are t...
time in my life. After doing **research**, [a non-profit **research** institute] because we're st...
y US employer [a non-profit **research** by looking at how literatur...
rested in supplementing this **research**. What City of Quartz doe...
ig. Not for travel, for a kind of **research** (both internet and olde-fa...
o high in minerals? Extended **research** for a friend who needs to...
it. I've been asked to do some **research** to middle schoolers. Mc...
eachers to use when teaching **research**, and wonder why anyo...
ernet is the ultimate source for **research** thing and now I'm com...
... since I did the whole **research**, I've tried to make my...

Acquiring documents

- **Which kinds of texts are suitable for automated content analysis?**
 - Unit of analysis is usually on *document level* (other units possible)
 - *Focussed* documents preferable (depending on your theoretical concepts)
 - Sufficient *number of words* required (depending on the applied method)
- **Typical text sources**
 - *Existing corpora*: Other social science projects or linguistic resources
 - *Online databases*: e.g. LexisNexis, Factiva, Gale Cengage (newspapers and press agencies); governments, parliaments, international orgs; etc...
 - *Web scraping*: Press releases, news sites, blogs, Twitter, etc.
(see Munzert et al., 2015, Wiley)
 - *Scan / OCR* of printed matter
- **Store documents with consistent formats and document names**
 - Plain txt files work best (converters freely available)
 - UTF-8 encoding standard for Latin alphabet

WZB

Readily available corpora – Examples I

- **UNGD corpus**
 - Baturo, Dasandi, Mikhaylov (2017, *Research and Politics*)
 - Full texts of all 8,093 speeches by national representatives in the *United Nations General assembly* between 1970 and 2018 (includes auto-translated versions to English)
 - <https://doi.org/10.7910/DVN/OTJX8Y>
- **EUspeech**
 - Schumacher, Schoonvelde, Traber, Dahiya, de Vries (2016)
 - 17,184 speeches of national, EU, and IMF executives during the 2007–2015 period (includes auto-translated versions)
 - <https://doi.org/10.7910/DVN/GKABNU>

3

3

WZB

Readily available corpora – Examples II

- **ParlSpeech**
 - Rauh, De Wilde, Schwalbach (2017, Harvard Dataverse)
 - 3.9 million plenary speeches in key parliamentary chambers of seven European states (including CZ/PSP 1993–2016)
 - www.bit.ly/ParlSpeech
- **Manifesto corpus**
 - Full text of partisan election manifestos from 50 countries
 - <https://manifesto-project.wzb.eu/information/documents/corpus>
- **US SOTU corpus (1790–2018)**
- **Annual IO reports, UNSC resolutions, EU legislation or:
*scrape/build your own!***

4

4

WZB

Pre-processing: Turning text into data (Typical, but sensitive and not universally applicable steps!)

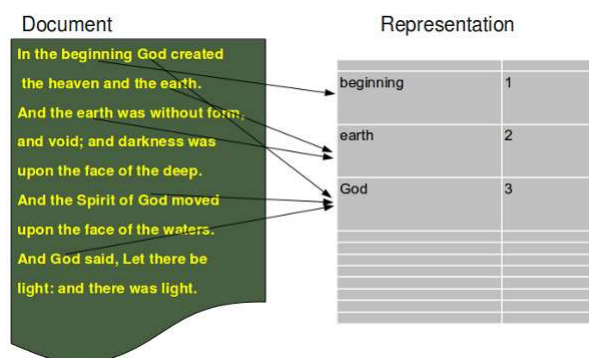
- **Remove...**
 - ... document “boilerplate” (info not part of the analysed message)
 - ... punctuation, capitalization, numbers
 - ... very common and very uncommon terms (“stop words”, <1% of docs)
 - **Lemmatization / Stemming**
 - Words referring to the same concept mapped to a single root
 - {economy, economic, economically} → economi
 - **Turn documents into “bags of words”**
 - Discards the order in which words occur!
 - Unigrams, bigrams ... n-grams
- ⇒ Document frequency matrix

5

5

WZB

“Bags of words” (illustration w/out stemming)



Source: python-course.eu

6

6

Document frequency matrix (illustration)

Label	Titles
c1	Human machine interface for Lab ABC computer applications
c2	A survey of user opinions of computer system response time
c3	The EPS user interface management system
c4	System and human systems engineering: testing of EPS
c5	Relation of user-perceived response time to error measurement
m1	The generation of random, binary, unordered trees
m2	The intersection graph of paths in trees
m3	Graph minors IV: Widths of trees and well-quasi-ordering
m4	Graph minors: A survey

	Terms	Documents									
		c1	c2	c3	c4	c5	m1	m2	m3	m4	
	computer	1	1	0	0	0	0	0	0	0	
	EPS	0	0	1	1	0	0	0	0	0	
	human	1	0	0	1	0	0	0	0	0	
	interface	1	0	1	0	0	0	0	0	0	
	response	0	1	0	0	1	0	0	0	0	
	system	0	1	1	2	0	0	0	0	0	
	time	0	1	0	0	1	0	0	0	0	
	user	0	1	1	0	1	0	0	0	0	
	graph	0	0	0	0	0	0	1	1	1	
	minors	0	0	0	0	0	0	0	1	1	
	survey	0	1	0	0	0	0	0	0	1	
	trees	0	0	0	0	0	1	1	1	0	

Source: <http://web.eecs.utk.edu/~mberry/order/node4.html>

7

7

The potential of discovery

- Even if you do not want to apply statistical analyses to your corpus, a look at the aggregated term frequencies may:
 - ... show unknown temporal patterns
 - ... provide contextual information for specific concepts (co-locations, keyword-in-context, synonyms, ...)
 - ... guide selection of individual texts for further human interpretation and coding
 - ... give you an aggregated perspective on the discourse helping to contextualize individual documents therein

8

8

WZB

Introducing our running example

▪ Speeches in the United Nations General Assembly

Based on the UNGD corpus assembled by *Baturo, Dasandi, Mikhaylov (2017, R&P)*

▪ What is the context of these documents we need to have in mind along the sender-message-recipient framework?

- Who speaks when? With what purpose?
- The examples will (try to) make inferences about the 'senders', assuming that speeches reflect state positions along the words used

▪ Climate change as the political issue of interest

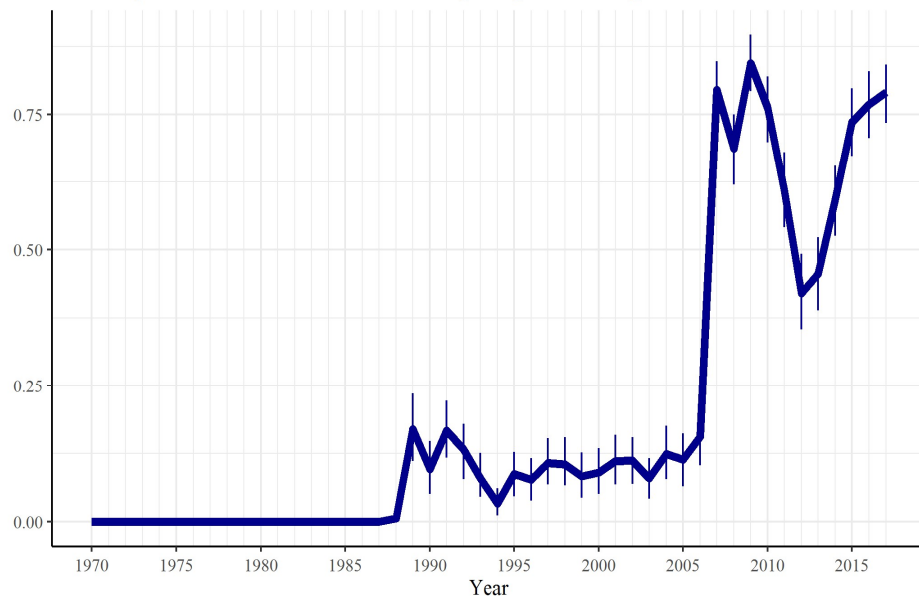
- Identified by speeches referring literally to 'climate(-/)change' or 'global(-/)warming'
- 100-term window around these references to see how and what national delegates say about or associate with climate change
- For a 'real' analysis, more fine-tuning will most likely be needed, bear with me...

9

9

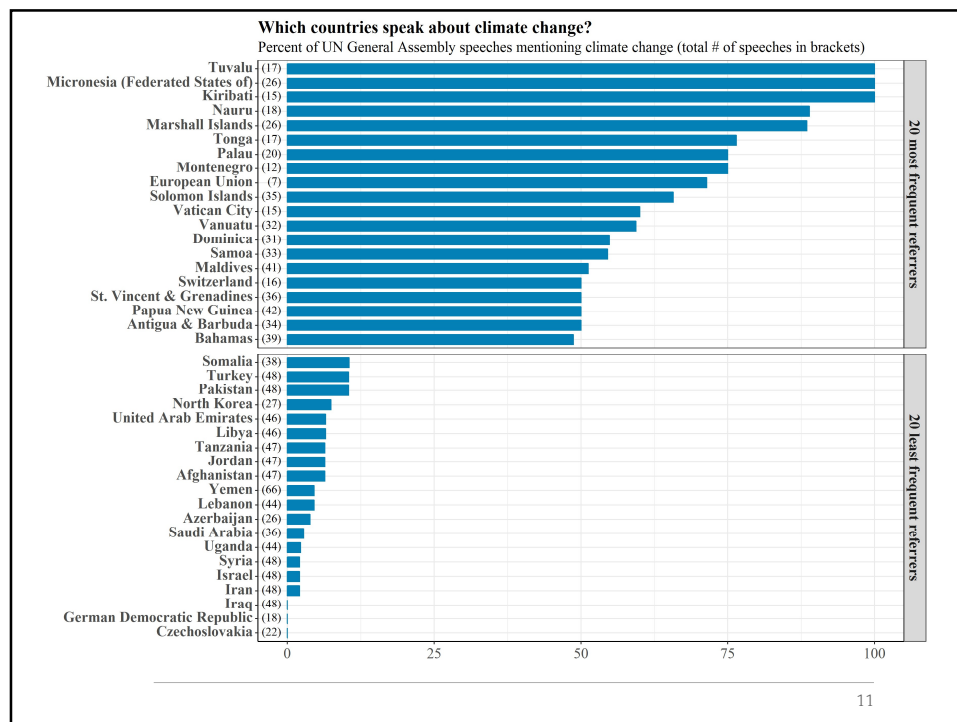
Climate change salience in UN General Assembly speeches over time

Share of speeches with references to 'climate change' or 'global warming'

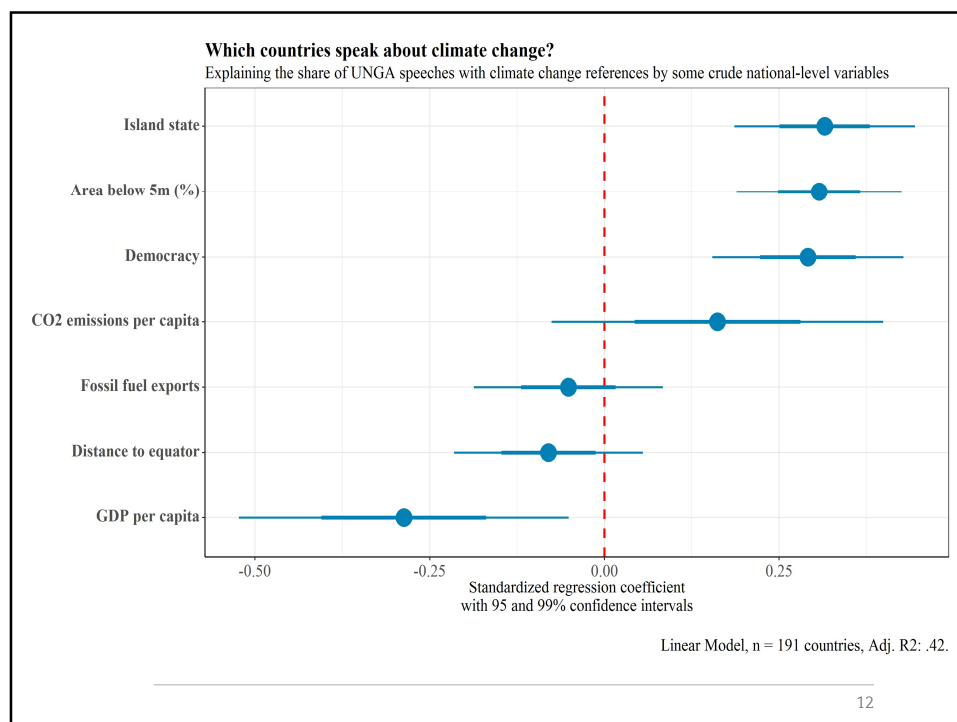


10

10



11



12

[illegible]

15