

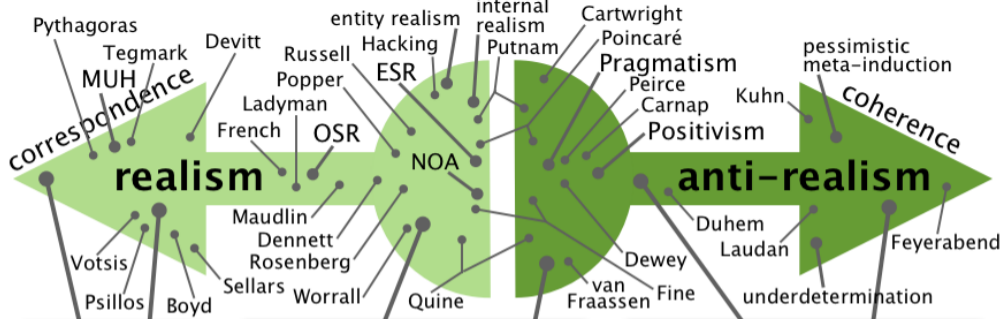
Text as Data

Juraj Medzihorsky



2016-11-28

?



Naive Realism

The world I see is real. What are you all arguing about?

Structural Realism

Science has identified real patterns, relationships, and structures (at least within a regime) in nature.

Instrumentalism

Theoretical concepts may have use in predicting observations, but we have no ontological commitments to them.

Scientific Realism

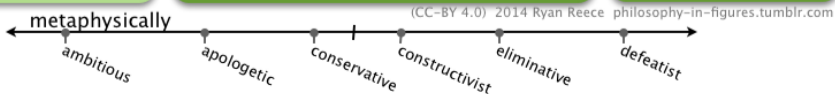
Science makes real progress in describing real features of the world.

Constructive Empiricism

Science aims to give us theories which are empirically adequate, but does not justify metaphysical claims about reality.

Relativism

Social constructivism. Epistemological anarchism.



Motivation

?

*all models are false,
but some are useful*

George E. P. Box

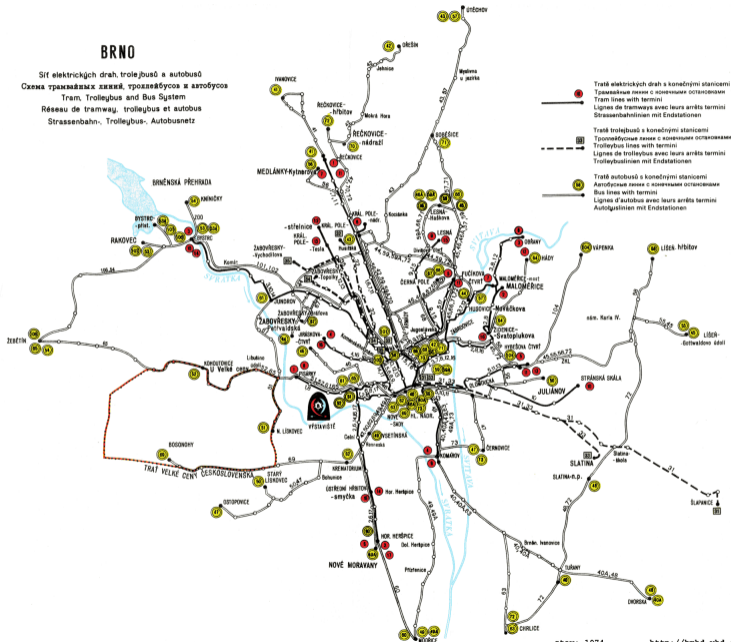


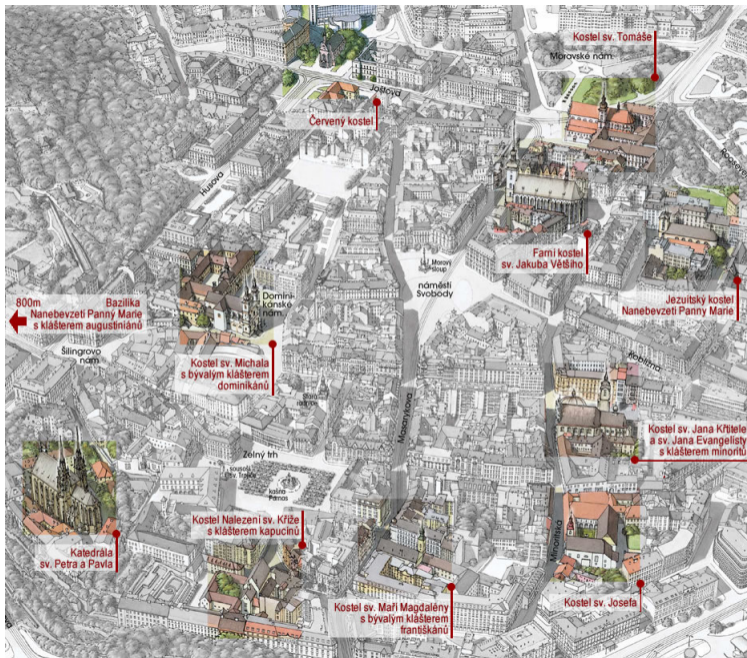
Brno

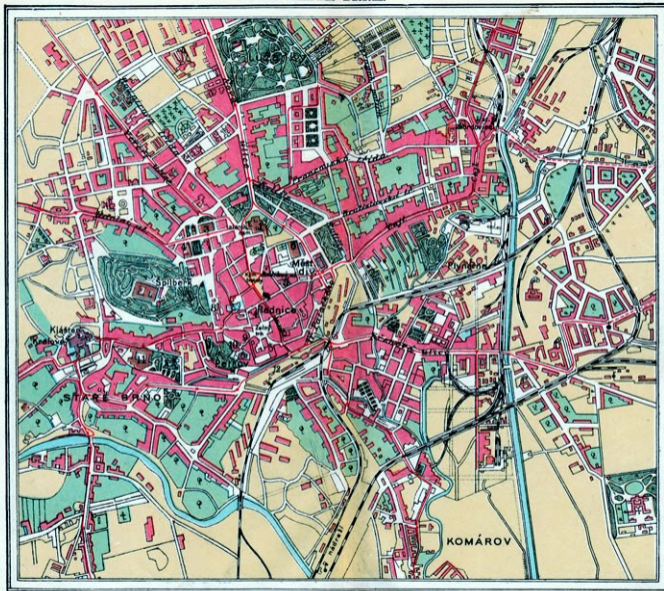


BRNO

Sif elektrických drah, trolejbusů a autobusů
Схема трамвайных линий, троллейбусов и автобусов
Tram, Trolleybus and Bus System
Réseau de tramway, trolleybus et autobus
Strassenbahn-, Trolleybus-, Autobusnetz







Měřítko 1:25.000

Význačné budovy:

Reprodukce Voj. zeměp. ústavu v Praze.

- 1 Zem. správa politická.
- 2 Národní divadlo.
- 3 Starý zemský dům.

- 4 Rádnice.
- 5 Chrám sv. Jakuba.
- 6 Nejvyš. s. dvůr.

- 7 Ředitelství pošt a telegrafů.
- 8 Průmyslové museum.
- 9 Zemské museum.
- 10 Dům.
- 11 Zem. nemocnice.
- 12 Státní nádraží.

Document Scraping



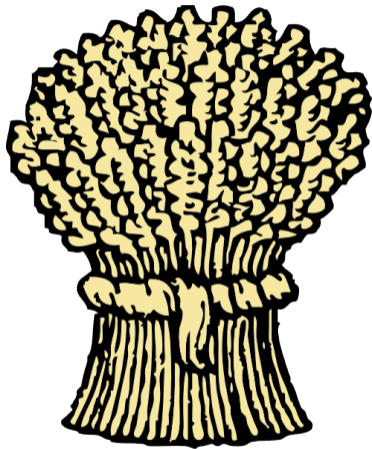
Document Scraping

- Numbers and text in files
- Local
- Web





- eXtensible Markup Language (XML)
- APIs (e.g. for Twitter)



Text Analysis

눈물 아롱아롱
피리 불고 가신 님의 밟으신 길은
진달래 꽃비 오는 서역 삼만리
흰 옷깃 여며 여며 가옵신 님의
다시 오진 못하는 파촉 삼만 리

Handwritten text in a medieval script, possibly Gothic or similar, arranged in a grid pattern. The text is written in red ink on a light-colored parchment or paper. The grid is formed by horizontal and vertical red lines. The text is arranged in approximately 10 rows and 15 columns. The script is highly stylized and dense. The text is arranged in a grid pattern, with each character fitting into a square cell. The text is arranged in a grid pattern, with each character fitting into a square cell. The text is arranged in a grid pattern, with each character fitting into a square cell.

IF YZIDAI NCH IYI NE QYMLQ

IE' DFKAI LEQT E'G' LGI DCKDEPDEI XEI' GP
PEY A SER PCHIQ' DEITXN' GP NCHN' Q'
DEITXN' FGE' YPK' LEQT

QYMLAI EGI

'ZIQZ' SGIYZI FGE' GP PEYI' A PCHIQEN
DEITX RGI DFK DEITPK' A SEPE'Y' KFYGIYAI
HEP DELEQY' MEI' QFC FEL HQLZ' HFEPI'K'Y' A
PER' GP' LCHDEN NCH' EY MCH DCKLQ' HPIZF'Y' A
PER'



golland ceas ceethy ollas g beag ollas ceedg llor g ollceedg gollas & ceedg
 gollceedg gllceedg gollas of ceedg gollceedg gollceedg gollceedg ceedg sand ceedg
 dar ceedg ceethy gollas gollceedg ceedg gollceedg gollceedg gollceedg ceedg sand ceedg of
 the llor ceedg gollceedg ollas ceedg gollas ceethy gllceedg gollceedg ceedg
 ceedg gollceedg olland ceethy ceedg ollceedg gollg gollceedg of ollceedg
 gollas ceethy ollg dar ceethy ceedg gollceedg ollceedg



gollceedg ceethy ollceedg dar golland o gollceedg of and gollg
 goll of gollg gollg of gollceedg ceethy olland ceedg gollceedg
 gollceedg ollceedg dar gollceedg gollceedg gollceedg gollceedg gollceedg
 gollceedg gollceedg of ceethy gollceedg gollceedg gollceedg gollceedg
 gollceedg of gollceedg gollceedg gollceedg gollceedg gollceedg gollceedg



?

Text Analysis: The Big Picture

Text Analysis

- Discourse
- Content

Content Analysis

- 'Manual' Text Analysis
- Computer-Assisted Text Analysis (CATA)

'Manual' Text Analysis

- Humans do most of the work
- Expensive
- Slow
- Reliability issues

Computer-Assisted Text Analysis

- Computers do most of the work
- Boom
 - Huge amount of digitized text available
 - Cheap computing power
 - New methods – CS & PS

Political Text in CATA

Examples:

- Manifestos & platforms
- Press releases
- Social media content
- Floor and debate speeches



With /

TOW

'Bag of Words'

- Common assumption in CATA
- Order of words (**n-ngrams** of words) does not matter
- Texts as vectors of word counts

Table 2: A word-frequency matrix for two randomly selected sentences from the corpus: a) ‘Let’s enforce the laws already on the book.’ (Cain in 2012); b) ‘Now, the administration has \$800 million on hand right now, cash on hand.’ (Hunter in 2008)

Stem	Sentence	
	a	b
administr	0	1
book	1	0
cash	0	1
enforc	1	0
hand	0	2
law	1	0
million	0	1

CATA in Political Science

- Grimmer, J., & Stewart, B. M. (2013). Text as data: The promise and pitfalls of automatic content analysis methods for political texts. *Political Analysis*, mps028.

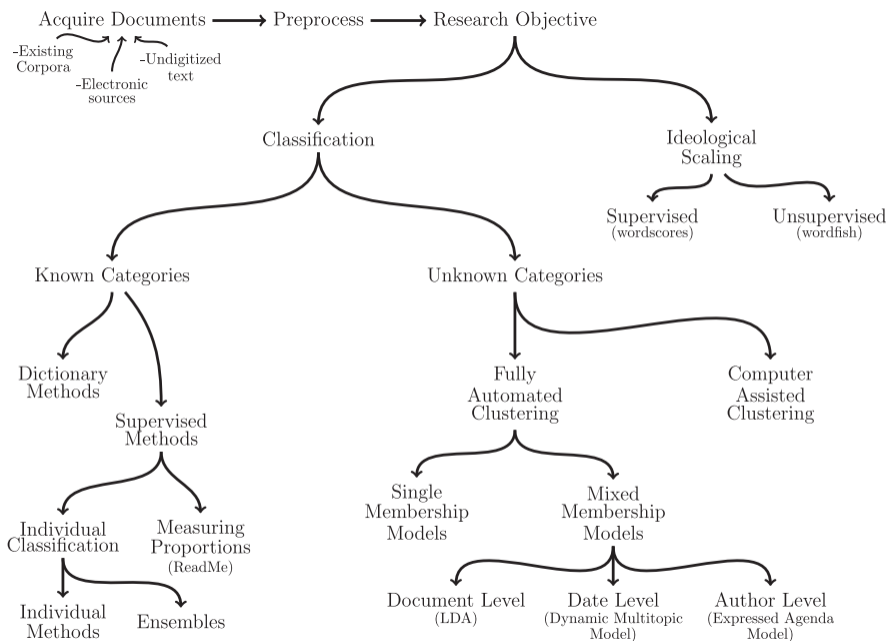


Fig. 1 An overview of text as data methods.

Dictionary-Based Methods

Google N-Grams

Dictionary-Based Methods

- Build a dictionary
- Statistical models are not necessary

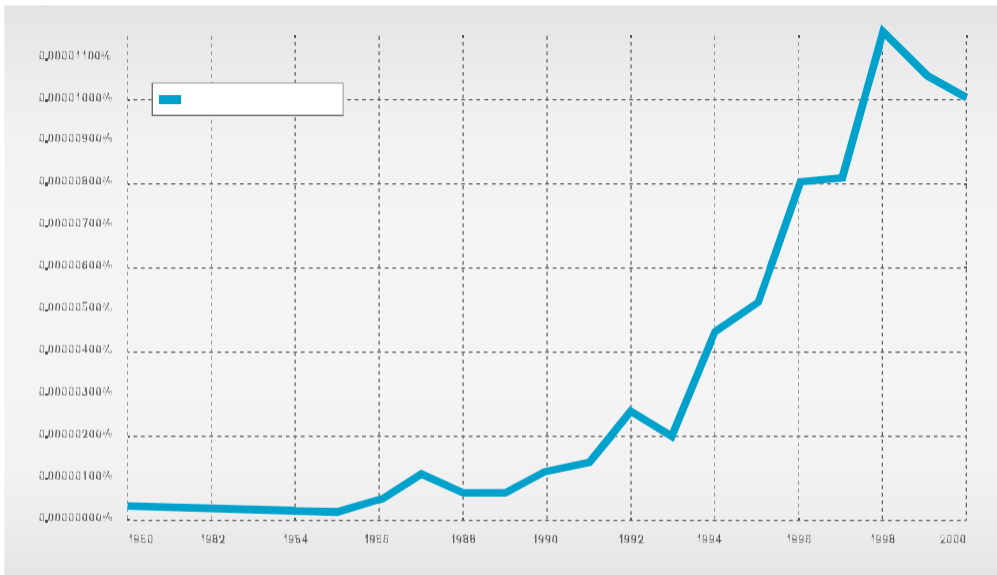
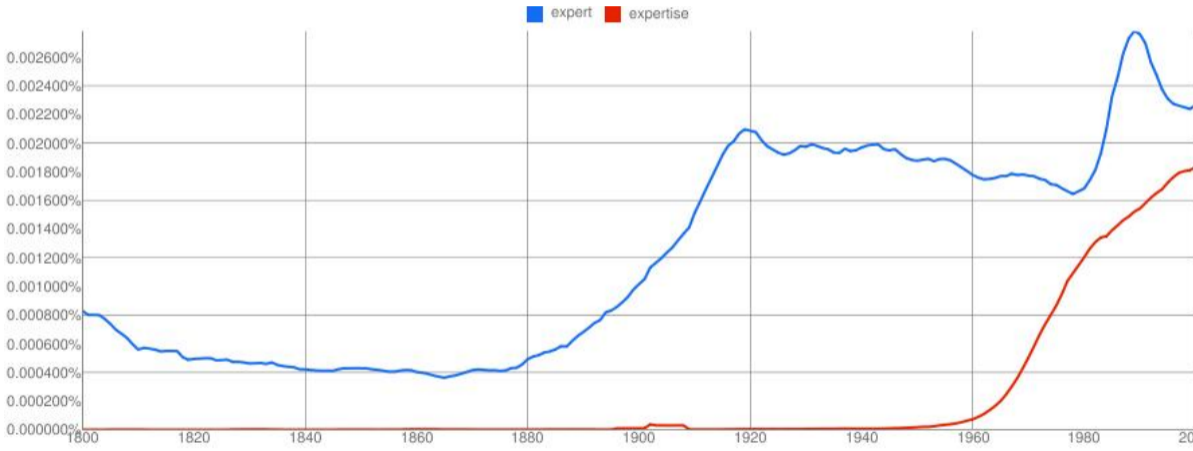


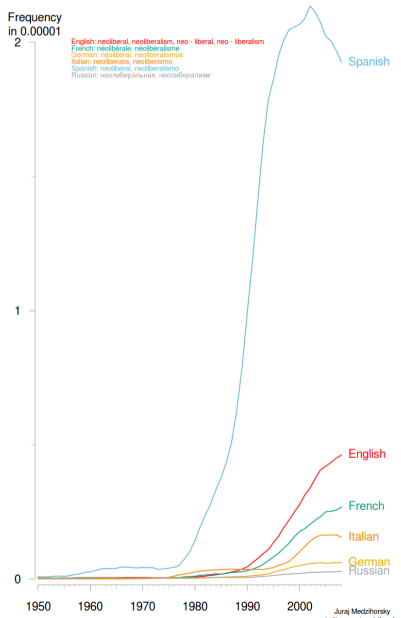
Figure 1: Frequency of appearance of “public intellectual” in Google Books from 1980 to 2000

Source: Graph. *Google Books*. Google, 18 June 2012. Web. 18 June 2012.

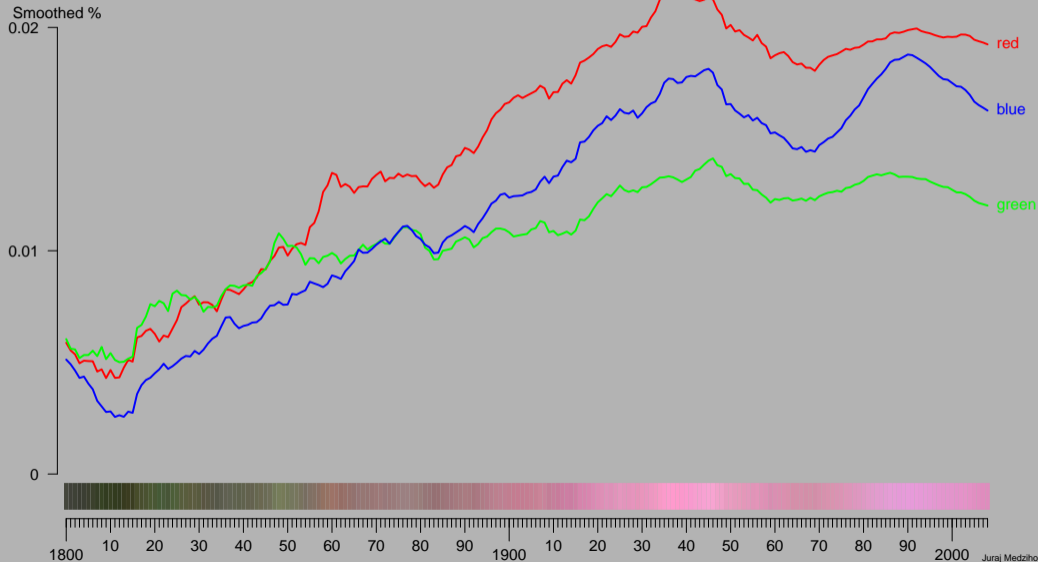
Figure 1: Frequency of appearance of “Expertise” and “Expert” in Google Books from 1800 to 2000



'Neoliberalism' in Google Book Corpora



Red, Green, and Blue in Google Books English Fiction Corpus



Scaling

Scaling Goals

- One or more dimensions
- Place documents (texts, speeches) in a space
- Place words in the same space

Common Scaling Methods

Supervised: Wordscores

Unsupervised: Wordfish

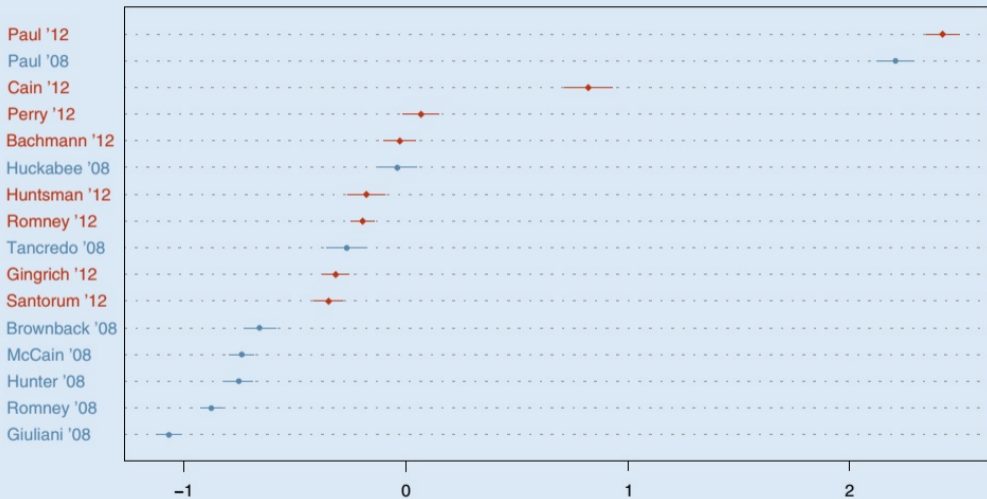
Unsupervised: Correspondence Analysis

A Scaling Example

- 2008 and 2012 Republican presidential primaries
- Debate transcripts from a UCSB website
- Expected move towards Tea Party positions
- Unsupervised scaling: Wordfish

Figure 1

Candidate Positions



Candidate positions extracted from their pre-lowa debate speeches with bootstrapped 95% confidence intervals (1,000 replications). 2008 candidates denoted by circles (blue) and 2012 candidates by diamonds (red).

Table 4

Selected Five-Sentence Sequences Spoken by a Single Candidate in a Single Debate

NEGATIVE, -1 ± 0.1

"I have joined together across the aisle on a number of pieces of legislation, many of them very important. I'm proud of my legislative record of conserving my ideals and my conservative principles and getting things done in Washington. And I am proud of that, and I will continue to hold to those ideals. But I will reach across the aisle to the Democrats who I have worked with, who know me, and we know we can work together for the good of this country. Let's raise the level of dialogue and discussion and debate in this campaign." (McCain on December 12, 2007; score: -1.1)

"It's the one place I found to agree with President Obama. If every parent in America had a choice of the school their child went to, if that school had to report its scores, if there was a real opportunity, you'd have a dramatic improvement. I visited schools where, three years earlier, there were fights, there were dropouts, there was no hope. They were taken over by a charter school in downtown Philadelphia, and all of a sudden the kids didn't fight anymore, because they were disciplined. They were all asked every day, what college are you going to? Not are you going to go to college, what college are you going." (Gingrich on September 7, 2011; score: -1)

"I can tell you a good union, the Steel Workers Union. When last year, Chris, we had a strike in a Kansas plant that made the tires for our humvees, I called up the president of the Steelworkers and the president of Goodyear, and within a very short period of time, they were working together, they got that thing done for the good of the country. A union is a receptacle of power, just like management. But those folks love this country, they love their family, and they helped to build a middle class, which has been important for America and for our party. We need to work with unions to win this presidency." (Hunter on October 9, 2007; score: -0.9)

POSITIVE, $+1 \pm 0.1$

"Repeal Dodd-Frank, repeal Obamacare. It really isn't that tough if you try. It is easy to turn around this economy, just have the backbone to do it. Well, as president of the United States, I would not be reappointing Ben Bernanke, but I want to say this. During the bailout, the \$700 billion bailout, I worked behind the scenes against the bailout, because one of the things that I saw from the Federal Reserve, the enabling act legislation is written so broadly that, quite literally, Congress has given the Federal Reserve almost unlimited power over the economy." (Bachmann on September 12, 2011; score: 1.9)

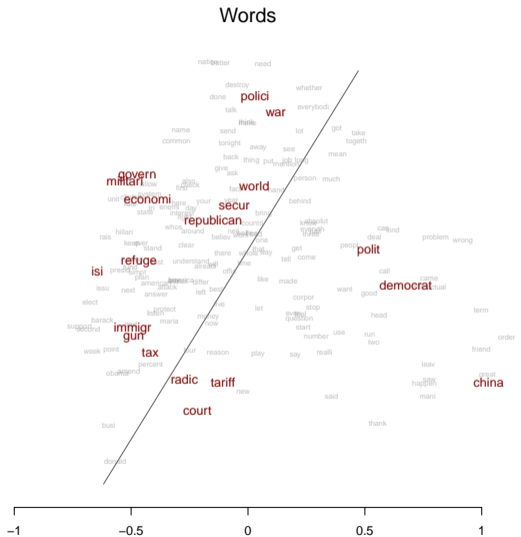
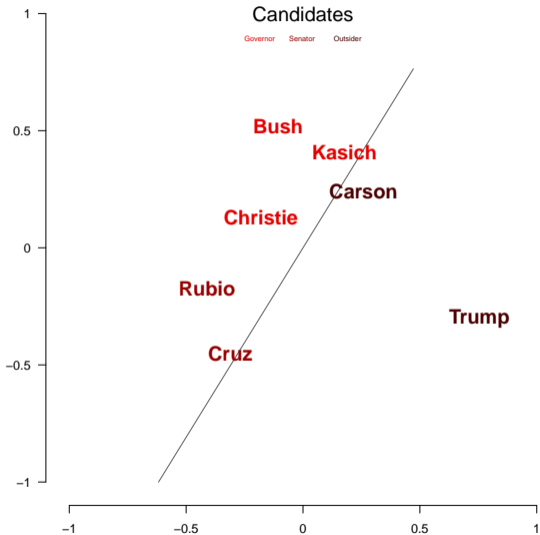
"If we look for it, you'll realize that our national sovereignty is under threat. Yes, and I would like to state that, to the statement earlier made that we all went to Washington to change Washington and Washington changed us, I don't think that applies to me; Washington did not change me. I would like to change Washington, and we could by cutting three programs, such as the Department of Education—Ronald Reagan used to talk about that—Department of Energy, Department of Homeland Security is the biggest bureaucracy we ever had. And besides, what we can do is we can have a stronger national defense by changing our foreign policy. Our foreign policy is costing us a trillion dollars, and we can spend most of that or a lot of that money home if we would bring our troops home." (Paul on November 28, 2007; score: 2)

"There's a responsible way for the federal government to do the things that it should do. Running organizations like the TSA, I would agree with Representative Paul, no. Having the federal government responsible for trying to micromanage Medicare, no, trying to micromanage education, no. The federal government is not good at micromanaging anything. This is why I believe in empowering the states to do more and limit what the federal government does with regard to those kinds of program." (Cain on August 11, 2011; score: 2.1)

Another Scaling Example

- Transcripts of US presidential debates
- Unsupervised scaling with correspondence analysis
- Two dimensions

Sixth GOP'16 Presidential Debate



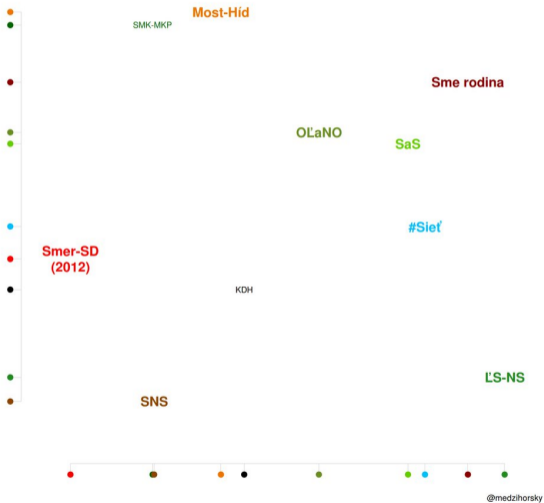
Yet Another Scaling Example

- Slovakian party manifestos: text and CMP codes
- Unsupervised scaling with correspondence analysis
- Two dimensions

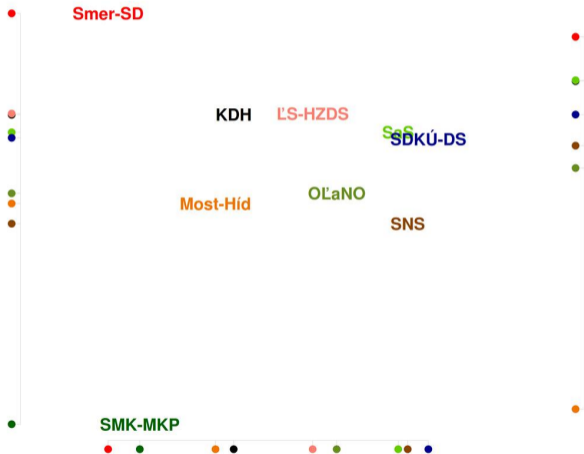
Slovakian Elections 2016: Manifesto Vocabularies

* Smer-SD ran without a new manifesto

Parties with seats in bold type

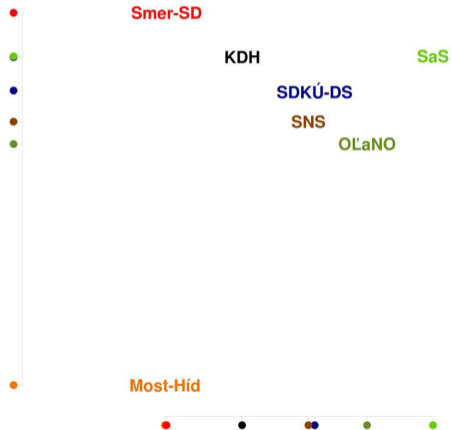


Voľby 2012: programy



@medzhorsky

Voľby 2012: CMP kódy



@medzhorsky

And One More Scaling Example

- Nielsen's (2013) dissertation
- 25,000 + documents by \sim 100 clerics
- *Jihad Score*
- Supervised scaling with a training set

If a **person arrives** while the **Imam** is preaching at **Friday** prayers, he should **pray** two brief prostrations and sit without **greeting** anyone as greeting people in this circumstance is **forbidden because** the Prophet, peace be upon him, says, "If your friend **speaks** to you during the **Friday** prayers, silence him while the **Imam** preaches because it is idle talk."

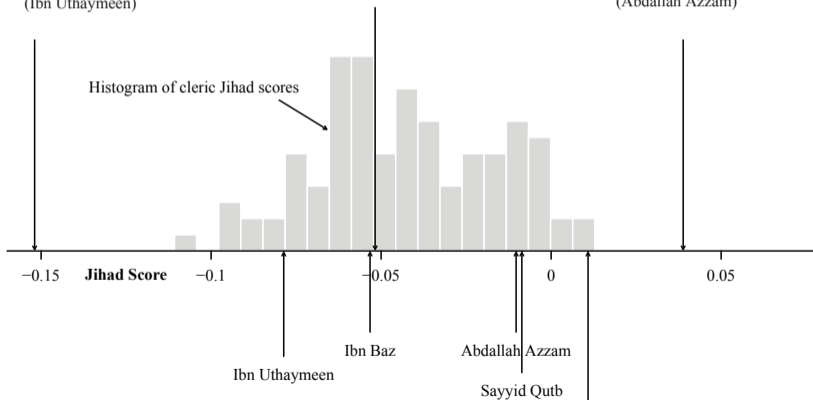
(Ibn Uthaymeen)

There is a **fundamental fact** about the **nature** of this religion and the **way** it works in **people's lives**. A **fundamental, simple fact, but** although it is simple, it is **often** forgotten or not realized **at all**. Forgetting this fact, or **failing** to recognize it arises from a serious **omission** from **views** of this religion: its **truthfulness** and **historical, present, and future reality**.

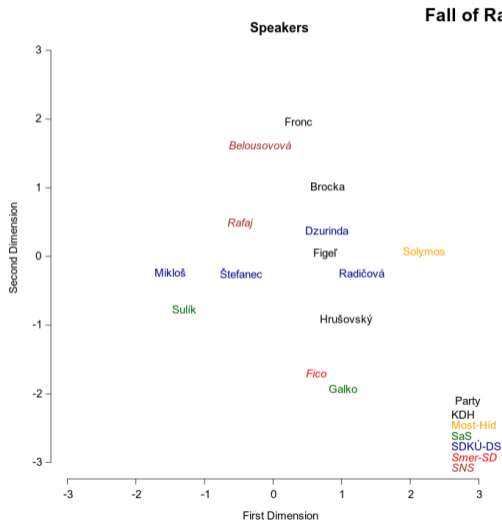
(Sayyid Qutb)

Ruling on **Fighting** Now in **Palestine** and **Afghanistan**. The foregoing **has** clarified that if an inch of Muslim lands are attacked, then **Jihad** is obligatory for the people of that area, and those **near** by. If they do not succeed or are incapable or lazy, the **individual obligation** widens to those behind them and then gradually the **individual obligation** expands until it is general for the whole **land**, from East to West.

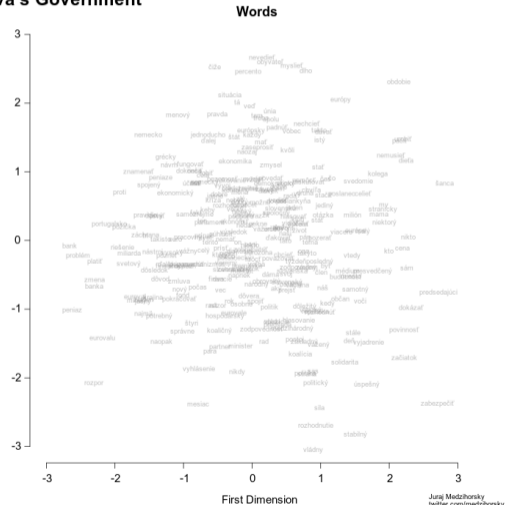
(Abdallah Azzam)



Failed Unsupervised Scaling



Fall of Radičová's Government



'Topic' Models

Topics

gene 0.04
dna 0.02
genetic 0.01
...

life 0.02
evolve 0.01
organism 0.01
...

brain 0.04
neuron 0.02
nerve 0.01
...

data 0.02
number 0.02
computer 0.01
...

Documents

Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many **genes** does an **organism** need to **survive**? Last week at the genome meeting here,* two genome researchers with radically different approaches presented complementary views of the basic genes needed for **life**. One research team, using **computer** analyses to compare known **genomes**, concluded that today's **organisms** can be sustained with just 250 genes, and that the earliest life forms required a mere 128 **genes**. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those **predictions**

"are not all that far apart," especially in comparison to the 75,000 **genes** in the human genome, notes Siv Andersson at Uppsala University in Sweden. "You arrived at the 800 number. But coming up with a consensus answer may be more than just a **genetic** numbers game, particularly as more and more **genomes** are compactly mapped and sequenced. "It may be a way of organizing any newly **sequenced genome**," explains Arcady Mushegian, a **computational** molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an



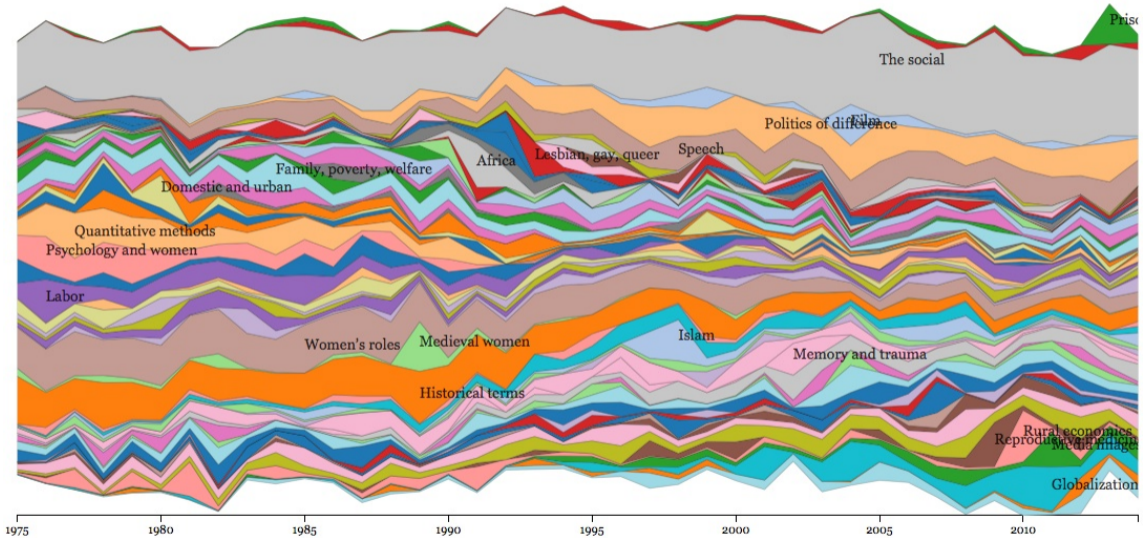
* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

Stripping down. Computer analysis yields an estimate of the minimum modern and ancient genomes.

SCIENCE • VOL. 272 • 24 MAY 1996

Topic proportions and assignments





VisArgue - A Visual Text Analytics Framework for the Study of Deliberative Communication

Mennatallah El-Assady¹, Valentin Gold², Annette Hautli-Janisz³,
Wolfgang Jentner¹, Miriam Butt², Katharina Holzinger², Daniel Keim³

¹Department of Computer and Information Science

²Department of Politics and Public Administration

³Department of Linguistics

University of Konstanz, Germany

`valentin.gold@uni-konstanz.de*`

Abstract

For the last two decades, deliberative democracy has been intensively debated within political science and other related fields. Only recently, deliberation research has experienced a computational turn. In this paper, we present a linguistic and visual framework for the study of deliberative communication. The framework

manding and time-consuming resulting in a limited set of debate corpora. Moreover, the coding is often subjective making it subject to critical judgments of other researchers (King, 2009; Black et al., 2010; Dacombe, 2013). As a result, manual coding poses challenges with respect to both validity and reliability.

Only recently, the computational turn in deliberation research allows to analyze large quantities of

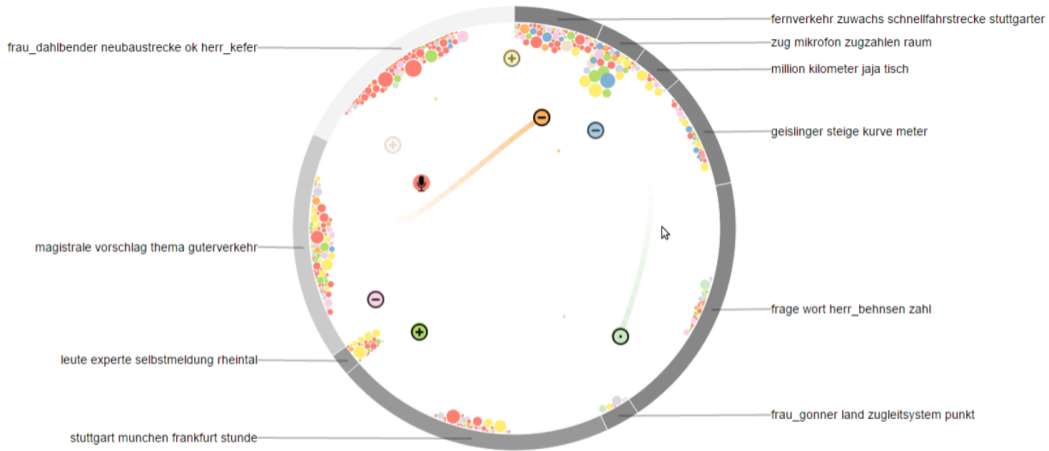


Figure 2: ConToVi Visualization

Back to the Big Picture

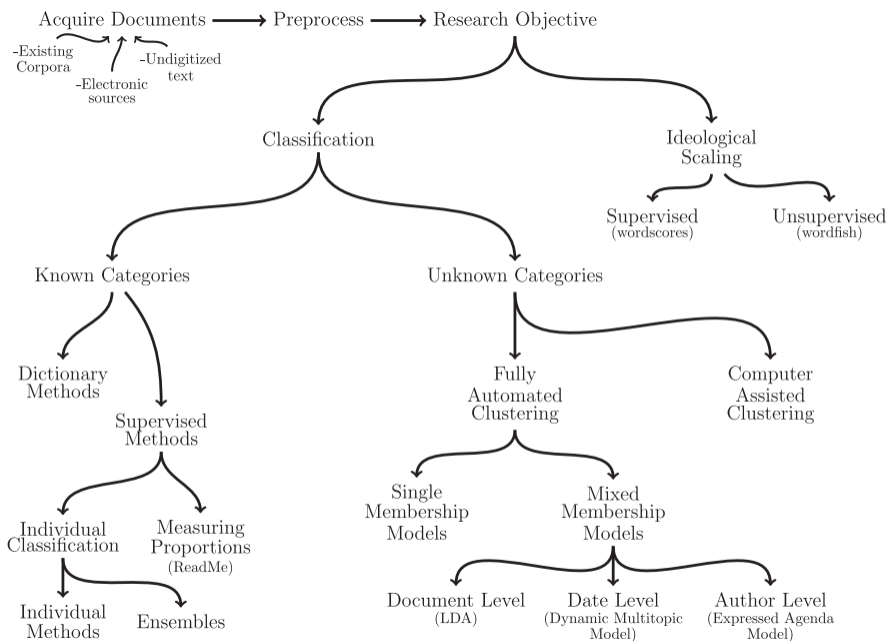


Fig. 1 An overview of text as data methods.

+

Beyond Money in Politics: Automatic Detection of Legislative Text Re-Use

Eugenia Giraudy

UC Berkeley & YouGov

eugenia.giraudy@gmail.com

Matthew Burgess

University of Michigan

mattburg@umich.edu

Julian Katz Samuels

University of Michigan

jkatzsam@umich.edu

Joe Walsh

University of Chicago

jtwalsh@uchicago.edu

Abstract

State legislatures introduce at least 45,000 bills each year. A large number of these bills are not drafted by legislators but by interest groups. However, existing approaches to detect the source of these bills are slow, biased, and incomplete.

This paper presents the Legislative Influence Detector (LID). LID uses the Smith-

drafted by interest groups. While state legislators normally adapt some of the language to meet their state's legal needs, such as references to existing state code, it is common practice to introduce bills that look the same as the document they borrowed from.

While legislative text re-use is a common phenomenon across states, Americans rely on watchdogs to find legislation of questionable origin, but with at least 45,000 bills introduced in state can-

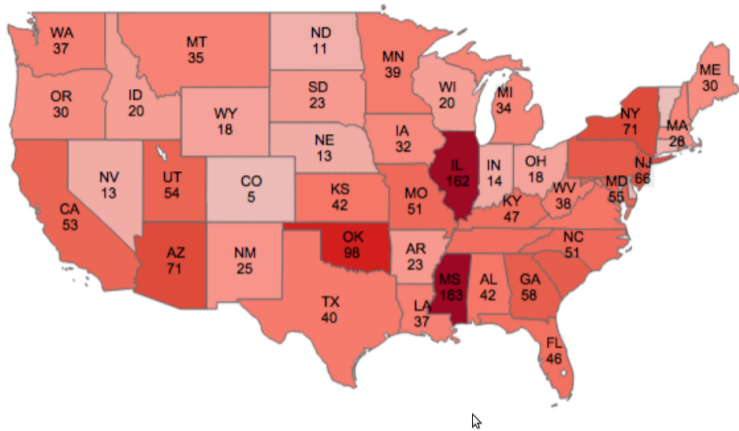


Figure 2: Introduced bills by state from ALEC model legislation

crease
ing to
cessib
nalists
dates.
every
LID e
view.
of pot
Sec
match
match
other s

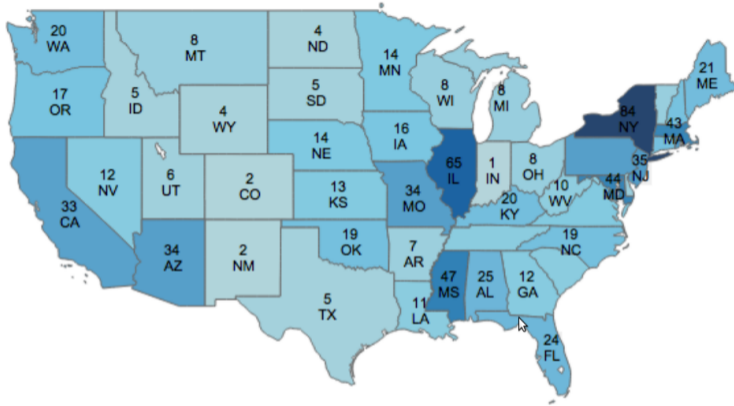


Figure 3: Introduced bills by state from ALICE model legislation

as m
lobby
lation
to the
This
At the
tential
own se
are de
feedba
that in
For ex
whose
a user

!