

QUANTITATIVE TEXT ANALYSIS

William Lowe

Hertie School of Governance

17th September 2020

TEXT AS DATA: THE APPROACH

Four broad approaches to studying text data

- Just read it and think a bit, e.g. op-eds, punditry, kremlinology, grand strategy, etc.
- Discourse Analysis
- Natural Language Processing (NLP)
- Text as Data (TADA)

the last two are, broadly, Computational Linguistics, but with a different focus

NOT DISCOURSE ANALYSIS

Although discourse analysis can be applied to all areas of research, it cannot be used with all kinds of theoretical framework. Crucially, it is not to be used as a method of analysis detached from its theoretical and methodological foundations. Each approach to discourse analysis that we present is not just a method for data analysis, but a theoretical and methodological whole - a complete package. [...] In discourse and analysis theory and method are intertwined and researchers must accept the basic philosophical premises in order to use discourse analysis as their method of empirical study.

(Jørgensen & Phillips, 2002)

Apparent differences are theoretical. The important difference for us is that

→ Discourse analysis *tightly couples* theory and measurement

Substantive theory \neq textual measurement...but they do have implications for one another

NOT (JUST) NLP

Overlapping NLP tasks

- Segmentation / tokenization: Locating words and sentences
- Part of Speech (POS) tagging: Associating grammatical roles with words (noun, verb, determiner, preposition, etc.)
- Parsing: grammatical structure from sentences

Distinctly NLP tasks

- Named Entity Recognition (NER): Identifying people, places, and things
- Information Extraction (IE): Extracting 'facts' (who did what to whom, when)

TEXT AS DATA: THE APPROACH

We are the measurement component for social science theory

- Theory provides the things to be measured
- Words and sometimes other things provide the data to measure them
- Language agnostic, behaviourist, structurally indifferent, shamelessly opportunistic
- obsessed with counting words

If Discourse analysis offers close reading, we will offer *distant reading*

Advantages

- Scales well
- Easy to integrate into existing models
- Can guide close reading later

TRANSCENDENTAL QUESTION

What are the *conditions for the possibility* for taking a TADA approach

In plainer language:

→ How could this possibly work?

BIG PICTURE

There is a *message* or *content* that cannot be directly observed, e.g.

- the topic of this lecture
- my position on some political issue
- the importance of defence issues to a some political party

and *behaviour*, including *linguistic behaviour*, e.g.

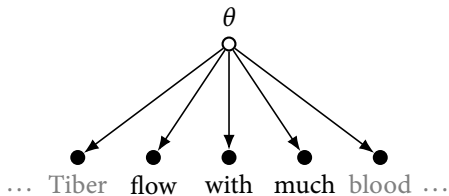
- yelling, writing, lecturing

which *can* be directly observed.

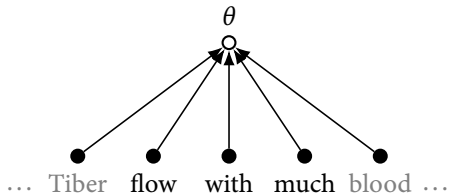
Although language can do things directly (Austin, 1962), we'll focus on the *expressed message* and the *words*...

COMMUNICATION

To *communicate* a message θ – to inform, persuade, demand, or threaten – a producer (the speaker or writer) *generates* words of different kinds in different quantities



To *understand* a message the consumer (the hearer, reader, coder) uses those words to *reconstruct* the message



COMMUNICATION

This process is

- stable (Grice, 1993; Searle, 1995)
- conventional (Lewis, 2011)
- disruptible (Riker et al., 1996)
- empirically underdetermined (Davidson, 1985; Quine, 1960)

How to model this without having to solve the problems of linguistics (psychology, politics) first?

Rely on:

- instrumentality
- reflexivity
- randomness

COMMUNICATION: INSTRUMENTALITY

Instrumentality from 'them': Language use is a form of action (Austin, 1962; Krebs & Dawkins, 1984; Wittgenstein, 1958)

Note the distinction between

X means Y

X is used to mean Y

Instrumentality from us:

- we aren't actually interested in words themselves; that's for linguists
- we aren't actually interested in what's in the head; that's for psychologists

Except as they help explain things we are interested in. Text is just data

COMMUNICATION: REFLEXIVITY

Politicians are often nice enough to talk as if they really do communicate this way

My theme here has, as it were, four heads. [...] The first is articulated by the word “opportunity” [...] the second is expressed by the word “choice” [...] the third theme is summed up by the word “strength” [and] my fourth theme is expressed well by the word “renewal”.

(Note however, these words occur 2, 7, 2, and 8 times in 4431 words)

COMMUNICATION: REFLEXIVITY

Politicians are often nice enough to talk as if they really do communicate this way

My theme here has, as it were, four heads. [...] The first is articulated by the word “opportunity” [...] the second is expressed by the word “choice” [...] the third theme is summed up by the word “strength” [and] my fourth theme is expressed well by the word “renewal”.

(Note however, these words occur 2, 7, 2, and 8 times in 4431 words)

A couple months ago we weren't expected to win this one, you know that, right? We weren't...Of course if you listen to the pundits, we weren't expected to win too much. And now we're winning, winning, winning the country – and soon the country is going to start winning, winning, winning.



COMMUNICATION AND COMPARABILITY

Quantitative text analysis works best when language usage is stable, conventionalized, and instrumental.

Implicitly, that means *institutional language*, e.g.

- courts
- legislatures
- op-eds
- financial reporting

Institution-specificity inevitably creates a *comparability* problem, e.g.

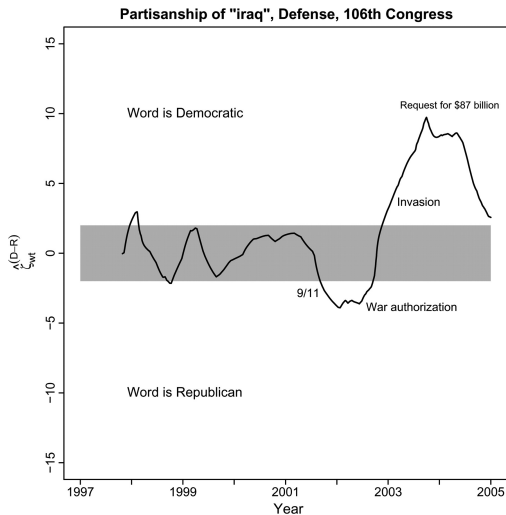
- upper vs lower chamber vs parliamentary hearings
- bureaucracy vs lobby groups (Klüver, 2009)
- European languages (Proksch et al., 2019)

INSTABILITY

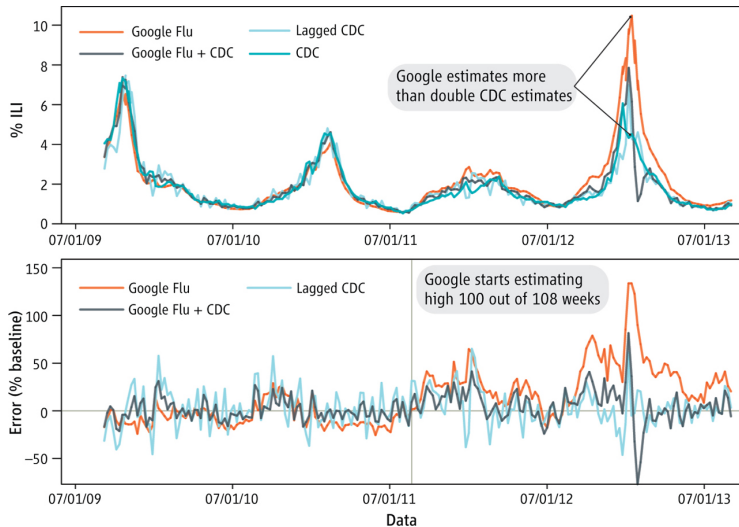
We are going to design instruments to measure θ and are going to assume that the $\theta \rightarrow W$ relationships are institutionally stable

What if they aren't?

RHETORICAL INSTABILITY



ALGORITHMIC INSTABILITY



REFLEXIVE SOLUTIONS

Sometimes these actors are happy to solve comparability problems for us, e.g.

- Lower court opinions (Corley et al., 2011) or Amicus briefs (Collins et al., 2015) *embedded in* Supreme Court opinions
- ALEC model bills *embedded in* state bills (Garrett & Jansa, 2015)

A perfect jobs for *text-reuse* algorithms...

COMMUNICATION: RANDOMNESS

Why randomness?

You almost never *say exactly the same words twice*, even when you haven't changed your mind about the message.

Hence words are the result of some kind of *sampling process*.

We model this process as random because we don't know or care about all the causes of variation (and because we're all secretly Bayesians)

Note: this is randomness *conditional on the institution*



WORDS AS DATA

What do we know about words as data?

They are *difficult*

- High dimensional
- Sparsely distributed (with skew)
- Not equally informative

DIFFICULT WORDS

Example: Conservative party 2017 manifesto compared to other parties over four elections:

- *High dimensional*. 3784 word types (adult native english speakers know 20-35,000)
- *Sparse*. Of 16083, word types in total, the Conservatives only used 3784
- *Skewed*. Of these 1731 words appeared exactly once and the most frequent word 1757 times

DIFFICULT WORDS

Example: Conservative party 2017 manifesto compared to other parties over four elections:

- *High dimensional*. 3784 word types (adult native english speakers know 20-35,000)
- *Sparse*. Of 16083, word types in total, the Conservatives only used 3784
- *Skewed*. Of these 1731 words appeared exactly once and the most frequent word 1757 times

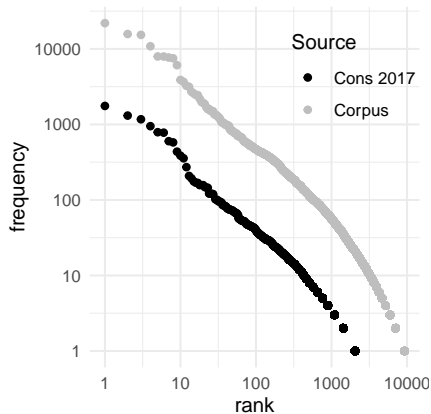
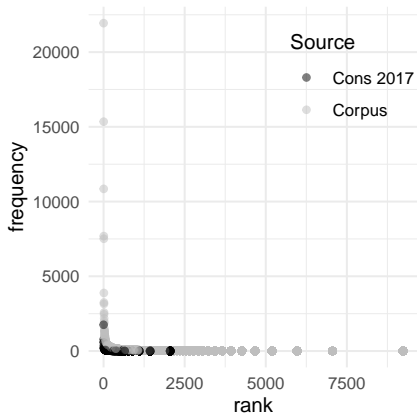
More generally: the Zipf-Mandelbrot law (Mandelbrot, 1966; Zipf, 1932)

$$F(W_i) \propto 1/\text{rank}(W_i)^\alpha$$

where $\text{rank}(\cdot)$ is the frequency *rank* of a word in the vocabulary and $\alpha \approx 1$

This is a Pareto distribution in disguise

DIFFICULT AT ALL SCALES



See Chater and Brown (1999) on scale invariance.

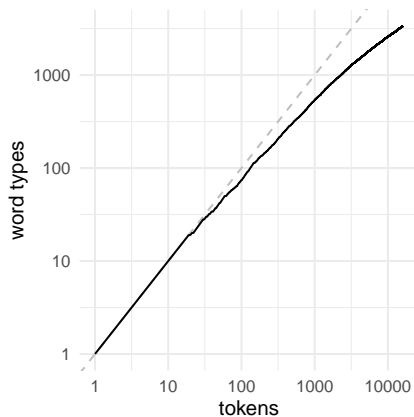
TYPES AND TOKENS

More generally: the Heaps-Herdan Law states that the number of word types appearing for the first time after n tokens is

$$D(n) = Kn^{\beta}$$

where K is between 10 and 100 and $\beta \approx 0.5$ for English.

(All the party manifestos shown here)



FREQUENCY AND INTERESTINGNESS

Frequency is inversely proportional to substantive interestingness

	Word	Freq.
1	the	21939
2	and	15747
3	to	15347
4	of	10850
5	we	7943
6	will	7930

Top 10

	Word	Freq.
16078	1.83	1
16079	2.20	1
16080	1.35	1
16081	33.34	1
16082	1.71	1
16083	rigation	1

Bottom ten

	Word	Freq.
20	people	1929
26	new	1507
27	government	1493
33	support	1212
34	work	1143
36	uk	1058

Top ten minus *stopwords*

DEALING WITH DIFFICULT WORDS

Removing stopwords, while standard in computer science, is not necessarily better...

Example:

- Standard collections contain, 'him', 'his', 'her' and 'she'.
- Words you'd want to keep when analyzing an abortion debates.

DEALING WITH DIFFICULT WORDS

For large amounts of text summaries are not enough.

We need a *model* to provide assumptions about

→ *equivalence*

→ *exchangeability*

Text as data started off making most use of equivalence, and ended up with increasingly sophisticated versions of exchangeability

Since ontogeny recapitulates phylogeny, let's walk through some standard text processing steps, asserting equivalences along the way...

PUNCTUATION INVARIANCE

As I look ahead I am filled with foreboding. Like the Roman I seem to see ‘the river Tiber flowing with much blood’...”

(Powell, 1968)

PUNCTUATION INVARIANCE

As I look ahead I am filled with foreboding. Like the Roman I seem to see ‘the river Tiber flowing with much blood’...”

(Powell, 1968)

index	token	index	token
1	as	1	like
2	i	2	the
3	look	3	roman
4	ahead	4	i
5	i	5	seem
6	am	6	to
7	...	7	...

LEXICAL UNIVOCALITY

type	count
as	1
i	2
look	1
ahead	1
am	1
...	...

token	count
like	1
the	1
roman	1
i	1
seem	1
to	1
...	...

ORDER INVARIANCE

		unit	
		'doc' 1	'doc' 2
type	ahead	1	0
	am	1	0
	as	1	0
	i	2	1
	like	0	1
	look	1	0
	roman	0	1
	seem	0	1
	the	0	1
	to	0	1

COUNT DATA

We have turned a corpus into a *contingency table*.

→ Or a term-document / document-term / document-feature matrix, in the lingo

Everything you learned in your categorical data analysis course applies

→ except that the variables of interest: θ are *not observed*

COUNT DATA

We have turned a corpus into a *contingency table*.

→ Or a term-document / document-term / document-feature matrix, in the lingo

	ahead	am	i	like	look		
doc 1	1	1	2	0	1	...	θ_{doc1}
doc 2	0	0	1	1	0	...	θ_{doc2}
	β_{ahead}	β_{am}	β_i	β_{like}	β_{look}		

Everything you learned in your categorical data analysis course applies

→ except that the variables of interest: θ are *not observed*

What are we going to assume about the cell contents?

STATISTICAL ASSUMPTIONS ABOUT WORDS

Word counts/rates are conditionally *Poisson*:

$$W_j \sim \text{Poisson}(\lambda_j)$$

Curiously

$$\mathbb{E}[W] = \text{Var}[W] = \lambda$$

Rate models are naturally *multiplicative*.

→ Rates increase / decrease by X%

Model assumptions are will turn on how λ is related to θ

STATISTICAL ASSUMPTIONS ABOUT WORDS

That means that for fixed document lengths, counts are conditionally *Multinomial*:

$$W_{i1} \dots W_{iV} \sim \text{Mult}(W_{i1} \dots W_{iV} \mid \pi_1 \dots \pi_V, N_i)$$

Here

$$E[W] = N\pi$$

and

$$\text{Cov}[W_i, W_j] = -N\pi_i\pi_j$$

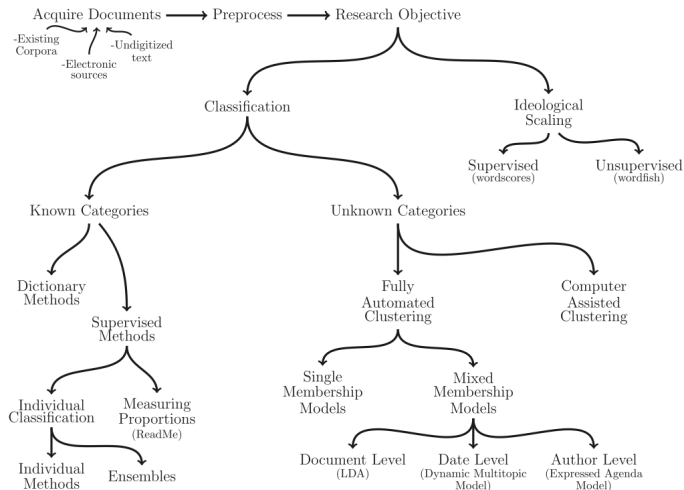
Negative covariance is due to the ‘budget constraint’

MODELLING DECISIONS

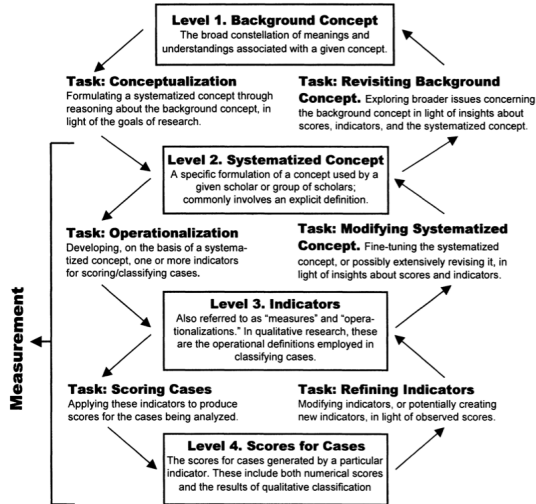
For each research problem involving content analysis we need to ask:

- What *structure* does θ have
- What counts as a *word*?
- What counts as a *document*?
- What is *observed*, what is *assumed*, and what is *inferred*?
- What is the *relationship* between θ and the words? The model

THE SPACE OF MODELS



IN THE GENERAL MEASUREMENT PROBLEM





EXCHANGEABILITY AND THE ‘BAG OF WORDS’



„Relief for upswing“

Patrick Döring

Reasonable relief of the working middle class and debt reduction do not exclude each other but are complementary. That proved to be true during the past four years. We oppose higher taxes on



„Stable currency for more prosperity“

Otto Fricke Jörg-Uwe Hahn

Safe money is a cornerstone of all free and fair social and economic orders. Inflation means destruction of savings and devaluation of what people have achieved in their lives.

THE MARGINAL CHALLENGE

The original:

Reasonable relief of the working middle class and debt reduction do not exclude each other but are complementary. That proved to be true during the past four years. We oppose higher taxes on citizens and businesses. They prevent growth and kill jobs, thus putting at risk the very existence of countless workers and their families.

Given the *marginal distribution of word types* a.k.a. the ‘bag of words’ how different could the meaning of the resulting document be?

SLIGHTLY DIFFERENT...

A possible reconstruction

We oppose higher taxes because they prevent growth and kill jobs, risking the very existence of countless families and businesses. Debt reduction is complementary to reasonable relief of the working middle class; they do not exclude each other. That proved to be true in the last four years.

But much the same sense

LOOSER CONSTRAINTS

We can make this a more semantic challenge by only demanding the substantively interesting word margins are maintained.

*businesses reasonable existence countless reduction families prevent risking working exclude
oppose higher growth relief middle proved taxes class years kill jobs debt true last four com-
plementary*

Removing stopwords mostly just removes *grammatical* constraints

To the extent this set of words characterizes what the FDP wanted to express in their platform, the ‘bag of words’ assumption is reasonable.

BUT NEGATION!

If we can add grammatical functors at will, could we make the opposite meaning by negating everything?

In principle (and in practice for some discursive forms) yes, yes we could.

And the bag of words assumption would fail

BUT NEGATION!

If we can add grammatical functors at will, could we make the opposite meaning by negating everything?

In principle (and in practice for some discursive forms) yes, yes we could.

And the bag of words assumption would fail

An interesting *empirical fact* about political discourse is that actors do not tend to disagree by negation but by redirection or diversion.

- Simple version: You talk about the environment, I talk about economic growth
- Sophisticated version: The 'heresthetic' (Riker et al., 1996)

EXCHANGEABILITY AND THE ‘BAG OF WORDS’

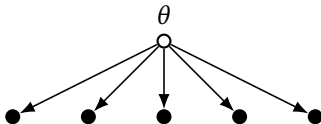
Technically a sequence of random variables is *exchangeable* if their *joint probability distribution* is invariant to reordering them

→ a ‘weaker’ relationship than identical and independently distributed (iid)

de Finetti (2008) showed that if we assume $\{W_i\}$ are exchangeable, then this is equivalent to assuming that they are (loosely speaking) generated by a distribution parameterized by some θ :

$$P(W_1 \dots W_V) = \int \prod_v P(W_v | \theta) P(\theta) d\theta$$

This is ‘the Representation Theorem’ and motivates our measurement model assumptions earlier



REFERENCES

- Austin, J. L. (1962). 'How to do things with words'. Clarendon Press.
- Chater, N. & Brown, G. D. A. (1999). 'Scale-invariance as a unifying psychological principle.' *Cognition*, 69(3), B17–24.
- Collins, P. M., Corley, P. C. & Hamner, J. (2015). 'The influence of amicus curiae briefs on u.s. supreme court opinion content: The influence of amicus curiae'. *Law & Society Review*, 49(4), 917–944.
- Corley, P. C., Collins, P. M. & Calvin, B. (2011). 'Lower court influence on u.s. supreme court opinion content'. *The Journal of Politics*, 73(1), 31–44.
- Davidson, D. (1985). 'Inquiries into truth and interpretation'. Clarendon Press.
- de Finetti, B. (2008). 'Philosophical lectures on probability' (A. Mura, Ed.; H. Hosni, Trans.). Springer.
- Garrett, K. N. & Jansa, J. M. (2015). 'Interest group influence in policy diffusion networks'. *State Politics & Policy Quarterly*, 15(3), 387–417.

REFERENCES

- Grice, P. (1993). 'Studies in the way of words' (3. print). Harvard Univ. Press.
- Jørgensen, M. & Phillips, L. (2002). 'Discourse analysis as theory and method'. Sage Publications.
- Klüver, H. (2009). 'Measuring interest group influence using quantitative text analysis'. *European Union Politics*, 10(4), 535–549.
- Krebs, J. & Dawkins, R. (1984). Animal signals: Mind-reading and manipulation. In J. Krebs & N. B. Davies (Eds.), *Behavioural ecology: An evolutionary approach* (2nd ed., pp. 380–402). Blackwell Science.
- Lazer, D., Kennedy, R., King, G. & Vespignani, A. (2014). 'The parable of google flu: Traps in big data analysis'. *Science*, 343(6176), 1203–1205.
- Lewis, D. K. (2011). 'Convention: A philosophical study' (Nachdr.). Blackwell.
- Mandelbrot, B. (1966). Information theory and psycholinguistics: A theory of word frequencies. In P. Lazarsfeld & N. Henry (Eds.), *Readings in mathematical social science*. MIT Press.

REFERENCES

- Monroe, B. L., Colaresi, M. & Quinn, K. M. (2008). 'Fightin' words: Lexical feature selection and evaluation for identifying the content of political conflict'. *Political Analysis*, 16(4), 372–403.
- Powell, E. (1968). 'Speech delivered to the conservative association'.
- Proksch, S.-O., Lowe, W., Wäckerle, J. & Soroka, S. (2019). 'Multilingual sentiment analysis: A new approach to measuring conflict in legislative speeches'. *Legislative Studies Quarterly*, 44(1), 97–131.
- Quine, W. v. O. (1960). 'Word and object'. MIT Press.
- Riker, W. H., Calvert, R. L., Mueller, J. E. & Wilson, R. K. (1996). 'The strategy of rhetoric: Campaigning for the american constitution'. Yale University Press.
- Searle, J. R. (1995). 'The construction of social reality'. Free Press.
- Wittgenstein, L. (1958). 'Philosophical investigations' (G. E. M. Anscombe, Trans.). Blackwell.

REFERENCES

Zipf, G. K. (1932). 'Selected studies of the principle of relative frequency in language'. Oxford University Press.