

Text Mining

Itamar Caspi

May 26, 2019 (updated: 2019-05-26)

Outline

1. Representing Text as Data
2. Text Regressions
3. Dictionary-based Methods
4. Topic Modeling
5. Hands-on

Representing Text as Data

Where to Start

A great introduction to this topic, along with many empirical examples from the social sciences appear in Gentzkow, Kelly, and Taddy's **"Text as Data"** (JEL forthcoming).

Basic notation

Definitions:

- A *corpus* is a collection of D documents (emails, tweets, speeches, articles, etc.)
- A *vocabulary* is a complete list of unique words that appear in the corpus.
- \mathbf{X} is a numerical array representation of text. Rows correspond to documents $i = 1, \dots, D$ and columns to words $j = 1, \dots, N$.
- \mathbf{Y} is a vector of predicted outcome (e.g., spam/ham, trump/not trump, etc.), one outcome per document.
- \mathbf{F} is a low-dimensional representation of \mathbf{X} .

Document term matrix (DTM)

In most applications raw text is represented as a numerical array \mathbf{X} where the elements of the array, X_{ij} , are counts of words (or more generally, *tokens*. More on that later.)

Illustration: Spam vs. Ham

Consider the task of spam detection:

ham	Go until jurong point, crazy.. Available only in bugis n great world la e buffet... Cine there got amore wat...
ham	Ok lar... Joking wif u oni...
spam	Free entry in 2 a wkly comp to win FA Cup final tkts 21st May 2005. Text FA to 87121 to receive entry question(std txt rate)T&C's apply 08452810075over18's
ham	U dun say so early hor... U c already then say...
ham	Nah I don't think he goes to usf, he lives around here though
spam	FreeMsg Hey there darling it's been 3 week's now and no word back! I'd like some fun you up for it still? Tb ok! XxX std chgs to send, å£1.50 to rcv
ham	Even my brother is not like to speak with me. They treat me like aids patent.
ham	As per your request 'Melle Melle (Oru Minnaminunginte Nurungu Vettam)' has been set as your callertune for all Callers. Press *9 to copy your friends Callertune
spam	WINNER!! As a valued network customer you have been selected to receivea å£900 prize reward! To claim call 09061701461. Claim code KL341. Valid 12 hours only.
spam	Had your mobile 11 months or more? U R entitled to Update to the latest colour mobiles with camera for Free! Call The Mobile Update Co FREE on 08002986030

In this case

- Documents are emails.
- Vocabulary includes words that appear in *each and every* emails.

NOTE: Spam detection is clearly a supervised learning task where $Y_i = \{\text{spam}, \text{ham}\}$.

Transforming a corpus to a DTM

Consider the following corpus (\$M=3\$):

```
txt <- c(doc1 = "Shipment of gold damaged in a fire",
        doc2 = "Delivery of silver arrived in a silver truck",
        doc3 = "Shipment of gold arrived in a truck" )

txt %>% quantda::dfm() # transform text as a document term matrix
```

```
## Document-feature matrix of: 3 documents, 11 features (36.4% sparse).
## 3 x 11 sparse Matrix of class "dfm"
##      features
## docs shipment of gold damaged in a fire delivery silver arrived truck
## doc1      1  1  1      1  1  1      1      0      0      0      0
## doc2      0  1  0      0  1  1      0      1      2      1      1
## doc3      1  1  1      0  1  1      0      0      0      1      1
```

- `nrow(dfm_txt)` returns the number of documents in the corpus.
- `ncol(dfm_txt)` returns the number of words in the vocabulary.
- `colnames(dfm_txt)` returns the number of words in the vocabulary

Does every words matter? ヽ(ツ)ノ/

We can significantly reduce the dimension of \mathbf{X} by

- filtering out very common ("stop words") and uncommon words.
- dropping numbers and punctuation.
- stemming, i.e., replacing words by their root. (*economi* instead of *economics*, *economists*, *economy*)
- convert to lower case

WARNING: Use text prepossessing steps with care. These steps should be application specific.

Example

Here, we remove stop words, punctuation, numbers, and stem words:

```
txt <- c(doc1 = "Shipment of gold damaged in a fire",
        doc2 = "Delivery of silver arrived in a silver truck",
        doc3 = "Shipment of gold arrived in a truck" )

txt %>% dfm(remove = stopwords("english"),
           remove_punct = TRUE,
           remove_numbers = TRUE,
           stem = TRUE)
```

```
## Document-feature matrix of: 3 documents, 8 features (50.0% sparse).
```

```
## 3 x 8 sparse Matrix of class "dfm"
```

```
##      features
```

## docs	shipment	gold	damag	fire	deliveri	silver	arriv	truck
## doc1	1	1	1	1	0	0	0	0
## doc2	0	0	0	0	1	2	1	1
## doc3	1	1	0	0	0	0	1	1

Words are generally not independent

Sometimes we might care about multiword expressions, e.g., "not guilty", "labor market", etc.

We can define tokens (the basic element of text) as n -gram - a sequence of n words from a given sample of text.

NOTE: Using n -gram with $n > 2$ is typically impractical due to the rapid growth in the dimensions of \mathbf{X} .

Example

Here is our sample text (just 2 "documents" in this example), where tokens are defined as *bi-gram* (a sequence of two words):

```
txt <- c(doc1 = "Shipment of gold damaged in a fire",  
        doc2 = "Delivery of silver arrived in a silver truck")
```

```
txt %>% dfm(ngrams = 2L)
```

```
## Document-feature matrix of: 2 documents, 12 features (45.8% sparse).
```

```
## 2 x 12 sparse Matrix of class "dfm"
```

```
##      features
```

```
## docs  shipment_of of_gold delivery_of gold_damaged damaged_in of_silver
```

```
## doc1      1      1      0      1      1      0
```

```
## doc2      0      0      1      0      0      1
```

```
##      features
```

```
## docs  in_a a_fire silver_arrived arrived_in a_silver silver_truck
```

```
## doc1  1      1      0      0      0      0
```

```
## doc2  1      0      1      1      1      1
```

The (social science) textmining play book

1. Collect text and generate a corpus.
2. Represent corpus as a DTM \mathbf{X} .
3. Then, proceed according to one of the following steps:
 - Use \mathbf{X} to predict an outcome \mathbf{Y} using high dimensional methods (e.g., lasso, Ridge, etc.). In some cases, proceed with $\hat{\mathbf{Y}}$ to subsequent analysis.
 - Apply dimensionality reduction techniques (dictionary, PCA, LDA, etc.) to \mathbf{X} and proceed with the output \mathbf{F} to subsequent analysis.

Remember:

"Text information is usually best as part of a larger system. Use text data to fill in the cracks around what you know. Don't ignore good variables with stronger signal than text!" (Matt Taddy)

Text Regression

Sounds familiar...

- We are interested in predicting some Y using \mathbf{X} .
- Clearly, with text as data, we are facing the high-dimensionality problem. \mathbf{X} has $M \times N$ elements.
- Classical methods such as OLS won't do the trick \Rightarrow need the ML touch.
- An obvious choice would be penalized linear/non-linear regression (e.g. Lasso, ridge, etc.). Other methods such as random forest can work too.

EXAMPLE: Lasso text regression `glmnet(Y, X)` where

$$\hat{\beta} = \underset{\beta \in \mathbb{R}^N}{\operatorname{argmin}} \sum_{i=1}^N (Y_i - X_i \beta)^2 + \lambda \|\beta\|_1$$

Can be easily extended to binary / categorical Y , e.g. `glmnet(X, Y, family = "binomial")`

Practical notes about penalized text regression

- Typically, DTM entries count the number of times word i appears in document d . This provide "intuitive" interpretation for regression coefficients.
-
- Nevertheless, as usual, beware of giving causal interpretation to the Lasso's coefficients. (Recall the irrepresentability condition.)

Dictionary-based Methods

Reducing dimensionality using a dictionary

- Dictionary-based methods provide a low-dimensional representation of high-dimensional text.
- Essentially, Think of F as an unobserved characteristic of the text that we are trying to estimate. A Dictionary-based methods provides a mapping from \mathbf{X} onto F :

$$g : \mathbf{X} \rightarrow F$$

Example: Sentiment analysis

- A prominent example of dictionary-based methods is sentiment analysis
- The latent factor we are trying to estimate is the writer's attitude towards the topic in question.
- The most common approach is based on prespecified dictionaries that classify word according to some predefined sentiment class (e.g. "positive", "negative", and "neutral".)
- Typically, the sentiment *score* of each document is a function of the relative frequencies of positive, negative, neutral, etc., words.

REMARK: Sentiment analysis can be supervised as well. E.g., the availability of labeled movie reviews (1-5 stars) can be used to train a model and use its predictions to classify unlabeled reviews.

Example: Loughran and McDonald financial sentiment dictionary

A random list of words from the Loughran and McDonald (2011) financial sentiment dictionary (positive/negative/litigious/uncertainty/constraining):

```
library(tidytext)
sample_n(get_sentiments("loughran"),8)
```

```
## # A tibble: 8 x 2
##   word          sentiment
##   <chr>         <chr>
## 1 losing        negative
## 2 rumors        uncertainty
## 3 unsustainable negative
## 4 spectacularly positive
## 5 gain          positive
## 6 indefiniteness uncertainty
## 7 allegations   negative
## 8 evades        negative
```

Topic Modeling

Topic models

- Topic models extend unsupervised learning methods to text data.
- Topic modeling classifies documents and words to latent topics and is often followed by more conventional empirical methods.
- The workhorse of topic modeling is the Latent Dirichlet Allocation model (Blei, Ng, and Jordan, 2003), or LDA for short.

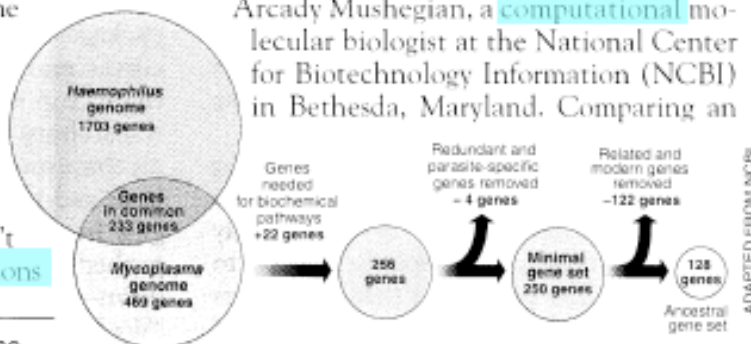
Intuition behind LDA

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

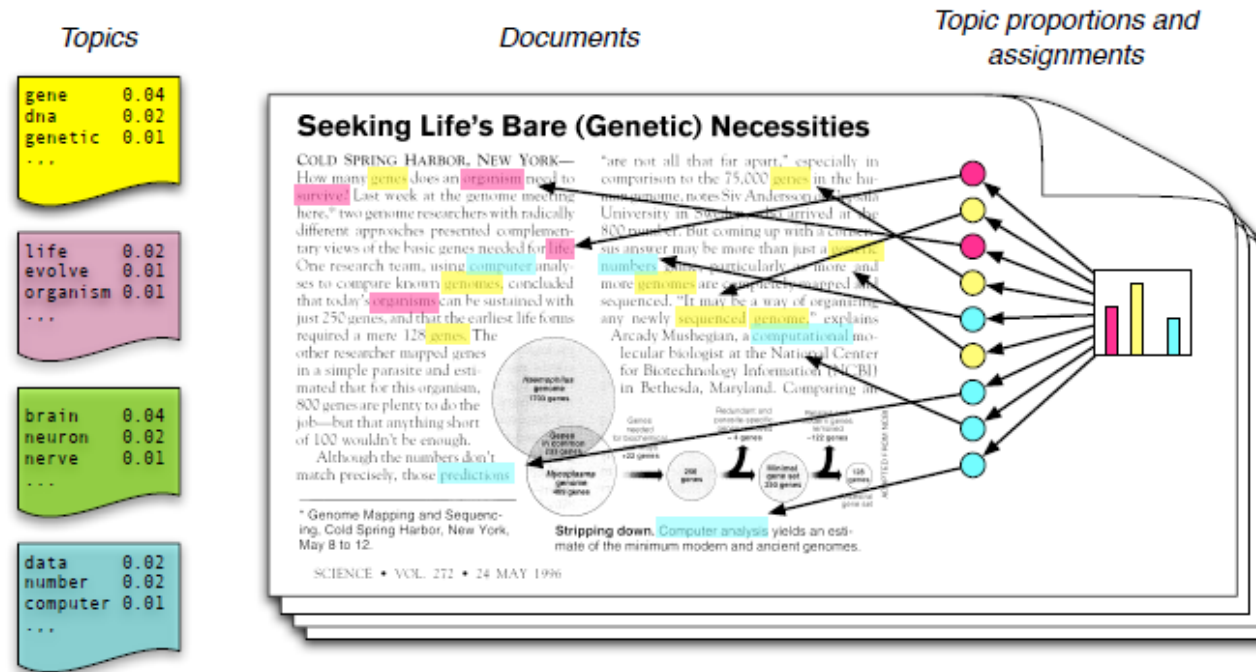
"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a genetic numbers game, particularly as more and more genomes are completely 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



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

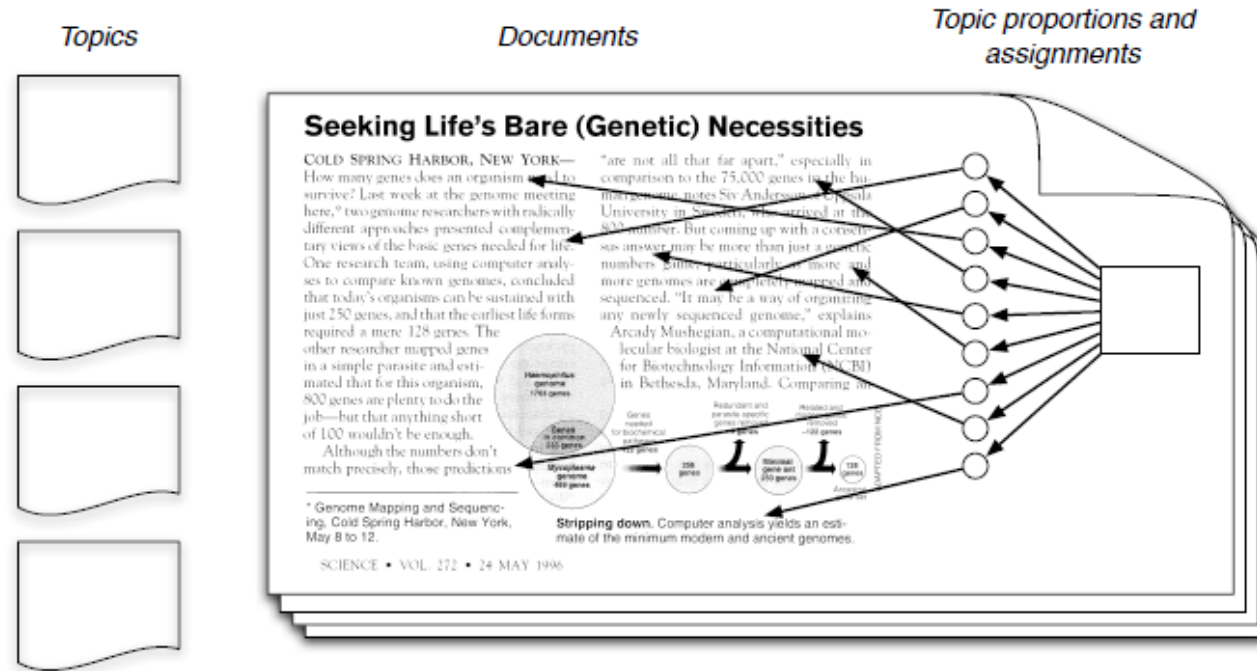
^{*} Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

Intuition behind LDA



- A topic is a distribution over *all* the words included in a *fixed* vocabulary.
- A word can have non-zero (yet different) probability multiple topics (e.g., bank)
- Each document is a mixture of topics
- Each word is drawn from one of the topics.

Intuition behind LDA



QUESTION: How realistic is the LDA setup? Does it matter? What's our goal here anyway?

Notation

- A *vocabulary* is a collection of words represented by the vector $\{1, \dots, V\}$
- Each *word* is represented by a unit vector $\boldsymbol{\delta}_v = (0, \dots, v, \dots, 0)'$
- A *document* is a sequence of N words denoted by $\mathbf{w} = (w_1, \dots, w_N)$.
- A *corpus* is a collection of M documents denoted by $\mathcal{D} = (\mathbf{w}_1, \dots, \mathbf{w}_M)$.

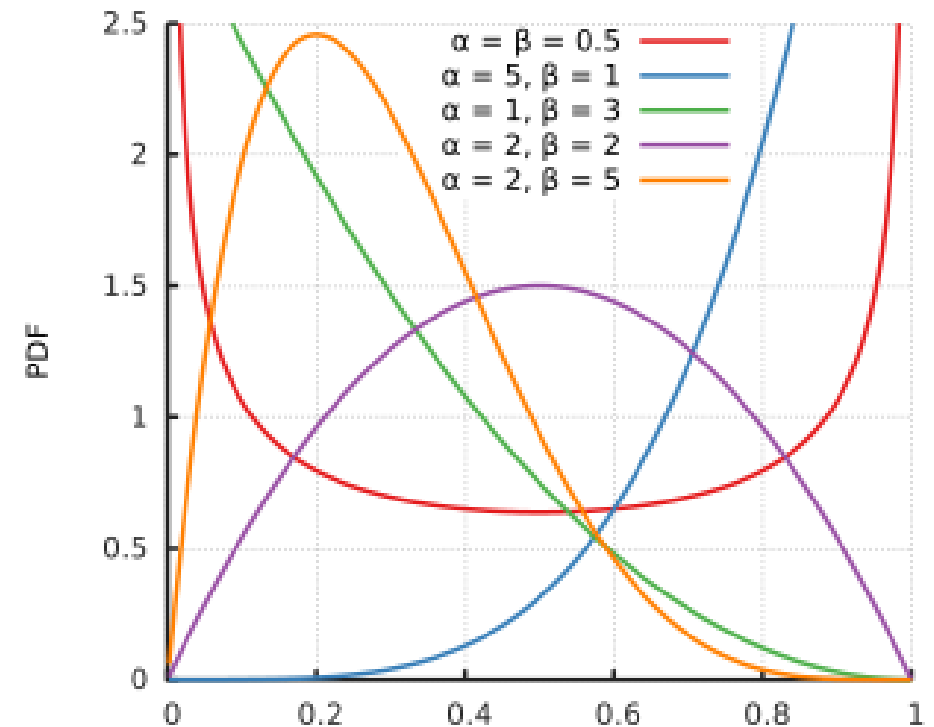
Prerequisite: The Beta distribution

The PDF for the Beta distribution, denoted as $B(\alpha, \beta)$ is

$$p(\theta|\alpha, \beta) \propto \theta^{\alpha-1}(1 - \theta)^{\beta-1}$$

for $\theta \in [0, 1]$ and $\alpha, \beta > 0$.

Because its properties, the Beta distribution is useful as a prior for probabilities.



The Dirichlet distribution

The Dirichlet distribution, denoted as $\text{Dir}(\boldsymbol{\alpha})$ is a multivariate generalization of the Beta distribution.

Let $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_K) \sim \text{Dir}(\boldsymbol{\alpha})$.

The PDF for a K -dimensional Dirichlet distribution is

$$p(\boldsymbol{\theta}|\boldsymbol{\alpha}) \propto \prod_{i=1}^K \theta_i^{\alpha_i-1}$$

where $K \geq 2$ is the number of categories, $\alpha_i > 0$ and $\theta_i \in (0, 1)$ for all i and $\sum_{i=1}^K \theta_i = 1$.

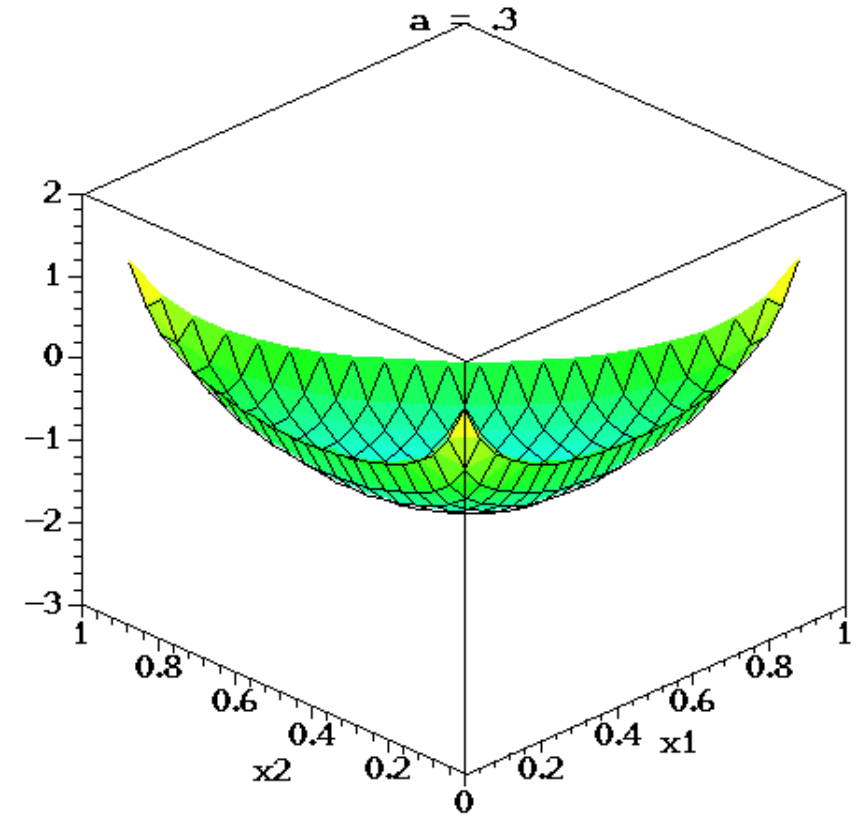
REMARK: The parameter $\boldsymbol{\alpha}$ controls the sparsity of $\boldsymbol{\theta}$

BOTTOM LINE: Vectors drawn from a Dirichlet distribution represent probabilities.

On the right:

The change in the density function ($K = 3$) as the vector α changes from $\alpha = (0.3, 0.3, 0.3)$ to $(2.0, 2.0, 2.0)$, while keeping $\alpha_1 = \alpha_2 = \alpha_3$.

REMARK: Placing $\alpha = (1, 1, 1)$ results in a uniform distribution over the simplex.



By [Initial version](#) by [Panos Ipeirotis](#), later modified by [Love Sun and Dreams - \[1\]](#), [CC BY 3.0](#), [Link](#)

The data generating process behind LDA

Assumption: The number of topics K and the size of the vocabulary V are fixed.

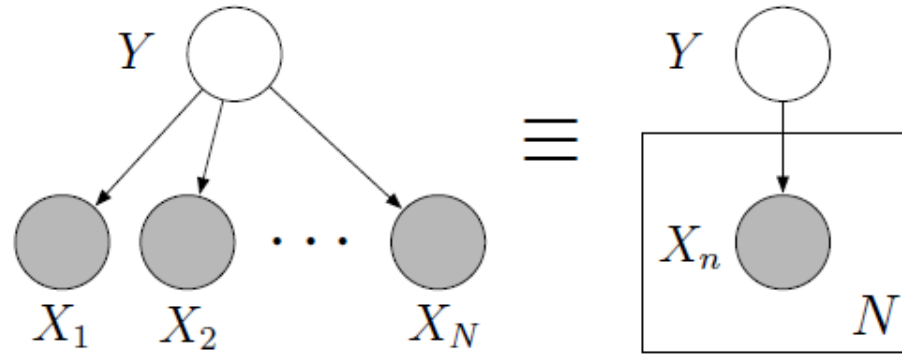
The DGP:

For each document $d = 1, \dots, \mathcal{D}$

1. Choose topic proportions $\theta_d \sim \text{Dir}(\alpha)$
2. For each word $n = 1, \dots, N$
 - 2.1. Choose a topic assignment $Z_{dn} \sim \text{Mult}(\theta_d)$.
 - 2.2. Choose a word $W_{dn} \sim \text{Mult}(\beta_{z_{dn}})$.

REMARK: Note the "factor model" aspects of LDA, where topics are factors and word probabilities are loadings, and both affect the probability of choosing a word.

Aside: Plate notation

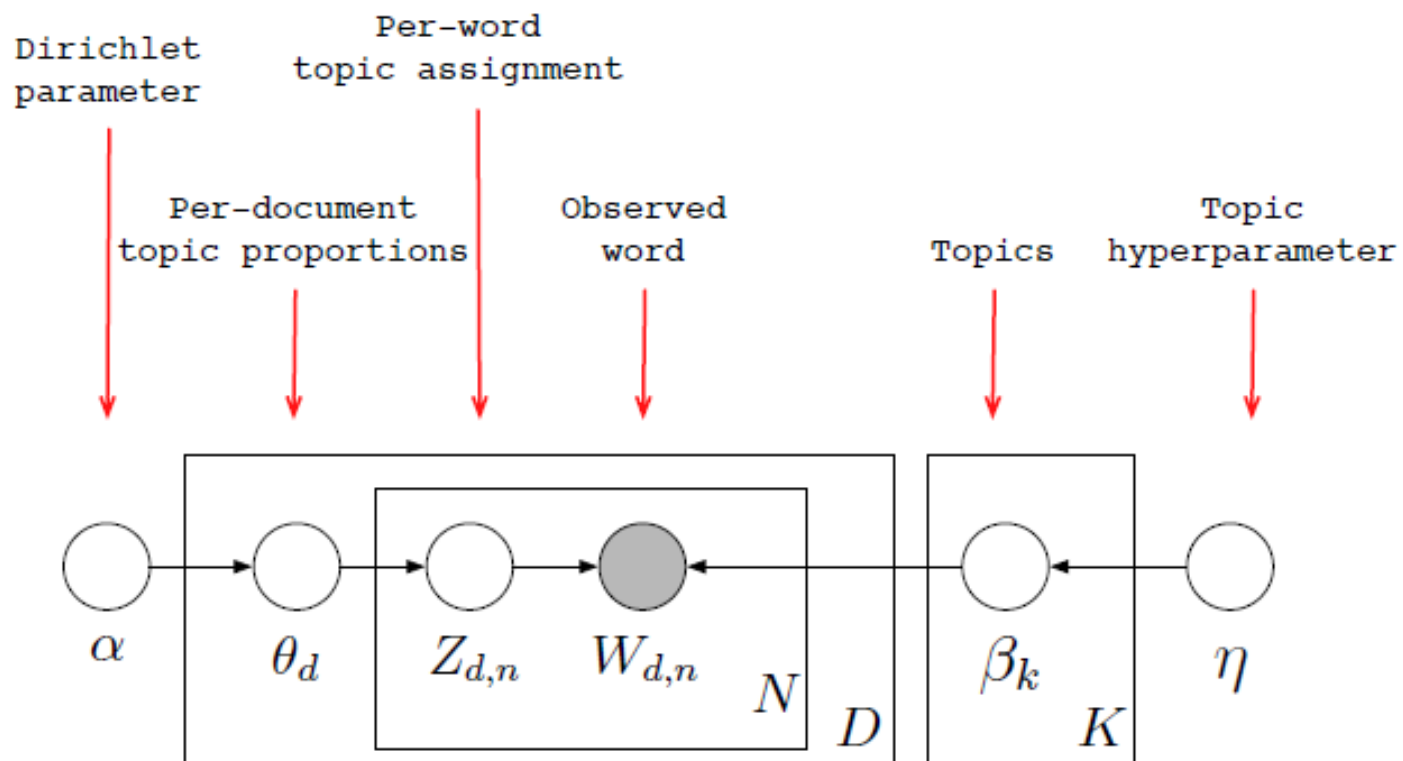


- Each *node* is a random variable
- *Shaded* nodes are observables
- *Edges* denote dependence
- *plates* denote replicated structures

The above graph corresponds to

$$p(y, x_1, \dots, x_N) = p(y) \prod_{n=1}^N p(x_n | y)$$

LDA in plate notation



Source: http://videlectures.net/mlss09uk_blei_tm/#.

Aside: Conjugate priors

The Dirichlet distribution is a conjugate prior for the Multinomial.

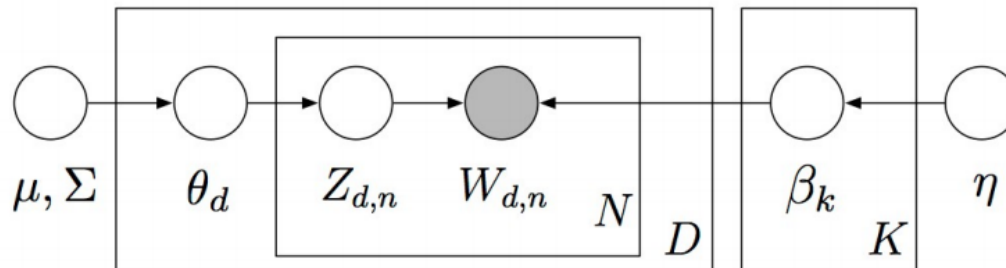
Let $n(Z_i)$ denote the count of topic i .

$$\boldsymbol{\theta} | Z_1, \dots, Z_N \sim \text{Dir}(\boldsymbol{\alpha} + n(Z_1, \dots, Z_N))$$

i.e., as the number of times we see topic i increases, our posterior becomes "peakier" at its i^{th} component.

Extension #1: Correlated topic models (Lafferty and Blei, 2005)

- LDA assumes that topics independently cooccur in documents.
- This is clearly wrong.
- For example, a document about *economics* is more likely to also be about *politics* than it is to be about *cooking*.
- Lafferty and Blei relax independence by drawing topic proportions from a logistic normal, which allows correlations between topic proportions:



where μ and Σ are priors for the logistic normal distribution.

Extension #2: Dynamic LDA (Blei and Lafferty, 2006)

Dynamic topic modeling takes into account the ordering of the documents and gives a richer posterior topical structure than LDA

In dynamic topic modeling, a topic is a *sequence* of distributions over words. Topics evolve systematically over time. In particular, the vector of parameters for topic k in period t evolves with a Gaussian noise:

$$\beta_{t,k} | \beta_{t-1,k} \sim \mathcal{N}(\beta_{t-1,k}, \sigma^2 I).$$

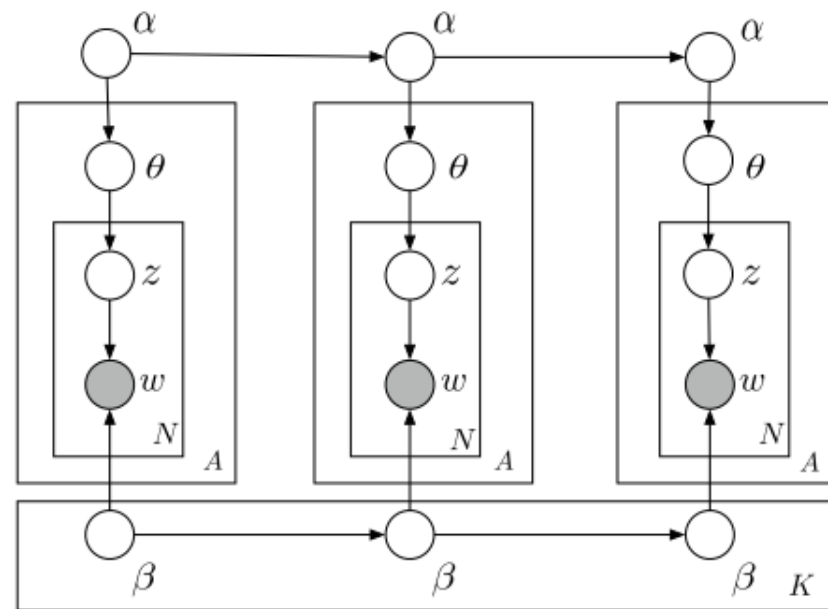
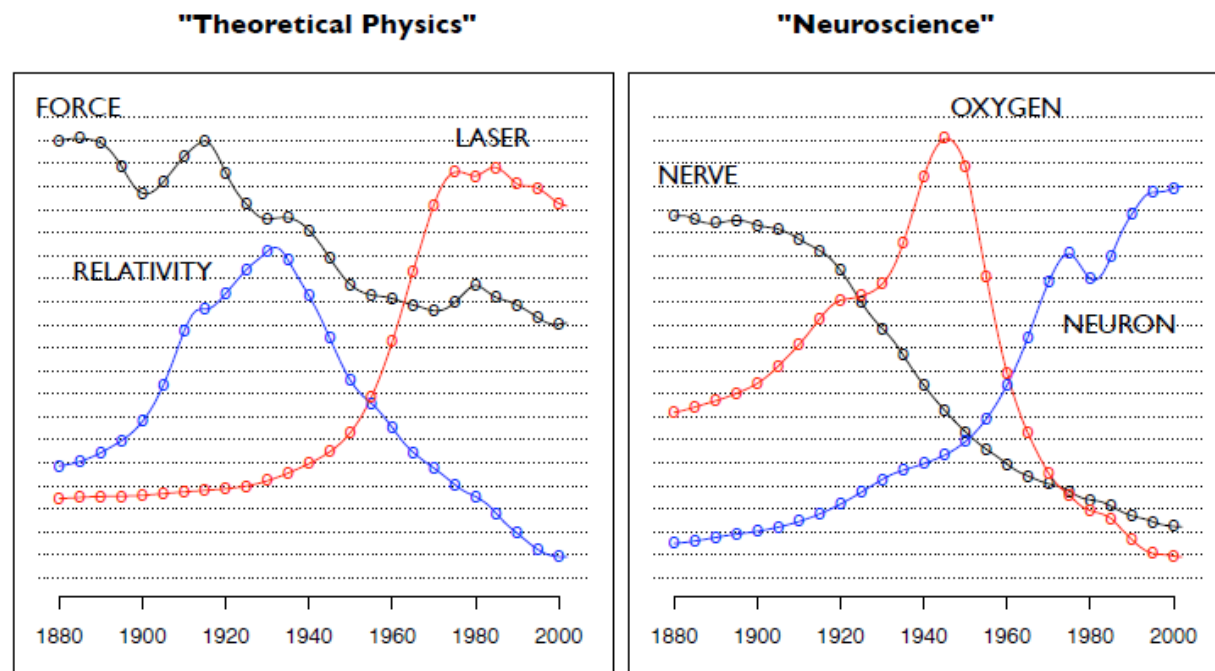


Figure 1. Graphical representation of a dynamic topic model (for three time slices). Each topic's natural parameters $\beta_{t,k}$ evolve over time, together with the mean parameters α_t of the logistic normal distribution for the topic proportions.

Dynamic LDA: Science 1881-1999

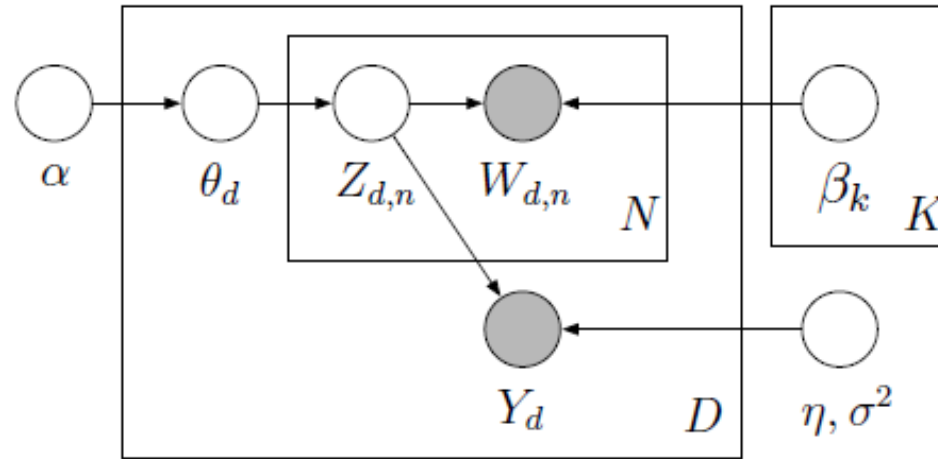
The posterior estimate of the frequency as a function of year of several words from the two topics: "Theoretical Physics" and "Neuroscience":



Source: Blei and Lafferty (2006).

Extension #3: Supervised Topic Model (McAuliffe and Blei, 2008)

add an extra connection between Z_{dn} to some observable attribute Y_d :



Source: McAuliffe and Blei (2008).

Structural Topic Models (Roberts, Stewart, and Tingley)

About the Structural Topic Model (STM):

"The Structural Topic Model is a general framework for topic modeling with document-level covariate information. The covariates can improve inference and qualitative interpretability and are allowed to affect topical prevalence, topical content or both."

In STM, topics are drawn from the following logistic normal distribution,

$$\boldsymbol{\theta}_d | \boldsymbol{X}_d \boldsymbol{\gamma}, \boldsymbol{\Sigma} \sim \text{LogisticNormal} (\boldsymbol{\mu} = \boldsymbol{X}_d \boldsymbol{\gamma}, \boldsymbol{\Sigma})$$

where \boldsymbol{X}_d is a vector of observed document covariates.

REMARK: In the case of no covariates, the STM reduces to a (fast) implementation of the Correlated Topic Model (Blei and Lafferty, 2007).

stm: R package for structural topic models

Roberts, Stewart, and Tingley (JSS, 2014)

About the stm R package:

"The software package implements the estimation algorithms for the model and also includes tools for every stage of a standard workflow from reading in and processing raw text through making publication quality figures."

The package is available on CRAN and can be installed using:

```
install.packages("stm")
```

To get started, see the [vignette](#) which includes several example analyses.

Applying topic models to measure the effect of transparency

Hansen, McMahon, and Prat (QJE 2017) study the effect of increasing Federal Open Market Committee (FOMC) transparency on debate during FOMC meetings.

- FOMC meetings have been tape recorded since the 1970s to prepare minutes.
 - Committee members believed that these tapes were erased afterward.
 - In October 1993, Fed chair Alan Greenspan ,discovered and revealed that before being erased the tapes had been transcribed and stored in archives all along.
 - Following Greenspan's revelation The Fed agreed to publish all past transcripts and extended that policy to cover all future transcripts with a five-year lag.
 - This gives Hansen et al. access to periods both when policy makers did and did not believe their deliberations would be public.
-

Topic modeling of FOMC meeting transcripts

Data:

- 149 FOMC meeting transcripts during, Alan Greenspan's tenure, before and after 1993.
- The unit of observation is a member-meeting.
- The outcomes of interest are
 - the proportion of words devoted to the K different topics
 - the concentration of these topic weights
 - the frequency of data citation.

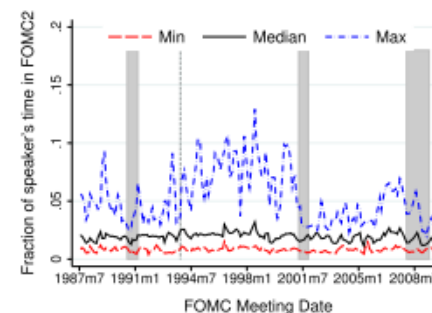
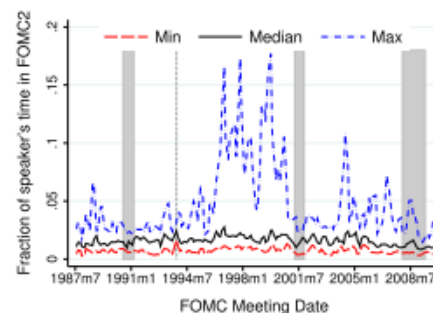
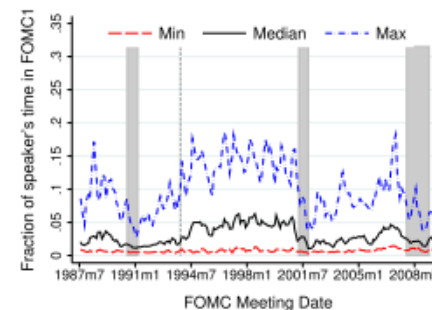
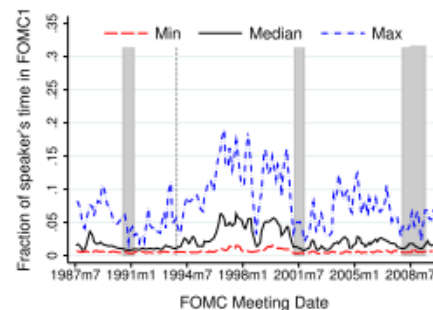
Estimation

- Estimate topics using LDA.
- Use LDA's output to construct outcomes of interest
- Difference / Difference-in-differences regressions that estimate the effects of the change in transparency on outcomes. For example, Hansen et al. estimate

$$y_{it} = \alpha_i + \gamma D(\text{Trans})_t + \lambda X_t + \varepsilon_{it}$$

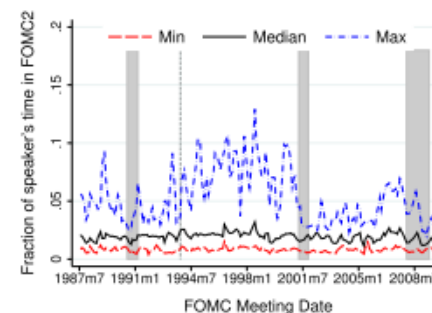
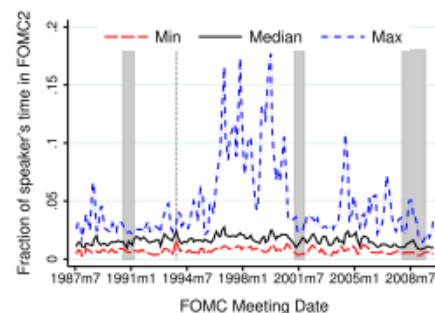
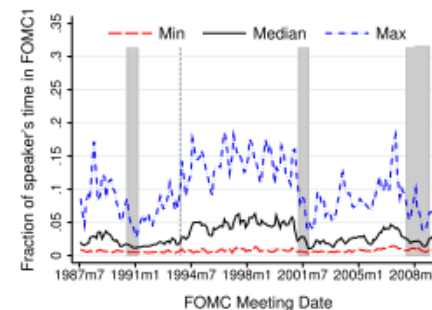
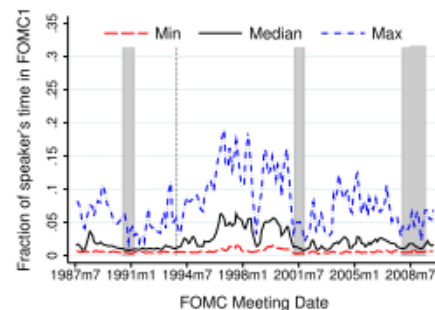
- where
 - y_{it} represents any of the communication measures for member i in time t .
 - $D(\text{Trans})$ is an indicator for being in the transparency regime (1 after November 1993, 0 before).
 - X_t is a vector of macro controls for the meeting at time t .

Pro-cyclical topics



Source: Hansen, McMahon, and Prat (QJE 2017).

Counter-cyclical topics



Source: Hansen, McMahon, and Prat (QJE 2017).

Increased accountability: More references to data

TABLE V
DIFFERENCE RESULTS FOR ECONOMIC SITUATION DISCUSSION (FOMC1):
COUNT MEASURES

Main regressors	Words (1)	Statements (2)	Questions (3)	Numbers (4)
D(Trans)	56.7* [.076]	−0.52 [.162]	−0.039 [.659]	3.71*** [.003]
D(Recession)	−1.95 [.952]	−0.69 [.159]	−0.19 [.314]	−0.71 [.488]
EPU index	0.30 [.186]	−0.00094 [.876]	0.00088 [.586]	0.0040 [.520]
D(2 day)	27.1 [.256]	1.36* [.085]	0.56* [.051]	1.28 [.188]
# of PhDs	6.68 [.561]	−0.45*** [.005]	−0.11*** [.009]	0.51 [.109]
Constant	528*** [.002]	10.0*** [.000]	2.44*** [.000]	1.50 [.740]
Unique members	19	19	19	19
Observations	903	903	903	903
Member FE	Yes	Yes	Yes	Yes
Time FE	No	No	No	No
Meeting section	FOMC1	FOMC1	FOMC1	FOMC1
Transparency effect	9.5*	−10	−2.5	53.2***

Source: Hansen, McMahon, and Prat (QJE 2017).

Increased conformity: increase in similarity

TABLE VI
DIFFERENCE RESULTS FOR ECONOMIC SITUATION DISCUSSION (FOMC1):
TOPIC MEASURES

Main regressors	Concentration (1)	Quant (2)	Avg Sim (B) (3)	Avg Sim (D) (4)	Avg Sim (KL) (5)
D(Trans)	0.0041 [.205]	-0.00027 [.831]	0.0082*** [.001]	0.0012 [.692]	0.032*** [.000]
D(Recession)	0.0061** [.028]	-0.000056 [.968]	0.0020 [.385]	0.015*** [.000]	-0.0017 [.758]
EPU index	3.7e-06 [.890]	-9.6e-06 [.541]	0.000050* [.077]	0.000029 [.300]	0.00015 [.109]
D(2 day)	-0.0040* [.093]	0.0042** [.024]	0.00044 [.802]	-0.0037*** [.001]	0.00051 [.914]
# of PhDs	0.0017 [.255]	-0.00063 [.292]	0.000097 [.885]	0.00079 [.671]	0.00018 [.928]
# Stems	0.000075*** [.000]	8.8e-06** [.049]	-3.5e-06 [.837]	0.000030*** [.001]	0.000049 [.284]
Constant	0.13*** [.000]	0.037*** [.000]	0.89*** [.000]	0.084*** [.001]	0.62*** [.000]
Unique members	19	19	19	19	19
Observation	903	903	903	903	903
Member FE	Yes	Yes	Yes	Yes	Yes
Time FE	No	No	No	No	No
Meeting section	FOMC1	FOMC1	FOMC1	FOMC1	FOMC1
Topics	P1	T4 & T23	P1	P1	P1
Similarity measure	—	—	Bhatta- charyya	Dot product	Kullback- Leibler
Transparency effect	2.5	-0.7	0.9***	1.1	4.9***

Source: Hansen, McMahon, and Prat (QJE 2017).

Hands-on

Go to:

11-hands-on-text-as-data.Rmd

```
slides %>% end()
```

 [Source code](#)

Selected references

Blei, D. M., Ng, A. Y., & Jordan, M. I. (2003). Latent Dirichlet allocation. *Journal of Machine Learning Research*, 3(Jan), 993-1022.

Blei, D. M., & Lafferty, J. D. (2006, June). Dynamic topic models. In *Proceedings of the 23rd international conference on Machine learning* (pp. 113-120). ACM.

Gentzkow, M., Kelly, B.T. and Taddy, M. (forthcoming). *The Quarterly Journal of Economics*.

Hansen, S., McMahon, M., & Prat, A. (2017). Transparency and Deliberation Within the FOMC: A Computational Linguistics Approach. *The Quarterly Journal of Economics*, 133(2), 801–870.

Lafferty, J. D., & Blei, D. M. (2006). Correlated topic models. In *Advances in neural information processing systems* (pp. 147-154).

Loughran, T. and McDonald, B., 2011. When is a liability not a liability? Textual analysis, dictionaries, and 10-Ks. *The Journal of Finance*, 66(1), pp.35-65.

Selected references

Roberts, M.E., Stewart, B.M., Tingley, D., Lucas, C., Leder-Luis, J., Gadarian, S.K., Albertson, B. and Rand, D.G., 2014. Structural topic models for open-ended survey responses. *American Journal of Political Science*, 58(4), pp.1064-1082.

Roberts, M.E., Stewart, B.M. and Tingley, D., 2014. stm: R package for structural topic models. *Journal of Statistical Software*, 10(2), pp.1-40.