

Twitter in the Parliament - A Text-based Analysis of German Political Entities

Patrick Schulze, Simon Wiegrebe

Project partners: *Prof. Dr. Paul W. Thurner, Sandra Wankmüller*
(Geschwister Scholl Institute of Political Science, LMU)

Supervisors: *Prof. Dr. Christian Heumann, Matthias Aßenmacher*

July 16, 2020

Introduction

- Huge amounts of data, especially text, produced by social media
- Field of particular interest in the context of social media and big data:
Politics
 - e.g., Brexit, 2016 presidential election in the US, Facebook data scandal
- Tools of analysis for such data simultaneously provided by advances in
Natural Language Processing (NLP)
- *Topic analysis*: analytical tool for discovery and exploration of latent thematic clusters within text

Introduction

- Key contributions of this project:
 - Construction of dataset containing Twitter posts by members of the German Bundestag and a variety of metadata
 - Application of the *Structural Topic Model* (STM), introduced by **roberts2016model**, to German MPs' Twitter communication
 - Development of new tools for estimation of relationship between topic proportions and metadata
 - Development and application of STM-specific train-test split to enable causal inference

Topic Modeling: Motivation and Theory

Motivation

- Motivating example: excerpt from a scientific article
blei2012presentation

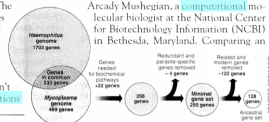
Seeking Life's Bare (Genetic) Necessities

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



* 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

- Question at hand: how to assign colored words to topics?

Topic Modeling: Motivation and Theory

Notation and Terminology (I)

- *Words* w : instances of a vocabulary of V unique *terms*
- *Documents* $d \in \{1, \dots, D\}$: sequences of words of length N_d ; $w_{d,n}$ denoting n -th word of document d
- *Corpus*: collection (or set) of D documents
- *Topics* $k \in \{1, \dots, K\}$: latent thematic clusters within a text corpus; (implicit) representation of a corpus
- *Topic-word distributions* β : probability distributions over words; β_k denoting the word distribution corresponding to the k -th topic

Topic Modeling: Motivation and Theory

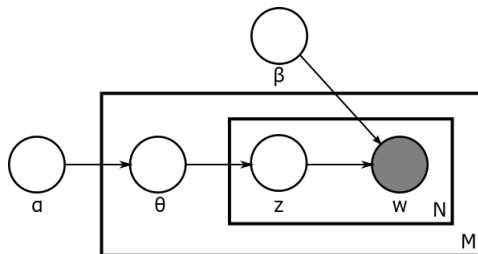
Notation and Terminology (II)

- *Topic assignments* $\mathbf{z}_{d,n}$: assignment of $w_{d,n}$ to a specific topic $k \in \{1, \dots, K\}$; $\beta_{d,n}$ representing the (assigned) word distribution for $w_{d,n}$
 - *Topic proportions* θ_d : proportions of document d 's words assigned to each of the topics; $\sum_{k=1}^K \theta_{d,k} = 1$, for all $d \in \{1, \dots, D\}$
 - *Bag-of-words* assumption: only words themselves meaningful, unlike word order or grammar; equivalent to assuming *exchangeability*
- aldous1985exchangeability**

Topic Modeling: Motivation and Theory

Latent Dirichlet Allocation (LDA) (I)

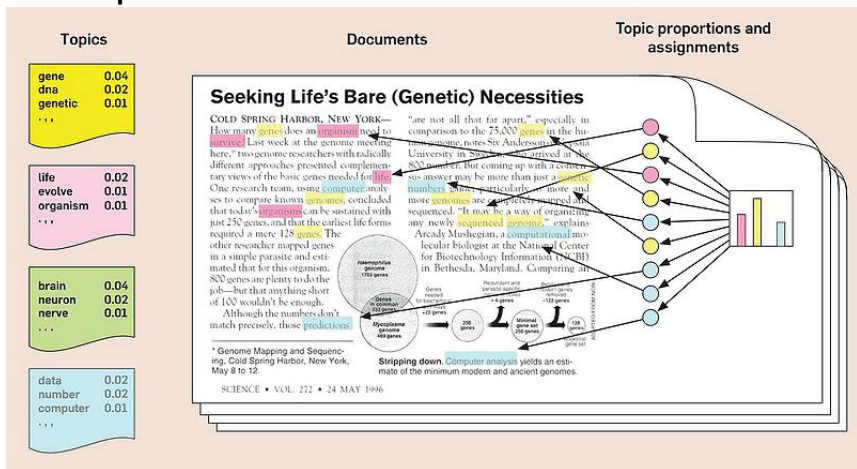
- First topic model with entirely probabilistic generating process: LDA **blei2003latent**
- Generative process for each document $d \in \{1, \dots, D\}$:
 - 1 Draw topic proportions $\theta_d \sim \text{Dir}_K(\alpha)$.
 - 2 For each word $n \in \{1, \dots, N_d\}$:
 - a Draw a topic assignment $z_{d,n} \sim \text{Multinomial}_K(\theta_d)$.
 - b Draw a word $w_{d,n} \sim \text{Multinomial}_V(\beta_{d,n})$.
- Graphical model representation of LDA: **blei2003latent**



Topic Modeling: Motivation and Theory

Latent Dirichlet Allocation (LDA) (II)

- Illustration of topic assignment for the words of a document:
blei2012probabilistic



Topic Modeling: Motivation and Theory

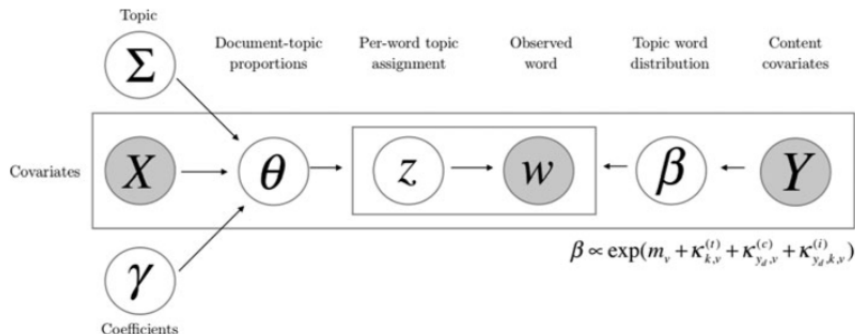
Structural Topic Model (STM)

- Topic model that incorporates document-level metadata:
 - *Topical prevalence* covariates $\mathbf{X} = [\mathbf{x}_1 | \dots | \mathbf{x}_D]^T \in \mathbb{R}^{D \times P}$
 - Categorical *topical content* variable $\mathbf{Y} \in \mathbb{R}^D$ with A levels, i.e., $Y_d \in \{1, \dots, A\}$, for all $d \in \{1, \dots, D\}$
- Generative process for each document $d \in \{1, \dots, D\}$:
 - ① Draw $\boldsymbol{\eta}_d \sim \mathcal{N}_{K-1}(\boldsymbol{\Gamma}^T \mathbf{x}_d^T, \boldsymbol{\Sigma})$, with $\eta_{d,K} = 0$ for model identifiability.
 - ② Normalize $\boldsymbol{\eta}_d$, for all $k \in \{1, \dots, K\}$: $\theta_{d,k} = \frac{\exp(\eta_{d,k})}{\sum_{j=1}^K \exp(\eta_{d,j})}$.
 - ③ For each word $n \in \{1, \dots, N_d\}$:
 - a) Draw topic assignment $\mathbf{z}_{d,n} \sim \text{Multinomial}_K(\boldsymbol{\theta}_d)$.
 - b) If no topical content variable specified: $w_{d,n} \sim \text{Multinomial}_V(\boldsymbol{\beta}_{d,n})$. Otherwise, determine document-specific word distributions $\mathbf{B}_a := [\boldsymbol{\beta}_1^a | \dots | \boldsymbol{\beta}_K^a]$ based on $Y_d = a$, for all topics $k \in \{1, \dots, K\}$; select $\boldsymbol{\beta}_{d,n} := \mathbf{B}_a \mathbf{z}_{d,n}$; and draw word $w_{d,n} \sim \text{Multinomial}_V(\boldsymbol{\beta}_{d,n})$.

Topic Modeling: Motivation and Theory

Graphical Model of the STM

- Visualization of the generative process again through graphical model **roberts2016model**:



Topic Modeling: Motivation and Theory

Inference and Parameter Estimation

- (Hierarchical) Bayesian model \Rightarrow exact inference impossible due to marginal distributions in the denominator of posterior distribution p
- Variational inference: positing a simple distribution family q for latent variables θ_d and \mathbf{z}_d
- Mean-field variational inference: positing full factorizability of approximating posterior q , i.e., $q(\theta_d, \mathbf{z}_{d,n}) = q(\theta_d) * q(\mathbf{z}_{d,n})$
- Then: minimizing Kullback-Leibler divergence between q and p
- STM uses a mean-field variational EM algorithm:
 - E-step: update posterior distributions of latent variables θ_d and $\mathbf{z}_{d,n}$
 - M-step: update model parameters Γ , Σ , and - if present - topical content parameters

Data

Data Collection (I)

- MP-level data: from www.bundestag.de/abgeordnete using Python's *BeautifulSoup* and a *selenium* web driver **van1995python richardson2007beautiful**

Philipp Amthor, CDU/CSU
Jurist

Abgeordnetenbüro
Deutscher Bundestag
Platz der Republik 1
11011 Berlin
Kontakt

Profile im Internet
philipp-amthor.de
[Facebook](#)

Biografie Reden Abstimmungen

Gelesen am 10. November 1992 in Ueckermünde

2011 Abitur am Greifen-Gymnasium Ueckermünde, 2012 bis 2017 Studium der Rechtswissenschaften an der Ernst-Moritz-Arnst Universität Greifswald (Studienschluss mit Praktikum, Stipendiat der Konrad-Adenauer-Stiftung, Kollegat am Jungen Kolleg des Alfred Krupp Wissenschaftskollegs, nebenberuflich u.a. Mitarbeiter verschiedener Abgeordneter des Deutschen Bundestages und des Landtages Mecklenburg-Vorpommern; seit 2017 Doktorand und wissenschaftlicher Mitarbeiter an der Ernst-Moritz-Arnst-Universität Greifswald und zugleich Mitarbeiter einer internationalen Wirtschaftskanzlei in Berlin.

Seit 2008 Mitglied der CDU und der Jungen Union, seit 2010 Mitglied im Landesvorstand der Jungen Union Mecklenburg-Vorpommern; seit 2012 Kreisvorsitzender der Jungen Union Vorpommern-Greifswald; seit 2014 Mitglied des Sozialausschusses des Kreistages Vorpommern-Greifswald; seit 2017 Vorsitzender des CDU-Stadtverbandes Ueckermünde.

Direkt gewählt

Mecklenburg-Vorpommern
 Wahlkreis 016: Mecklenburgische Seenplatte I – Vorpommern-Greifswald II

Mitgliedschaften und Ämter im Bundestag

Ordentliches Mitglied
 > Ausschuss für die Angelegenheiten der Europäischen Union
 > Ausschuss für Innere und Heimat

Stellvertretendes Mitglied
 > Ausschuss für Recht und Verbraucherschutz

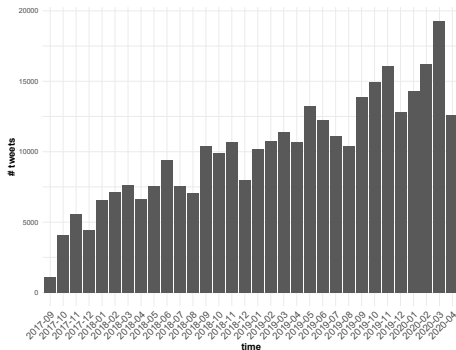
Veröffentlichungspflichtige Angaben

- Twitter profiles: from official party homepages
- Socioeconomic data and 2017 German federal election results: from www.bundeswahlleiter.de.

Data

Data Collection (II)

- Tweets (and further Twitter features): via the official Twitter API using Python's *tweepy* library **roesslein2020tweepy**
- Monthly tweets (after dropping MPs without electoral district) for our period of analysis, September 24, 2017 through April 24, 2020:



- In the following: grouping each MP's tweets on a monthly basis

Data

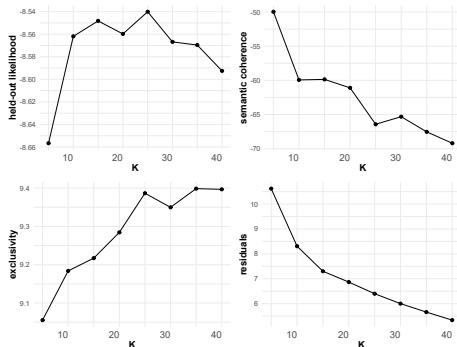
Data Preprocessing

- Preprocessing: in R **R**, using the *quanteda* package **quanteda**
- Transcription of German umlauts (e.g. $\ddot{a} \rightarrow a$) and ligature ($\beta \rightarrow ss$)
- Removal of hyphens: relevant for compound words (e.g., *Corona-Krise* vs *Coronakrise*)
- Transformation of text data into document-feature matrix (DFM); conversion to lowercase; removal of stopwords, units (*kg*, *uhr*), interjections (*aaahhh*, *ufff*), etc.
- Word stemming, i.e., cutting off word endings (e.g., *politisch* \rightarrow *polit*) **lucas2015computer**

Model Selection and Global Characteristics

Model Selection

- Model evaluation metrics for hyperparameter K (number of topics):



- "Best" trade-off: $K = 15$

Model Selection and Global Characteristics

Labeling (I)

- Three-step procedure for labeling
- First step: top words for different weighting methodologies

Topic 1 Top Words:

Highest Prob: buerg, link, merkel, frau, sich

FREX: altpartei, islam, linksextremist, asylbewerb, linksextrem

Lift: eitan, 22jaehrig, abdelsamad, abgehalftert, afdforder

Score: altpartei, linksextremist, frauenkongress, islamist, boehring

Topic 3 Top Words:

Highest Prob: brauch, wichtig, leid, dank, klar

FREX: emissionshandel, soli, marktwirtschaft, feedback, co2steu

Lift: aequivalenz, altersvorsorgeprodukt, bildungsqualitaet, co2limit, co2meng

Score: emissionshandel, co2limit, basisrent, euet, technologieoff

Topic 4 Top Words:

Highest Prob: sozial, miet, kind, arbeit, brauch

FREX: mindestlohn, miet, wohnungsbau, mieterinn, loehn

Lift: auseinanderfaellt, baugipfel, bestandsmiet, billigflieg, binnennachfrag

Score: miet, mieterinn, mietendeckel, grundsicher, bezahlbar

Topic 6 Top Words:

Highest Prob: gruen, klimaschutz, brauch, klar, euro

FREX: fossil, erneuerbar, kohleausstieg, verkehrsminist, verkehrsw

Lift: abgasbetrug, abgebagert, abschalteinricht, abschaltet, ammoniak

Score: erneuerbar, fossil, zdebel, verkehrsminist, klimaschutz

Model Selection and Global Characteristics

Labeling (II)

- Word cloud of **Highest Prob** top words (for topic 1):



- Word size corresponding to word frequency in topic 1

Model Selection and Global Characteristics

Labeling (III)

- Second step: look at documents (i.e., original tweets) with highest proportion of topic 1



Martin Hess
@Martin_Hess_AfD

Ehem. Verfassungsrichter bestätigt AfD-Forderung:
Zurückweisung illegaler Migranten dringend geboten.
Gegenwärtige Politik widerspricht dem Verstand und
auch der Verfassung. Wir müssen zurück zu Recht &
Ordnung, wie die #AfD seit fast 3 Jahren fordert!



Hans-Jürgen Papier hält Zurückweisung von Migranten an deutscher Grenze für ...
Im Asylstreit meldet sich nun Ex-Verfassungsrichter Papier zu Wort. Die
Zurückweisung von Migranten an den Grenzen sei zwingend nötig, schreibt er in...
welt.de

Model Selection and Global Characteristics

Labeling (IV)

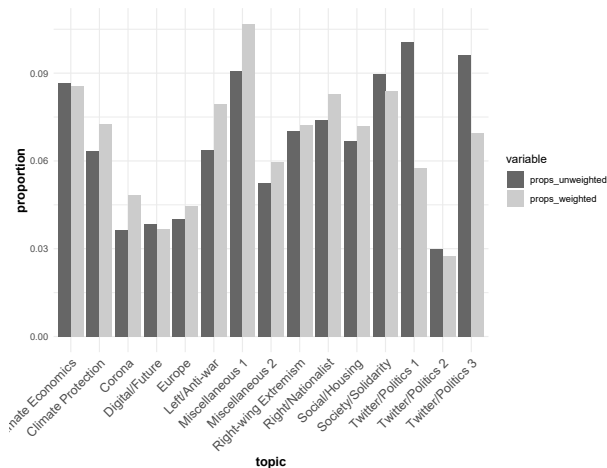
- Third step: assigning labels

Topic 1	Right/Nationalist
Topic 2	Miscellaneous 1
Topic 3	Climate Economics
Topic 4	Social/Housing
Topic 5	Digital/Future
Topic 6	Climate Protection
Topic 7	Europe
Topic 8	Corona
Topic 9	Left/Anti-war
Topic 10	Twitter/Politics 1
Topic 11	Twitter/Politics 2
Topic 12	Miscellaneous 2
Topic 13	Twitter/Politics 3
Topic 14	Right-wing Extremism
Topic 15	Society/Solidarity

Model Selection and Global Characteristics

Global Topic Proportions

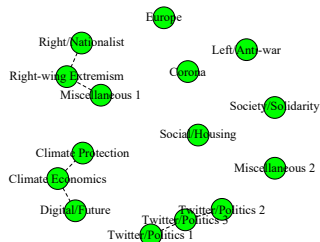
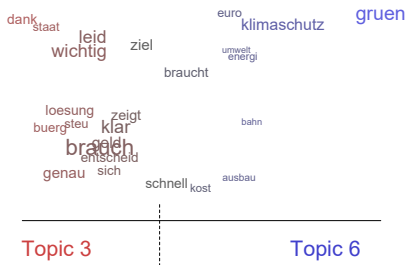
- Illustration of **global** topic proportions:



Model Selection and Global Characteristics

Global Topic Correlations

- Vocabulary overlap (left) and topic correlations (right):



Covariate-level Topic Analysis

Overview

- Exploration of estimated topical structure with respect to different dimensions, e.g. membership in political party, time, ...
- Specifically: examining relationship between document-level prevalence covariates \mathbf{x}_d and topic proportions θ_d
- Natural idea: regress topic proportions on prevalence covariates
 - In standard regression analysis, dependent variable as realization of random variable
 - In STM, however: posterior of topic proportions θ_d accessible
 - Loss of information if "naïvely" using mean/mode of this posterior as dependent variable of regression
 - Solution: performing sampling technique known as "method of composition" in social sciences
- Alternatively: direct assessment of logistic normal distribution with estimated topical prevalence parameters $\hat{\Gamma}$ and $\hat{\Sigma}$

Covariate-level Topic Analysis

Method of Composition: Usage within R Package *stm*

- Notation:

- $\boldsymbol{\theta}_{(k)} := (\theta_{1,k}, \dots, \theta_{D,k})^T \in [0, 1]^D$: proportion of k -th topic for all D documents
- $q(\boldsymbol{\theta}_{(k)} | \mathbf{X}, \mathbf{W})$: approximate variational posterior of $\boldsymbol{\theta}_{(k)}$
- $q(\hat{\boldsymbol{\xi}} | \mathbf{X}, \boldsymbol{\theta}_{(k)})$: (normal) distribution of estimated regression coefficients $\hat{\boldsymbol{\xi}}$ from OLS regression $\boldsymbol{\theta}_{(k)} = \mathbf{X}\boldsymbol{\xi} + \epsilon$, where $\epsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$

- Method of composition:

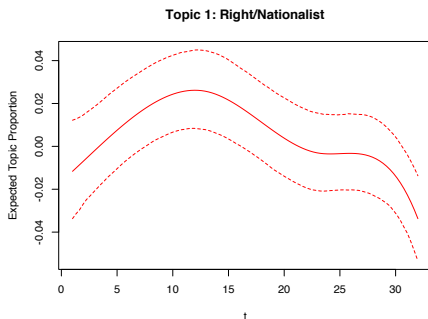
- ① Draw $\boldsymbol{\theta}_{(k)}^* \sim q(\boldsymbol{\theta}_{(k)} | \mathbf{X}, \mathbf{W})$.
 - ② Draw $\hat{\boldsymbol{\xi}}^* \sim q(\hat{\boldsymbol{\xi}} | \mathbf{X}, \boldsymbol{\theta}_{(k)}^*)$.
- Then: $\hat{\boldsymbol{\xi}}_1^*, \dots, \hat{\boldsymbol{\xi}}_m^*$ is an i.i.d. sample from the marginal posterior of regression coefficients

$$q(\boldsymbol{\xi} | \mathbf{X}, \mathbf{W}) = \int_{\boldsymbol{\theta}_{(k)}} q(\boldsymbol{\xi} | \mathbf{X}, \boldsymbol{\theta}_{(k)}) q(\boldsymbol{\theta}_{(k)} | \mathbf{X}, \mathbf{W}) d\boldsymbol{\theta}_{(k)}$$

Covariate-level Topic Analysis

Method of Composition: Usage within R Package *stm*

- Problem: OLS regression not suitable for (sampled) proportions, which are restricted to interval (0,1)
- ⇒ Estimated relationship between proportions and prevalence covariates potentially producing negative estimated proportions



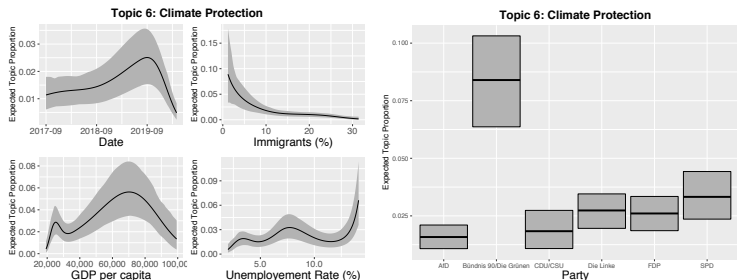
Topic 6: Climate Protection



Covariate-level Topic Analysis

Method of Composition: Extension of existing approach

- Instead of OLS regression, we can use a beta regression or a quasibinomial GLM (both with logit-link) to adequately model proportions
- In this case, regression coefficients are *asymptotically* normally distributed



Covariate-level Topic Analysis

Problem: Univariate Modeling of Proportions

- Recall, by assumption: $\theta_d \sim \text{LogisticNormal}(\Gamma^T \mathbf{x}_d^T, \Sigma)$
- Logistic normal distribution assuming high dependence among individual components
- However, *univariate* k -th topic proportion used as dependent variable in regression within method of composition
- Problem with this approach: dependence among components neglected \Rightarrow uncertainty estimates particularly unrealistic

Covariate-level Topic Analysis

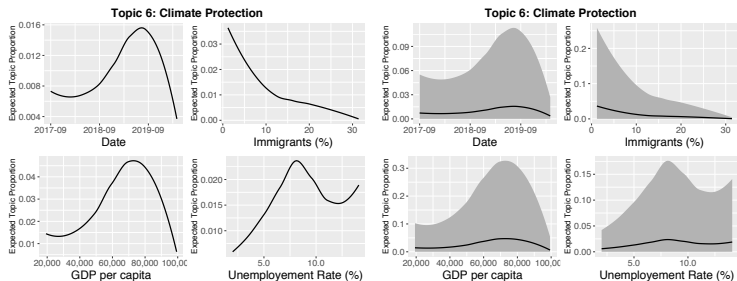
Multivariate Modeling via Logistic Normal Distribution (I)

- Inference within STM involves finding estimates $\hat{\Gamma}$ and $\hat{\Sigma}$
- Idea: plugging estimates into logistic normal distribution \Rightarrow for a given covariate value \mathbf{x}_d^* , "predicting" topic proportion as $\theta_d^* \sim \text{LogisticNormal}(\hat{\Gamma}^T (\mathbf{x}_d^*)^T, \hat{\Sigma})$
- Ideally: applying fully Bayesian approach and sampling from (variational) posterior of Γ (and updating Σ , obtained via MLE) \Rightarrow "Predictive Posterior" of topic proportions
- However, output obtained from R package *stm* not allow for simple implementation of such a procedure (i.e., sampling from variational posterior of Γ and updating Σ)
 - Yet, possible in theory!

Covariate-level Topic Analysis

Multivariate Modeling via Logistic Normal Distribution

- Still, our results suggest a high discrepancy between:
 - Distribution of topic proportions assumed in generative process of STM
 - Impression we gain of this distribution via separate modeling of topics.
- Fully Bayesian approach: most likely yielding even higher uncertainty



Causal Inference

Correlation vs. Causality

- In previous section: assessment of relationship between metadata and topic proportions
- Framework to be used to *explore* topics with respect to different dimensions
- In particular, *causal* interpretation of results generally not justified ("correlation vs. causality")
- When making causal inference, need to consider that topic proportions are *latent* variables
- Possible solution: conducting a train-test split

Causal Inference

Identification Problem and Overfitting

- Setup: two groups (treatment and control), individuals otherwise similar
- Objective: quantifying treatment effect, in our case effect of treatment on prevalence of specific topic.
- Necessary assumption: response of an individual depending only on their treatment
- *Identification problem*: estimating topic model to discover latent topic proportions can introduce additional dependency among individuals \Rightarrow response of each individual *not* only determined by treatment of that individual!
- *Overfitting*: fitted topic model might mistake noise for patterns in some way \Rightarrow response again not solely determined by treatment of an individual, but additionally by specific characteristics of other individuals

Causal Inference

Train-test split

- Idea: splitting data \mathcal{D} into training set $\mathcal{D}_{\text{train}}$ and test set $\mathcal{D}_{\text{test}}$
- Training set $\mathcal{D}_{\text{train}}$ used to determine a model that infers latent topic proportions from a given text
- Test set $\mathcal{D}_{\text{test}}$ used to assess relation between *predicted* test set topic proportions and test set prevalence covariates
- Identification problem solved: model used for prediction determined by training set observations \Rightarrow treatment of test set observations not dependent on other individuals' treatment from test set.
- Overfitting also solved: noise from training set very unlikely to be replicated on test set

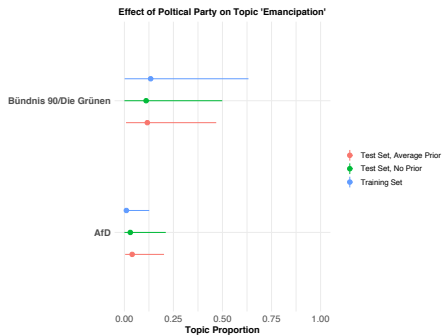
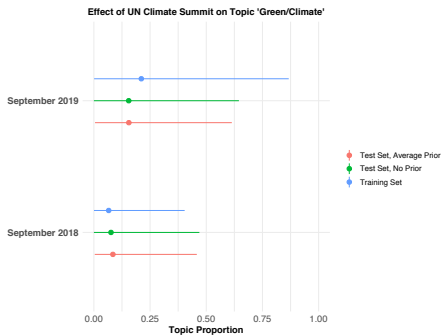
Causal Inference

Implementation within the STM

- Inputting documents, i.e., words and metadata from the training set $\mathcal{D}_{\text{train}}$, to obtain estimates $(\hat{\beta}_{\text{train}}, \hat{\Gamma}_{\text{train}}, \hat{\Sigma}_{\text{train}})$ using the STM
- Then, estimating (variational) posterior of test set topic proportions, conditional on the model parameters $(\hat{\beta}_{\text{train}}, \hat{\Gamma}_{\text{train}}, \hat{\Sigma}_{\text{train}})$ from training set $\mathcal{D}_{\text{train}}$ as well as words \mathbf{W}_{test} from test set $\mathcal{D}_{\text{test}}$
- Estimation of (variational) posterior conditional on data and training set parameters via E-step of (variational) EM algorithm
- Benefit of using the STM: covariate information from training set directly used to predict topic proportions on test set
- Important: Covariate information from test set must not be used!
 - Otherwise: predicting different topic proportions for two documents from test set with exact same words if prevalence covariates differ
 - However, causal effect should be zero in such a case!

Causal Inference

Results (I)



Causal Inference

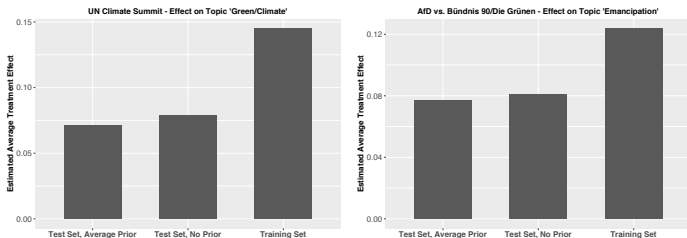
Results (II)

- UN Climate Action Summit 2019 held on September 23, 2019
- As observed, topic associated with climate issues much more prevalent during that time than the year before
- MAP estimates for different prior specifications on test set rather similar, yet estimated effect for training data much larger
- Similar results for effect of political party on topic labeled as 'Emancipation': average difference of estimated topic proportions between both parties larger for the training data
- Additionally: credible intervals on the training data different from those on the test data in both cases

Causal Inference

Results (III)

- Estimation of treatment effect: determining the average difference of predicted topic proportions between both groups



- Treatment effect larger if "naïvely" estimated solely on training data in both cases!

Discussion

- Use of dataset and results for future (political) analyses
- Topic-metadata relationship:
 - Room for methodological improvement
 - Applicability in predictive tasks
- Train-test split and causal inference:
 - Alternative model designs
 - Natural experiments

Bibliography