

# Political news coverage of massmedia

...

Franziska Löw\*

November 2020

## Introduction

In democracies, the media fulfill fundamental functions: They should inform the people, contribute to the formation of opinion through criticism and discussion and thus enable participation. In recent decades, however, concern has grown about the role of the media in politics in general and in election campaigns in particular. They are criticized for influencing election results through their reporting and for helping populist parties in particular to flourish. After the 2017 federal elections in Germany, for example, the media were accused of contributing to the success of the right-wing populist AfD by increasingly including the party's content and using the same language in their articles as the AfD. Representatives of these media houses strongly opposed this accusation. The purpose of this study is to examine whether there is evidence that support the accusation of biased media reporting, especially during election campaigns.

For advertising-financed media the battle for the attention of the recipients is at the center of economic decisions. Online media in particular, which offer their content to a large extent free of charge and generate their revenue through advertising space, compete for the scarce resource of attention. Consumers pay a non-monetary price providing their attention, which the media platform bundles and sells on to advertising customers. This business model corresponds to that of a platform market, in which media companies act as platforms that connect the market of advertising with the reader market to exploit the indirect network effects between them (Dewenter and Rösch 2014). A profit-maximizing publisher therefore directs its economic decisions according to what will attract the most attention.

This conclusion, derived from the economic theory of platform markets, corresponds to the notion of media logic, a central concept in the field of media and communication studies (Takens et al. 2013). The debate about media logic is embedded in the broader discussion about the interaction between the press, politics and the public. The underlying thesis is that the content of political news is the product of news values and narrative techniques that media use to attract audiences (Strömbäck 2008). According to Takens et al. (2013), three content attributes highly correspond with news values and influence how journalists interpret political events: 1) personalized content, i.e., the focus on individual politicians; 2) the framing of politics as a contest and 3) negative coverage. Similarly Blassnig et al. (2019) states that media primarily focus on news factors, i.e. the factors that turn an event into news worth reporting like conflict, drama, negativity, surprise or proximity. Likewise populist messages often co-occur with negative, emotionalized, or dramatized communication style, thus utilizing similar mechanisms as the media logic, respectively the attention economy. In fact, Blassnig et al. (2019) shows that populist key messages by political and media actors in news articles provoke more reader comments under these articles. Media competing for the attention of readers therefore have an incentive to pick up on the key messages of these parties.

Political parties want the media agenda to be congruent with their own agenda to define the issue-based criteria on which they will be evaluated by voters, especially during election campaigns (Eberl, Boomgaarden,

---

\*Institut für Industrieökonomik, Helmut Schmidt Universität. Email: .

and Wagner 2017). Parties instrumentalize their public relations in order to highlight issues that they are perceived to be competent on, that they “own” and that are important to their voters (Kepplinger and Maurer 2004).

But does increased reporting also lead to rising survey results? Especially if this reporting is largely negative, which is the case for reporting on AfD during election campaign phase. In political science, several studies have examined at least the first aspect of this question (see for example Druckman and Parkin (2005), Eberl (2018)). In general, it is assumed that smaller, non-established parties in particular benefit from placing their topics in the media in order to get them into the voters’ heads. Here, the tendency of the reporting is irrelevant but rather the quantity is decisive.

However, the causal relationship between reporting and voter preferences is not the subject of this study. Rather, it is intended to investigate whether differences exist in media coverage of different parties before and after the 2017 federal elections in Germany. In order to answer these and other media-related questions in the political context, quantifying the content of media is a prerequisite. One of the key challenges is to determine the features that are used to describe media content (audio, video, text). Studies that rely on quantifying media content for their analyses use, for example, visibility (how often political actors appear in the media) or tonality (how they are evaluated). Other studies examine the topics discussed or the language used in the media, in order to identify whether political actors are able to place their own policy positions in the media. Leading studies from economic literature, for example, examine how often a newspaper quotes the same think tanks (Groseclose and Milyo (2005), Lott and Hassett (2014)) or uses the same language (Gentzkow and Shapiro 2004) as members of Congress. Following this approach, the present paper compares topics discussed in media outlets with topics addressed in the press releases of the parties in the German “Bundestag”, to measure the content similarity between online news and parties press releases.<sup>1</sup> To discover the latent topics in the corpus of text data, the structural topic model (STM) developed by M. E. Roberts, Stewart, and Airoldi (2016) is applied. This probabilistic text model results in a probability distribution for each document across all topics, which is then aggregated to calculate the degree of difference between the news articles of different media providers and the press releases of the parties using a linear regression model.

## Literature review

Newspaper articles and their metadata, such as publisher and publication date have been subject of investigation in economic research.

## Background information

### The political situation in Germany (June 2017 - March 2018)

The articles analyzed in this paper cover a period from June 1, 2017 to March 1, 2018 and thus cover both the most important election campaign topics for the Bundestag elections on September 24, 2017 and the process of forming a government that lasted until February 2018. After four years in a grand coalition with the Social Democrats (SPD), German Chancellor Angela Merkel, member of the conservative party CDU/CSU (also known as Union), ran for re-election. The SPD nominated Martin Schulz as their candidate.

On the right side of the political spectrum, AfD (alternative for Germany) managed to be elected to the German Bundestag for the first time in 2017. The political debate about the high refugee numbers of the past years brought a political upswing to the AfD, which used the dissatisfaction of parts of the population to raise its own profile. In the course of the reporting on the federal elections, leading party members of the AfD as well as party supporters repeatedly accused the mass media of reporting unilaterally and intentionally presenting the AfD badly.

---

<sup>1</sup>For the sake of simplicity, both news articles and press releases will be referred to as documents in the following.

After the election, the formation of a government was difficult due to the large number of parties elected to the Bundestag and the considerable loss of votes by the major parties CDU/CSU and SPD. Since all parties rejected a coalition with the AfD, numerically only two coalitions with an absolute parliamentary majority were possible: a grand coalition (“GroKo” - from the German word Große Koalition) of CDU/CSU and SPD, and a Jamaica coalition (coalition of CDU/CSU, FDP (economic liberal party) and B90/Die Grünen (Bündnis 90/Die Grünen, green party)). The grand coalition was initially rejected by the SPD. The four-week exploratory talks on the possible formation of a Jamaica coalition officially failed on November 19, 2017 after the FDP announced its withdrawal from the negotiations. FDP party leader Christian Lindner said that there had been no trust between the parties during the negotiations. The main points of contention were climate and refugee policy. CDU and CSU regretted this result, while B90/Die Grünen sharply criticized the liberals’ withdrawal. The then Green leader Cem Özdemir accused the FDP of lacking the will to reach an agreement.

After the failure of the Jamaica coalition talks, a possible re-election or a minority government as alternatives were discussed in the media before the SPD decided to hold coalition talks with the CDU/CSU. This led to great resistance from the party base, which called for a party-internal referendum on a grand coalition. After the party members voted in favor of the grand coalition, a government was formed 171 days after the federal elections.

Figure 1 shows that support for the two major popular parties has been declining in recent months since August 2017, with the CDU/CSU again showing positive survey results since November 2017.<sup>2</sup> However, the poll results of the SPD have been falling since March 2017. At the same time, the AfD in particular has been recording increasingly positive survey results since June 2017.

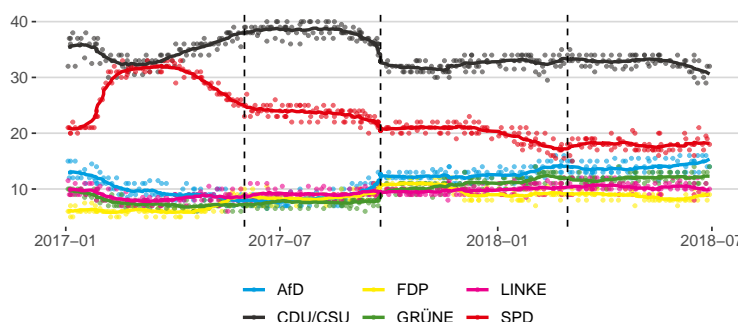


Figure 1: Election polls

## German online news market

The analysis performed in this paper is based on the news articles of the following news websites: Bild.de, DIE WELT, FOCUS ONLINE, Handelsblatt.com, SPIEGEL ONLINE, stern.de, ZEIT ONLINE. As can be seen from Figure 2(a), expect for Handelsblatt.com (position 53), these media outlets are among the top 30 German online news providers in the period under review in terms of visits.<sup>3</sup>

The main source of income for these privately managed media houses is digital advertising, even though paid content is playing an increasingly important role. However, according to a survey on digital news by the Reuters Institute (Newman et al. 2018) only 8% of respondents pay for online news. The online survey for German data was undertaken between 19th - 22nd January 2018 by the Hans Bredow Institute<sup>4</sup>

<sup>2</sup>For each party the survey results of the seven major institutes are considered. To calculate a smooth line for each party on each day, the moving average within 15 days (7 before the day, 7 after the day, and the day itself) is estimated. The data source is <https://www.wahlrecht.de/>.

<sup>3</sup>The term visit is used to describe the call to a website by a visitor. The visit begins as soon as a user generates a page impression (PI) within an offer and each additional PI, which the user generates within the offer, belongs to this visit.

<sup>4</sup><https://www.hans-bredow-institut.de/de/projekte/reuters-institute-digital-news-survey>

with a total sample size of 2038 adults (aged 18+) who access news once a month or more. Among other questions, participants were asked which news sources they use to access news online.<sup>5</sup> The results displayed in Figure 2(b) indicate that the media used for the analysis play a relevant role in their consumption.

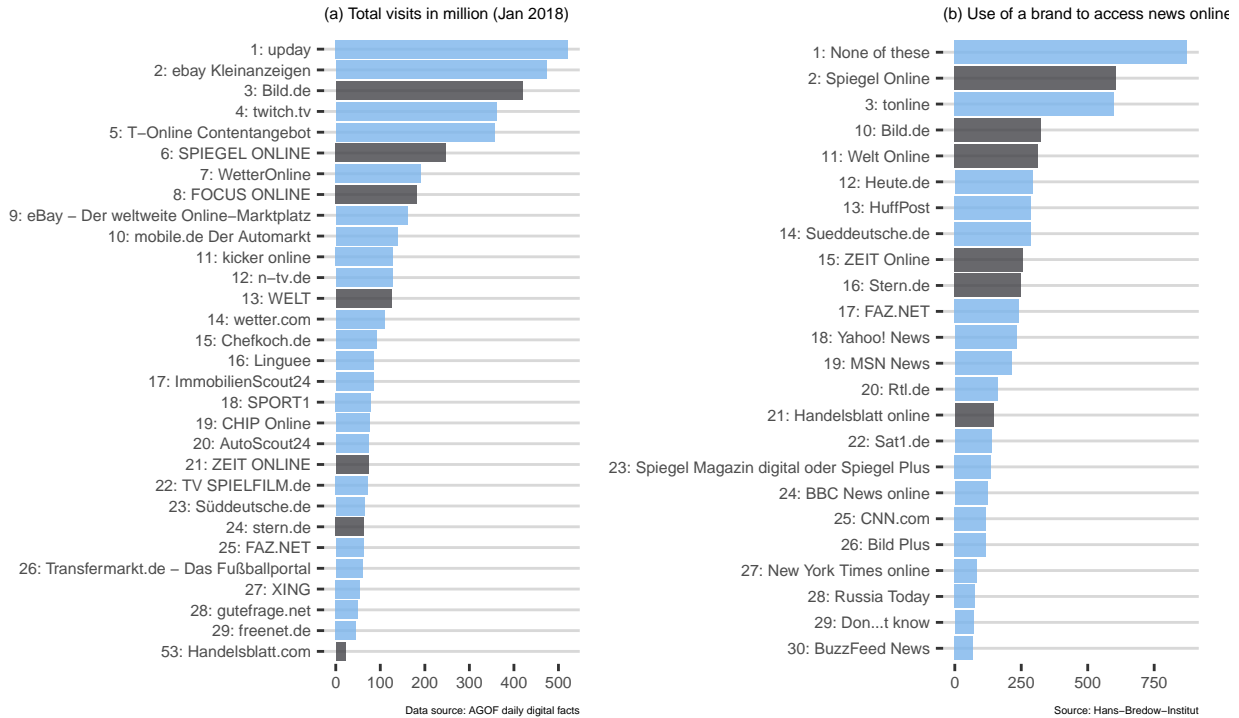


Figure 2: Selected german news brands

## Data

I conduct the estimation on a sample of 16,473 online news articles from the seven German news providers mentioned in the previous section<sup>6</sup> about domestic politics and press releases of the seven parties that have been in the Bundestag since the 2017 federal elections<sup>7</sup>. Both news articles and press releases are dated from June 1, 2017 to March 1, 2018.

News articles scraped from the Webhose.io API.<sup>8</sup> In order to consider only news about national politics, the articles were filtered based on their URL.

Figure 3 shows the distribution of the number of articles by date and media outlet. There is a high peak around the federal elections on September, 24th and another one shortly after the failure of the Jamaica coalition talks on November, 19th (indicated by the red dotted lines).<sup>9</sup> Furthermore, Figure 3 shows that *DIE WELT* published the most articles on domestic policy, followed by *stern.de*, *Handelsblatt* and *FOCUS ONLINE*.

<sup>5</sup>The exact question was: “Which of the following brands have you used to access news online in the last week (via websites, apps, social media, and other forms of Internet access)? Please select all that apply.”

<sup>6</sup>Bild.de, DIE WELT, FOCUS ONLINE, SPIEGEL ONLINE, stern.de, ZEIT ONLINE, Handelsblatt

<sup>7</sup>CDU, SPD, B90/Grüne, FDP, AfD, Die Linke

<sup>8</sup>For more information see <https://docs.webhose.io/reference#about-webhose>. The scraping code was written in Python and can be made available on request.

<sup>9</sup>The peak in July especially for *stern.de* is due to increased reporting about the G20 summit in Hamburg.

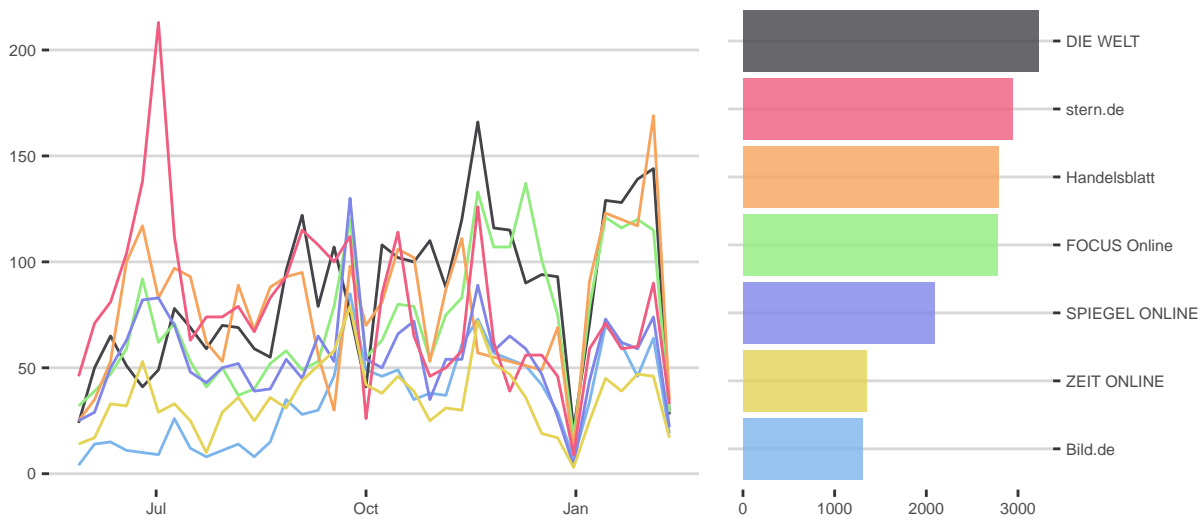


Figure 3: Distribution of news articles

The press releases were scraped from the public websites of the political parties and parliamentary groups using an automated script written in *Python*.<sup>10</sup>

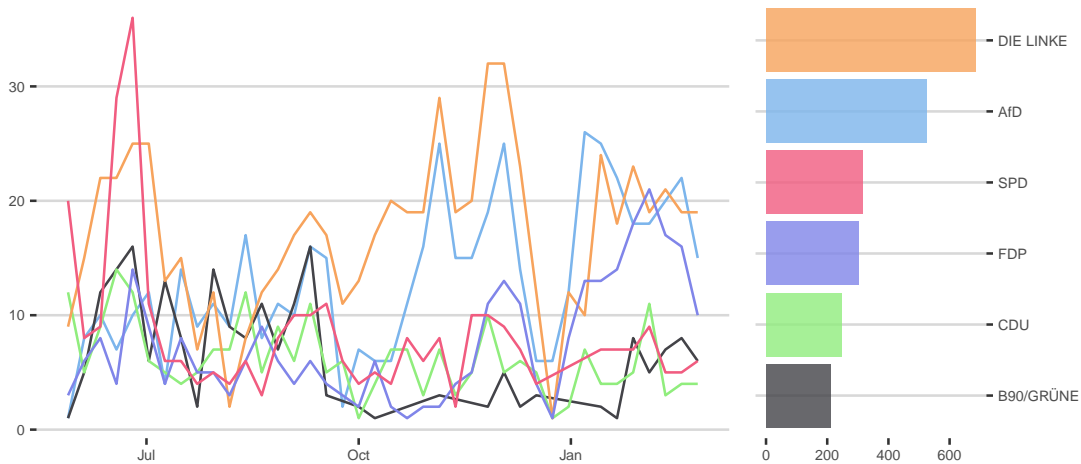


Figure 4: Distribution of press releases

Looking at the boxplots of text length (Figure 5), it becomes evident that:<sup>11</sup>

- mean and median text length of news articles is higher than in press releases
- Handelsblatt published new articles with the highest median text length (488) followed by ZEIT ONLINE (459), however DIE WELT has the article with the highest word count (14.507).
- press releases of CDU have the highest median (256), but also the highest standard dev. (106). press releases of FDP have the lowest median (144).

<sup>10</sup>The scraping code was written in Python and can be made available on request.

<sup>11</sup>See Table 9 for an overview of the summary statistics.

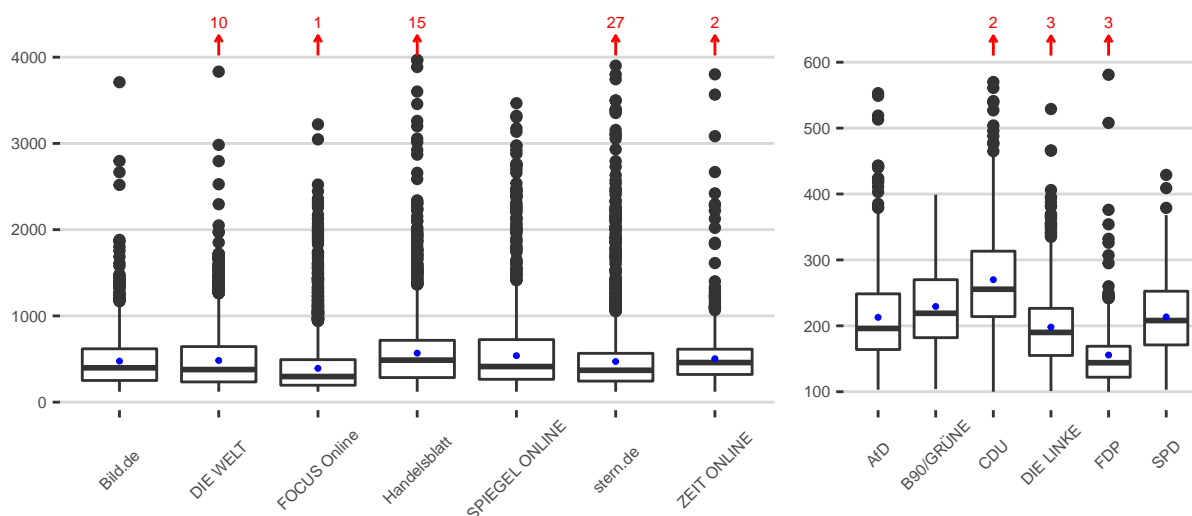


Figure 5: Text length

## Data preparation

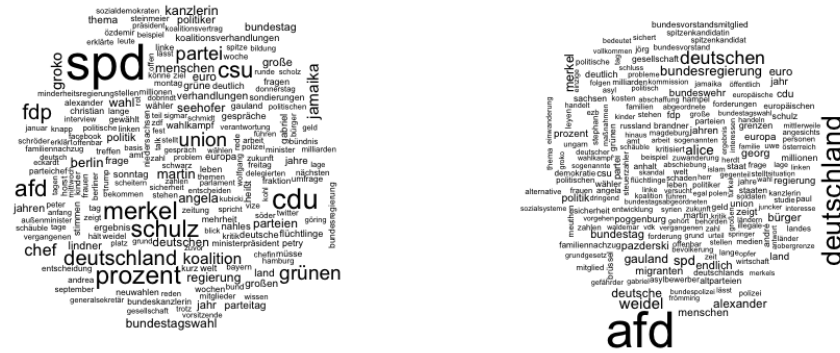
To use text as data for statistical analysis, different pre-processing steps have to be conducted. In fact, in order to use text as data and reduce the dimensionality to avoid unnecessary computational complexity and overfitting, pre-processing the text is a central task in text mining (Gentzkow, Kelly, and Taddy 2017, @bholat\_text\_2015). Intuitively the term frequency (tf) of a word is a measure of how important that word may be for the understanding of the text. To visualize these terms, word clouds are a commonly used technique in text mining as they translate the tf into the size of the term in the cloud. As can be seen in ??, problems arise with words that are highly frequent. For example “die”, or “der (eng.”the“),”und” (eng. “and”), and “ist” (eng. “is”) are extremely common but unrelated to the quantity of interest. These terms, often called stop words (Gentzkow, Kelly, and Taddy 2017), are important to the grammatical structure of a text, but typically don’t add any additional meaning and can therefore be neglected.



Figure 6: Wordcloud before pre-processing

To remove distorting words, the pre-defined stop word list from the Snowball project<sup>12</sup> is used together with a customized, domain-specific list of stop-words. Additionally punctuation character (e.g. ., ,, !, ?, etc.) and all numbers are removed from the data. A next step to reduce the dimensionality of text data is to apply an adequate stemming technique. Stemming is a process by which different morphological variants of a word are traced back to their common root. For example, “voting” and “vote” would be treated as two instances of the same token after the stemming process. There are many different techniques for the stemming process. I apply the widely used Porter-Stemmer algorithm, which is based on a set of shortening rules that are applied to a word until it has a minimum number of syllables.<sup>13</sup>

As an example, the following word clouds represent the most frequent words of the pre-processed articles for Bild.de and press releases of AfD. It becomes evident that these are texts discussing domestic policy issues. The SPD in particular seems to be highly frequent for *Bild.de*.



The next step is to divide the entire data set into individual documents and to represent these documents as a finite list of unique terms. In this setting, each news article and each press release represents a document  $d$ , whereby each of these documents can be assigned to a news website or a party. The sum of all documents forms what is called the corpus. For each document  $d \in \{1, \dots, D\}$  the number of occurrences of term  $v$  in document  $d$  is computed, in order to obtain the count  $x_{d,v}$ , where each unique term in the corpus is indexed by some  $v \in \{1, \dots, V\}$  and where  $V$  is the number of unique terms. The  $D \times V$  matrix  $\mathbf{X}$  of all such counts is called the document-term matrix. Each row in this matrix represents a document, where each entry in this row counts the occurrences of a unique term in that document. Table 1 provides a sample output of the document-term matrix used in this paper, where each document is represented by a unique id (the row name in the example below). This representation is often referred to as the bag of words model (Gentzkow, Kelly, and Taddy 2017), since the order in which words are used within a document is disregarded.

Table 1: Document-term matrix - sample values

	kandidatin	angestiegen	sonderpartei	tag	jahre	türkischen	genannten	steffen
154	0	0	0	6	0	0	0	0
4213	0	0	0	0	0	0	0	0
2174	0	0	0	0	0	0	0	0
18108	0	0	0	0	0	0	0	0
3723	0	0	0	1	0	0	0	0
1113	0	0	0	1	0	0	0	0
14221	0	0	0	0	0	0	0	0
5449	0	0	0	0	0	0	0	0
4949	0	0	0	0	0	0	0	0
11173	0	0	0	0	0	0	0	0

<sup>12</sup><http://snowball.tartarus.org/algorithms/german/stop.txt>

<sup>13</sup><https://tartarus.org/martin/PorterStemmer/>



# Estimate topic similarity of documents

## A structural topic model to identify the latent topics

To discover the latent topics in the corpus of press releases and news articles, a structural topic modeling (STM) developed by (M. E. Roberts, Stewart, and Airoldi 2016) is applied. In general, topic models formalize the idea that documents are formed by hidden variables (topics) that generate correlations among observed terms. They belong to the group of unsupervised generative models, meaning that the true attributes (topics) cannot be observed. The STM developed by (M. E. Roberts, Stewart, and Airoldi 2016) is a recent extension of the standard topic modelling technique, labeled as latent Dirichlet allocation (LDA), which refers to the Bayesian model in (Blei, Ng, and Jordan 2003) that treats each word in a topic and each topic in a document as generated from a Dirichlet - distributed prior.<sup>14</sup>

The underlying idea for these models suggests that each individual topic  $k$  potentially contains all of the unique terms within the vocabulary  $V$  with different probability. Therefore, each topic  $k$  can be represented as a probability vector  $\phi_k$  over all unique terms  $V$ . Simultaneously, each individual document  $d$  in the corpus can be represented as a probability distribution  $\theta_d$  over the  $K$  topics.

The STM is an extension of the LDA process since it allows covariates of interest (such as the publication date of a document or it's author) to be included in the prior distributions for both topic proportions ( $\theta$ ) and topic-word distributions ( $\phi$ ). This way, STM offers a method of 'structuring' the prior distributions in the topic model, including additional information in the statistical inference procedure, while LDA assumes that  $\theta$  Dirichlet( $\alpha$ ) and  $\phi$  Dirichlet( $\beta$ ), where  $\alpha$  and  $\beta$  are fitted with the model.

In order to include the covariates in the statistical inference procedure, two design matrices of covariates ( $X$  and  $Z$ ) are specified, where each row defines a vector of covariates for a specific document. In  $X$ , the covariates for topic prevalence are given, so that the probability of a topic for each document varies according to  $X$ , rather than resulting from a single common prior. The same applies to  $Z$ , in which the covariates for the word distribution within a topic are specified. The underlying data generating process to generate each individual word  $w_{d,n}$  in a document  $d$  for the  $n^{th}$  word-position can be described as follows:

- for each document  $i$ , draw its distribution of topics  $\theta_d$  depending on the metadata included in the model defined in  $X$ ;
- for each topic  $k$ , draw its distribution of words  $\phi_k$  depending on the metadata included in the model defined in  $Z$ ;
- for each word  $n$ , draw its topic  $z_n$  based on  $\theta_i$ ;
- for each word  $n$ , draw the term distribution for the selected topic  $\phi_{z_{d,n}}$ .

One crucial assumption to be made for topic models like LDA or STM is the number of topics ( $K$ ) that occur over the entire corpus. There is not a "right" answer to the number of topics that are appropriate for a given corpus (Grimmer and Stewart 2013). (M. Roberts, Stewart, and Tingley 2016b) propose to measure topic quality through a combination of semantic coherence and exclusivity of words to topics. Semantic coherence is a criterion developed by (Mimno et al. 2011) and is closely related to pointwise mutual information (Newman et al. 2010): it is maximized when the most probable words in a given topic frequently co-occur together.

Using the function *searchK* from the *stm* package several automated tests are performed to help choose the number of topics including the average exclusivity and semantic coherence as well as the held out likelihood (Wallach, Mimno, and McCallum 2009) and the residuals (Taddy 2012). This process revealed that a model with 40 topics best reflects the structure in the corpus. Furthermore, I use the author and bi-week dummies of a document as topical prevalence variable. In other words, I assume that the probability of a topic to be included in a news article or a press release depends on the author of that document and when it was published. I argue that these variables are best suited to capture temporal and publisher level variation in the documents.

---

<sup>14</sup>See also Griffiths and Steyvers (2002), Griffiths and Steyvers (2004) and Hofmann (1999)



In general inference of mixed-membership models, such as the one applied in this paper, has been a thread of research in applied statistics (Blei, Ng, and Jordan 2003, @erosheva\_mixed-membership\_2004, @braun\_variational\_2010). Topic models are usually imprecise as the function to be optimized has multiple modes, such that the model results can be sensitive to the starting values (e.g. the number of topics and the covariates influencing the prior distributions). Since an ex ante valuation of a model is hardly possible, I compute a variety of different models and compare their posterior probability. This enables me to check how results vary for different model solution (M. Roberts, Stewart, and Tingley 2016a). I then cross-checked some subset of assigned topic distributions to evaluate whether the estimates align with the concept of interest (Gentzkow, Kelly, and Taddy 2017). These manual audits are applied together with numeric optimization based on the topic coherence measure suggested by (Mimno et al. 2011).

```
## stm(documents = news_df_sparse, K = 40, prevalence = ~source +
##      s(year_biweek), data = covariates, init.type = "Spectral")
```

## Results of the STM

As stated above, the generative process of the STM results in a topic distribution  $\theta_d$  for each document  $d$  over all topics  $k$ . The average of each topic across all documents results in the expected probability of a topic across the whole corpus. Figure 7 displays the top 20 topics ordered by their average probability over the whole corpus. Since every topic is a probability distribution over words, top words may help to understand what each topic is about and are used as labels in Figure 7.<sup>15</sup> However, since those most probable words not necessarily the most exclusive words and they only represent a small fraction of the probability distribution, interpretation should be done very cautiously.

- 1: Topic 35 is the most common topic across the corpus
- 4: topic 38 about refugees
- 5: topic 39 about the jamaica coalition
- 6: followed by topic 8 about the “GroKo”



Figure 7: Top 20 topics by prevalence

Each document has a probability distribution over all topics. Figure 8 illustrates the topic distribution of two news paper articles and Figure 9 does the same for two press releases randomly chosen from the

<sup>15</sup>Table 10 gives an overview of the most probable terms for each topic.

corpus.<sup>16</sup>Figure 10 shows the distribution of mean topic probability of all topics for news articles and press releases before and after the election date.

Observations: - High peak both before and after election for topic 35 & 20 - no big difference between pre and post election for press releases, but for news articles

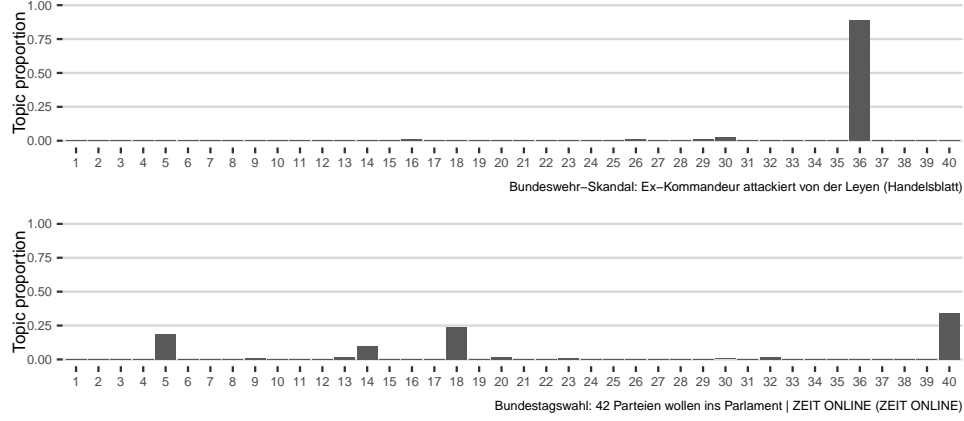


Figure 8: Topic probability of sample news articles

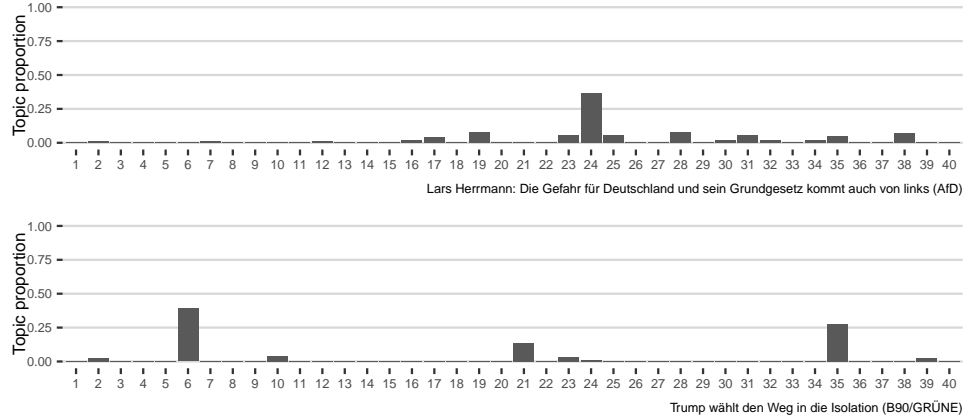


Figure 9: Topic probability of sample press releases

## Document Cosine similarity

Next, the cosine similarity measure is used in order to compare the retrieved topic distribution of every document. Cosine similarity is built on the geometric definition of the dot product of two vectors. It is a measure for the distance between two vectors and is defined between zero and one; values towards 1 indicate similarity.

$$\text{cosine similarity} = \cos(\theta) = \frac{a * b}{||a|| ||b||}$$

<sup>16</sup>Translations news articles: (1) Bundeswehr scandal: ex-commander attacks Von Der Leyen. (2) Bundestag elections: 42 parties want to be elected to parliament. Translations press releases: (1) Lars Herrmann: The danger for Germany and its Basic Law is also coming from the left. (2) Trump chooses the path to isolation.

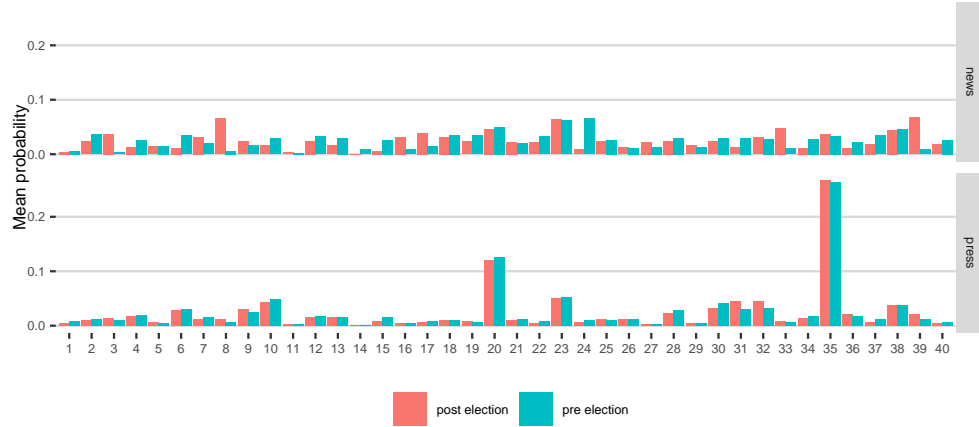


Figure 10: Mean topic prevalence

As topic proportions per document are vectors of the same length, the cosine similarity allows a comparison of the topic distribution between two documents.<sup>17</sup>

For example, Table 2 displays the most similar documents of the document with the title *Glyphosat-Streit vergiftet GroKo-Verhandlungen: Was wusste Merkel?* (Bild.de).

Table 2: Most similar documents

title	source	cos_sim
Alice Weidel: Undemokratisches Ausgrenzen der AfD schadet dem Amt des Bundespräsidenten	AFD	0.729
Fakten müssen für Zulassung von Pflanzenschutzmitteln maßgeblich sein	FDP	0.474
Dialog mit Russland garantiert Sicherheit in Europa	DIE LINKE	0.415
Keine Verlängerung für Glyphosatanwendung	SPD	0.31
Georg Pazderski: Rotrotgrün auf Bundesebene keine Option	AFD	0.163

In the next step, for each news paper, the cosine similarity between all topic-document distribution pairs between the news papers articles and the press releases is calculated if that press releases was published in a range of +/- 4 days from the publication date of the news article. This means the topic distribution of news article 1 is compared to press release 1, press release 2, 3 and so on as long as the press releases was published within 7 days before the news article.

- Only pairs where both documents are published either before or after the election date.
- See Table 3 for an illustration of the data

Table 3: Dataset structure step 1

title1	title2	cosine_sim	source1	source2	date1	date2
Mansfeld: Das sin...	Bundestagswahlkam...	0.02	DIE WELT	AfD	2017-09-17	4
Wir müssen reden:...	Pkw-Maut ist sch...	0.3	DIE WELT	FDP	2017-08-31	3
Überwachung: Die ...	SPD-Rentenkonzept...	0.02	DIE WELT	DIE LINKE	2017-06-10	3
Steuerpolitik: SP...	Medienpolitischer...	0.07	DIE WELT	SPD	2017-07-04	4
SPD-General: "Ein...	Gesundheitssyste...	0.21	DIE WELT	CDU	2018-02-11	3

Next, the mean cosine similarity for each date(1) and source(2) is estimated to obtain the final data frame (see Table 4)

<sup>17</sup>For applications of cosine similarity to compare of topic model outcomes see e.g. Rehs (2020) and Ramage, Dumais, and Liebling (2010)

Table 4: Final dataset structure

date1	source1	source2	cos_sim
2017-11-29	DIE WELT	CDU	0.16
2017-09-21	DIE WELT	FDP	0.1
2017-06-28	DIE WELT	CDU	0.1
2017-07-04	DIE WELT	B90/GRÜNE	0.16
2017-08-26	DIE WELT	FDP	0.15

## Model estimations

Figure 11 plots the data points for each news paper. The smoothed line shows the conditional mean for each party (cosine similarity between the news articles and each parties press releases conditional on the day). This helps to analyze the trend.

Observations:

- at the beginning, most news papers (except for Handelsblatt) have the highest similarity with the AfD.
- this similarity is diminishing over the course of time. However, for DIE WELT & Focus Online it is still the highest at the end of the time span.
- stern.de, Bild.de, ZEIT ONLINE & SPIEGEL ONLINE: Similarity with AfD is lower at the end of the time period

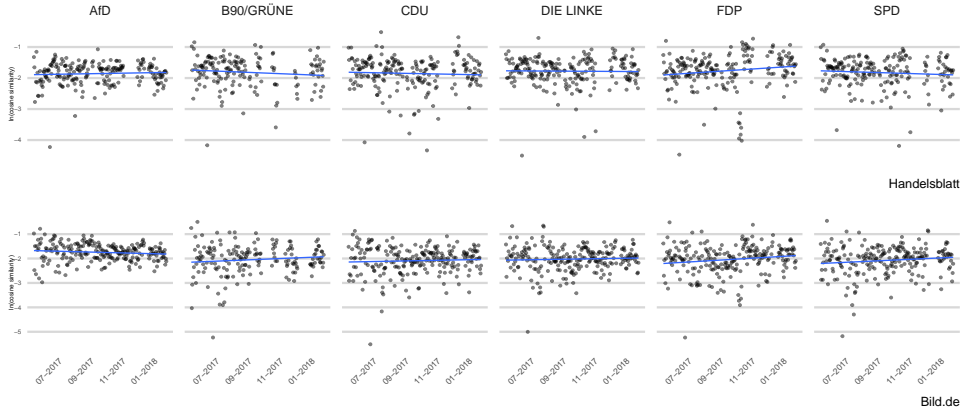


Figure 11: Log of daily mean cosine similarity between topic distributions in newspaper/press articles pairs

## OLS dummy regression

- for each news publisher, a simple OLS regression is estimated, where the similarity score ( $\ln(CS_i)$ ) on day  $i$  between the news articles of that news paper and the press releases of a political party  $k$  is the dependent variable and the political party dummies are the independent variables.

$$\ln(CS_i) = \beta_0 + \beta_n D_{i,k-1} + \epsilon_i$$

- $i$  = date (date1),  $k$  = political party (source2)

Table 5:

	Dependent variable:					
	DIE WELT (1)	stern.de (2)	ZEIT ONLINE (3)	Cosine similarity of topic distribution Handelsblatt (4)	FOCUS Online (5)	
source2B90/GRÜNE	-0.224*** (0.036)	-0.230*** (0.044)	-0.216*** (0.057)	0.044 (0.050)	-0.287*** (0.038)	-0.3
source2CDU	-0.268*** (0.033)	-0.200*** (0.042)	-0.277*** (0.054)	0.006 (0.047)	-0.295*** (0.036)	-0.3
source2DIE LINKE	-0.269*** (0.033)	-0.142*** (0.042)	-0.195*** (0.053)	0.079* (0.047)	-0.294*** (0.036)	-0.2
source2FDP	-0.243*** (0.033)	-0.134*** (0.042)	-0.195*** (0.054)	0.088* (0.047)	-0.280*** (0.036)	-0.2
source2SPD	-0.262*** (0.034)	-0.179*** (0.042)	-0.222*** (0.054)	0.032 (0.048)	-0.299*** (0.036)	-0.3
Constant	-1.765*** (0.024)	-1.990*** (0.029)	-1.766*** (0.038)	-1.858*** (0.033)	-1.829*** (0.025)	-1.7
Observations	1,364	1,370	1,395	1,197	1,419	
R <sup>2</sup>	0.068	0.025	0.022	0.005	0.072	
Adjusted R <sup>2</sup>	0.064	0.022	0.019	0.001	0.068	
Residual Std. Error	0.363 (df = 1358)	0.455 (df = 1364)	0.591 (df = 1389)	0.483 (df = 1191)	0.400 (df = 1413)	0.55
F Statistic	19.742*** (df = 5; 1358)	7.113*** (df = 5; 1364)	6.272*** (df = 5; 1389)	1.195 (df = 5; 1191)	21.824*** (df = 5; 1413)	12.411*

Note:

Interpretation:

When  $D$  switches from from 0 (AfD) to 1 (any other party), the % impact of  $D$  on  $Y$  can be estimated as  $\exp(\beta) - 1$

Example (B90/Gruene):

```
## source2B90/GRÜNE
## -0.2010668
```

- Compared to AfD, the similarity of topics of DIE WELT and B90/G is 20% lower.
- In general, the document similarity is significantly less in for all news papers except for Handelsblatt, where no significant difference can be found between the parties.

## Regression discontinuity model

- We assume that the topic similarity changes after the election day:

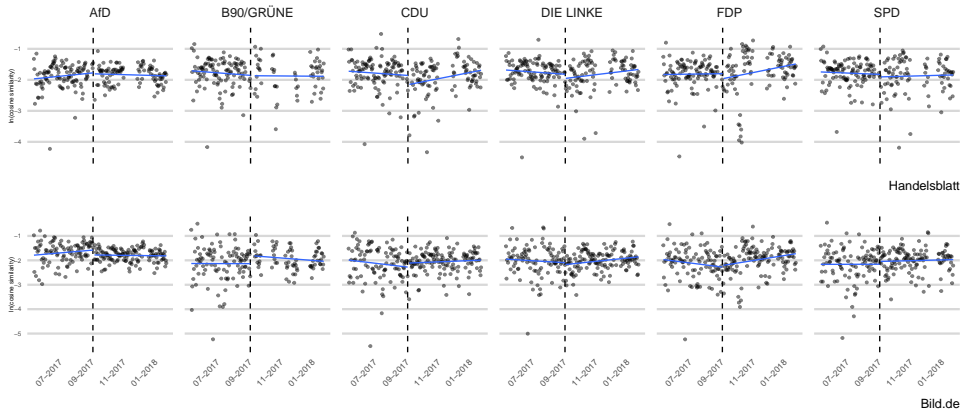


Figure 12: Log of mean cosine similarity between topic distributions in newspaper/press articles pairs - with cutoff value

- Initially introduced by Thistlethwaite and Campbell (1960) & formalized by Hahn, Todd, and Van der Klaauw (2001)

- It's a rigorous nonexperimental approach used to estimate treatment effects, where treatment is determined by whether an observed assignment variable (aka "forcing" or "running" variable) exceeds a known cutoff point.
- In our model, the running variable  $W$  is the time difference between date  $i$  and the election date (in days) and the cutoff point  $c$  is the election day
- The idea of regression discontinuity design is to use observations with a  $W_i$  close to  $c$  for the estimation of  $\beta_1$
- $\beta_1$  is the average treatment effect for observations with  $W_i = c$  which is assumed to be a good approximation to the overall treatment effect.
- In other words,  $\beta_1$  gives us the average change of the similarity between the content of news articles and press releases after the election day
- Including interaction terms  $T_i Dn_j$  in the second step allows us to estimate the treatment effect for each party  $\rightarrow$  since we are estimating an isolated model for each newspaper, we have the average treatment effect of the election day for each newspaper/party pair.

Calculate a regression discontinuity model for each newspaper  $i$ .

$$\ln(\text{CS}_i) = \beta_0 + \beta_1 T_i + \beta_2 W_i + \beta_n D_{i,k-1} + \epsilon_i$$

where

$$T_i = 1 \text{ if date } \geq \text{election date} \quad T_i = 0 \text{ if date } < \text{election date}$$

so that the receipt of treatment  $T_i$  is determined by the threshold  $c$  (election day) of the continuous (running) variable  $W_i$ .

### Without interaction terms

$$\text{CosineSimilarity}_i = \beta_0 + \beta_1 T_i + \beta_n Dn_j + \epsilon_i$$

Table 6:

	Dependent variable:				
	DIE WELT (1)	stern.de (2)	ZEIT ONLINE (3)	Cosine similarity of topic distribution Handelsblatt (4)	FOCUS Online (5)
treated	-0.115*** (0.040)	-0.099* (0.051)	-0.169*** (0.060)	-0.166*** (0.055)	-0.036 (0.042)
X_centered	0.0001 (0.0003)	0.001*** (0.0003)	-0.001** (0.0004)	0.001** (0.0004)	0.0001 (0.0003)
source2B90/GRÜNE	-0.237*** (0.035)	-0.229*** (0.044)	-0.250*** (0.056)	0.035 (0.050)	-0.290*** (0.039)
source2CDU	-0.269*** (0.033)	-0.199*** (0.042)	-0.278*** (0.052)	0.006 (0.047)	-0.296*** (0.036)
source2DIE LINKE	-0.268*** (0.033)	-0.141*** (0.041)	-0.195*** (0.052)	0.080* (0.047)	-0.294*** (0.036)
source2FDP	-0.243*** (0.033)	-0.134*** (0.042)	-0.197*** (0.052)	0.088* (0.047)	-0.280*** (0.036)
source2SPD	-0.264*** (0.033)	-0.177*** (0.042)	-0.228*** (0.052)	0.033 (0.048)	-0.299*** (0.036)
Constant	-1.706*** (0.030)	-1.949*** (0.037)	-1.665*** (0.047)	-1.780*** (0.042)	-1.811*** (0.032)
Observations	1,364	1,370	1,395	1,197	1,419
R <sup>2</sup>	0.085	0.033	0.078	0.012	0.072
Adjusted R <sup>2</sup>	0.080	0.028	0.073	0.007	0.068
Residual Std. Error	0.360 (df = 1356)	0.454 (df = 1362)	0.575 (df = 1387)	0.481 (df = 1189)	0.400 (df = 1411)
F Statistic	18.009*** (df = 7; 1356)	6.642*** (df = 7; 1362)	16.650*** (df = 7; 1387)	2.149** (df = 7; 1189)	15.738*** (df = 7; 1411)

Note:

### With interaction terms

The interaction term  $T_i * Dn$  means that the slope can vary on either side of the treatment threshold for each party.

- The coefficient  $\beta_1$  is how the intercept jumps (the RDD effect)
- $\beta_3$  is how the slope changes for each party

Table 7:

	Dependent variable:				
	DIE WELT (1)	stern.de (2)	ZEIT ONLINE (3)	Cosine similarity of topic distribution Handelsblatt (4)	FOCUS Online (5)
treated	-0.134** (0.058)	-0.065 (0.073)	-0.118 (0.090)	-0.078 (0.082)	-0.075 (0.082)
X_centered	0.0001 (0.0003)	0.001*** (0.0003)	-0.001** (0.0004)	0.001** (0.0004)	0.0001 (0.0003)
source2B90/GRÜNE	-0.267*** (0.048)	-0.216*** (0.060)	-0.180** (0.077)	0.105 (0.066)	-0.338*** (0.048)
source2CDU	-0.259*** (0.048)	-0.115* (0.060)	-0.232*** (0.077)	0.098 (0.066)	-0.288*** (0.048)
source2DIE LINKE	-0.251*** (0.048)	-0.109* (0.060)	-0.160** (0.077)	0.134** (0.066)	-0.270*** (0.048)
source2FDP	-0.291*** (0.048)	-0.195*** (0.060)	-0.201*** (0.077)	0.069 (0.066)	-0.355*** (0.048)
source2SPD	-0.278*** (0.048)	-0.143** (0.060)	-0.201*** (0.077)	0.098 (0.066)	-0.339*** (0.048)
treatedTRUE:source2B90/GRÜNE	0.069 (0.071)	-0.023 (0.090)	-0.151 (0.113)	-0.161 (0.103)	0.104 (0.071)
treatedTRUE:source2CDU	-0.020 (0.066)	-0.162* (0.083)	-0.086 (0.105)	-0.188** (0.094)	-0.014 (0.066)
treatedTRUE:source2DIE LINKE	-0.033 (0.066)	-0.063 (0.083)	-0.065 (0.104)	-0.110 (0.094)	-0.044 (0.066)
treatedTRUE:source2FDP	0.092 (0.066)	0.119 (0.083)	0.009 (0.105)	0.041 (0.094)	0.140* (0.066)
treatedTRUE:source2SPD	0.027 (0.067)	-0.066 (0.084)	-0.048 (0.105)	-0.135 (0.095)	0.074 (0.067)
Constant	-1.696*** (0.037)	-1.967*** (0.047)	-1.692*** (0.059)	-1.824*** (0.051)	-1.790*** (0.037)
Observations	1,364	1,370	1,395	1,197	1,419
R <sup>2</sup>	0.089	0.042	0.079	0.020	0.079
Adjusted R <sup>2</sup>	0.081	0.033	0.071	0.010	0.071
Residual Std. Error	0.360 (df = 1351)	0.453 (df = 1357)	0.575 (df = 1382)	0.480 (df = 1184)	0.399 (df = 1351)
F Statistic	10.955*** (df = 12; 1351)	4.925*** (df = 12; 1357)	9.927*** (df = 12; 1382)	2.018** (df = 12; 1184)	10.043*** (df = 12; 1351)

Note:

## Annex

Table 8: Online sources for press releases

	Party	Parliamentary Group
CDU	cdu.de	presseportal.de
SPD	spd.de	spdfraktion.de
FDP	fdp.de	fdpbt.de
B90/Die Grünen	gruene.de	gruene-bundestag.de
DIE LINKE	die-linke.de	die-linke.de/start/presse/aus-dem-bundestag
AfD	afd.de	afdbundestag.de

Table 9: Summary statistics of text length

source	n	mean	sd	median	min	max
AfD	523	212.83	72.16	196	103	553
B90/GRÜNE	211	229.32	63.37	219	104	399
Bild.de	1303	476.07	318.28	398	121	3710
CDU	248	274.54	106.08	256	100	1030
DIE LINKE	686	200.47	71.78	190	101	1048
DIE WELT	3222	509.57	612.06	380	121	14507
FDP	301	161.9	83.78	144	100	999
FOCUS Online	2780	393.89	317.05	297.5	121	5647
Handelsblatt	2785	589.51	495.82	488	121	6899
SPD	315	213.41	56.16	208	103	429
SPIEGEL ONLINE	2089	539.09	415.05	413	121	3466
stern.de	2943	514.66	616.55	373	121	9287
ZEIT ONLINE	1351	513.75	387.14	459	121	8015

% Table created by stargazer v.5.2.2 by Marek Hlavac, Harvard University. E-mail: hlavac at fas.harvard.edu  
 % Date and time: Sun, Feb 21, 2021 - 14:53:05

## References

Bholat, David M., Stephen Hansen, Pedro M. Santos, and Cheryl Schonhardt-Bailey. 2015. “Text Mining for Central Banks.” *SSRN Electronic Journal*, June 15 [http://www.academia.edu/13430482/Text\\_mining\\_for\\_central\\_banks](http://www.academia.edu/13430482/Text_mining_for_central_banks).

Blassnig, Sina, Sven Engesser, Nicole Ernst, and Frank Esser. 2019. “Hitting a Nerve: Populist News Articles Lead to More Frequent and More Populist Reader Comments.” *Political Communication*. August.



Table 10: 7 most probable terms per topic

	Top Terms
1	a, the, s, of, u, brexit, großbritannien
2	merkel, angela, kanzlerin, bundeskanzlerin, cdu, merkels, deutschland
3	spd, union, cdu, csu, koalitionsvertrag, koalitionsverhandlungen, schulz
4	afd, weidel, gauland, alice, alexander, politiker, äusserungen
5	stimmen, wahlkreis, kandidaten, afd, wahl, gewählt, fdp
6	trump, us, usa, deutschland, präsident, donald, berlin
7	cdu, union, peter, politiker, spahn, altmaier, schäuble
8	spd, koalition, union, groko, große, koalitionsverhandlungen, parteitag
9	afd, partei, sachsen, gauland, parteien, pazderski, höcke
10	diesel, unternehmen, deutschland, autos, deutschen, industrie, fahrverbote
11	ge, ten, be, le, ver, lambsdorff, te
12	gericht, prozess, urteil, richter, staatsanwaltschaft, verfahren, jahre
13	berlin, deutschen, osten, o, tag, jahr, millionen
14	august, cdu, spd, prozent, bundestagswahl, wahl, parteien
15	kohl, helmut, kohls, einheit, kanzler, tod, deutschen
16	spd, nahles, andrea, partei, scholz, schulz, schwesig
17	csu, seehofer, horst, söder, obergrenze, bayern, chef
18	prozent, umfrage, spd, union, fdp, cdu, afd
19	polizei, stadt, menschen, polizisten, täter, verletzt, angaben
20	euro, milliarden, jahr, millionen, prozent, bund, geld
21	grünen, linke, linken, özdemir, partei, wagenknecht, göring
22	cdu, niedersachsen, spd, grünen, rot, fdp, landtag
23	welt, politik, menschen, jahren, lange, frage, fragen
24	g, hamburg, gipfel, polizei, hamburger, demonstranten, scholz
25	deutschland, is, verfassungsschutz, syrien, gefährder, islamisten, staat
26	steinmeier, schmidt, russland, frank, bundespräsident, glyphosat, walter
27	afd, petry, partei, fraktion, frauhe, meuthen, gauland
28	berliner, berlin, amri, maizière, innenminister, behörden, daten
29	gabriel, sigmar, außenminister, spd, schröder, amt, gerhard
30	bundestag, spd, abgeordneten, abgeordnete, parlament, abstimmung, fraktion
31	türkei, erdoğan, türkischen, deutschland, bundesregierung, türkische, deutsche
32	frauen, deutschland, kinder, studie, eltern, muslimen, antisemitismus
33	fdp, jamaika, lindner, koalition, neuwahlen, spd, grünen
34	facebook, maas, twitter, gesetz, internet, netz, heiko
35	eu, deutschland, europa, bundesregierung, europäischen, deutschen, menschen
36	bundeswehr, soldaten, leyn, nato, ursula, einsatz, verteidigungsministerin
37	schulz, spd, martin, kanzlerkandidat, wahlkampf, bundestagswahl, partei
38	flüchtlinge, deutschland, menschen, zahl, flüchtlingen, familiennachzug, jahr
39	fdp, grünen, jamaika, csu, union, grüne, cdu
40	bundestagswahl, afd, wahl, prozent, partei, bundestag, parteien

- Druckman, James N., and Michael Parkin. 2005. "The Impact of Media Bias: How Editorial Slant Affects Voters." *The Journal of Politics* 67 (4): 1030–49. <https://doi.org/10.1111/j.1468-2508.2005.00349.x>.
- Eberl, Jakob-Moritz. 2018. "Lying Press: Three Levels of Perceived Media Bias and Their Relationship with Political Preferences." *Communications*, March. <https://doi.org/10.1515/commun-2018-0002>.
- Eberl, Jakob-Moritz, Hajo G. Boomgaarden, and Markus Wagner. 2017. "One Bias Fits All? Three Types of Media Bias and Their Effects on Party Preferences." *Communication Research* 44 (8): 1125–48. <https://doi.org/10.1177/0093650215614364>.
- Erosheva, Elena, Stephen Fienberg, and John Lafferty. 2004. "Mixed-Membership Models of Scientific Publications." *Proceedings of the National Academy of Sciences* 101 (suppl 1): 5220–7. <https://doi.org/10.1073/pnas.0307760101>.
- Gentzkow, Matthew A., and Jesse M. Shapiro. 2004. "Media, Education and Anti-Americanism in the Muslim World." *Journal of Economic Perspectives* 18 (3): 117–33. <https://doi.org/10.1257/0895330042162313>.
- Gentzkow, Matthew, Bryan T. Kelly, and Matt Taddy. 2017. "Text as Data." Working Paper 23276. National Bureau of Economic Research. <https://doi.org/10.3386/w23276>.
- Griffiths, Thomas L., and Mark Steyvers. 2002. "A Probabilistic Approach to Semantic Representation." *Proceedings of the Annual Meeting of the Cognitive Science Society* 24 (24). <https://escholarship.org/uc/item/44x9v7m7>.
- . 2004. "Finding Scientific Topics." *Proceedings of the National Academy of Sciences* 101 (suppl 1): 5228–35. <https://doi.org/10.1073/pnas.0307752101>.
- Grimmer, Justin, and Brandon Stewart. 2013. "Text as Data: The Promise and Pitfalls of Automatic Content Analysis Methods for Political Texts." *Political Analysis* 21: 267–97.
- Groseclose, Tim, and Jeffrey Milyo. 2005. "A Measure of Media Bias." *The Quarterly Journal of Economics* 120 (4): 1191–1237. <https://www.jstor.org/stable/25098770>.
- Hahn, Jinyong, Petra Todd, and Wilbert Van der Klaauw. 2001. "Identification and Estimation of Treatment Effects with a Regression-Discontinuity Design." *Econometrica* 69 (1): 201–9. <https://www.jstor.org/stable/2692190>.
- Hofmann, Thomas. 1999. "Probabilistic Latent Semantic Indexing." In *Proceedings of the 22Nd Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, 50–57. SIGIR '99. New York, NY, USA: ACM. <https://doi.org/10.1145/312624.312649>.
- Kepplinger, Hans Mathias, and Marcus Maurer. 2004. "Der Einfluss Der Pressemitteilungen Der Bundesparteien Auf Die Berichterstattung Im Bundestagswahlkampf 2002." In *Quo Vadis Public Relations? Auf Dem Weg Zum Kommunikationsmanagement: Bestandsaufnahmen Und Entwicklungen*, edited by Juliana Raupp and Joachim Klewes, 113–24. Wiesbaden: VS Verlag für Sozialwissenschaften. [https://doi.org/10.1007/978-3-322-83381-5\\_9](https://doi.org/10.1007/978-3-322-83381-5_9).
- Lott, John R., and Kevin A. Hassett. 2014. "Is Newspaper Coverage of Economic Events Politically Biased?" *Public Choice* 160 (1): 65–108. <https://doi.org/10.1007/s11127-014-0171-5>.
- Mimno, David, Hanna M. Wallach, Edmund Talley, Miriam Leenders, and Andrew McCallum. 2011. "Optimizing Semantic Coherence in Topic Models." In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, 262–72. EMNLP '11. Stroudsburg, PA, USA: Association for Computational Linguistics. <http://dl.acm.org/citation.cfm?id=2145432.2145462>.
- Newman, David, Jey Han Lau, Karl Grieser, and Timothy Baldwin. 2010. "Automatic Evaluation of Topic Coherence." In *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, 100–108. HLT '10. Stroudsburg, PA, USA: Association for Computational Linguistics. <http://dl.acm.org/citation.cfm?id=1857999.1858011>.
- Newman, Nic, Richard Fletcher, Antonis Kalogeropoulos, David Levy, and Rasmus Kleis Nielsen. 2018. "Reuters Institute Digital News Report 2018." Reuters Institute for the Study of Journalism. <http://media.digitalnewsreport.org/wp-content/uploads/2018/06/digital-news-report-2018.pdf?x89475>.

- Ramage, Daniel, Susan Dumais, and Daniel Liebling. 2010. *Characterizing Microblogs with Topic Models*.
- Rehs, Andreas. 2020. “A Structural Topic Model Approach to Scientific Reorientation of Economics and Chemistry After German Reunification.” *Scientometrics* 125 (2): 1229–51. <https://doi.org/10.1007/s11192-020-03640-0>.
- Roberts, Margaret E., Brandon M. Stewart, and Edoardo M. Airolidi. 2016. “A Model of Text for Experimentation in the Social Sciences.” *Journal of the American Statistical Association* 111 (515): 988–1003. <https://doi.org/10.1080/01621459.2016.1141684>.
- Roberts, Margaret, Brandon Stewart, and Dustin Tingley. 2016a. “Navigating the Local Modes of Big Data: The Case of Topic Models.” In *Computational Social Science: Discovery and Prediction*. New York: Cambridge University Press.
- . 2016b. “Stm: R Package for Structural Topic Models.” *Journal of Statistical Software* forthcoming (December).
- Strömbäck, Jesper. 2008. “Four Phases of Mediatization: An Analysis of the Mediatization of Politics.” *The International Journal of Press/Politics* 13 (3): 228–46. <https://doi.org/10.1177/1940161208319097>.
- Taddy, Matt. 2012. “On Estimation and Selection for Topic Models.” In *Proceedings of the 15th International Conference on Artificial Intelligence and Statistics*.
- Takens, Janet, Wouter Atteveldt, Anita van Hoof, and Jan Kleinnijenhuis. 2013. “Media Logic in Election Campaign Coverage.” *European Journal of Communication* 28 (June): 277–93. <https://doi.org/10.1177/0267323113478522>.
- Thistlethwaite, Donald L., and Donald T. Campbell. 1960. “Regression-Discontinuity Analysis: An Alternative to the Ex Post Facto Experiment.” *Journal of Educational Psychology* 51 (6): 309–17. <https://doi.org/10.1037/h0044319>.
- Wallach, Hanna M., David M. Mimno, and Andrew McCallum. 2009. “Rethinking LDA: Why Priors Matter.” In *Advances in Neural Information Processing Systems 22*, edited by Y. Bengio, D. Schuurmans, J. D. Lafferty, C. K. I. Williams, and A. Culotta, 1973–81. Curran Associates, Inc. <http://papers.nips.cc/paper/3854-rethinking-lda-why-priors-matter.pdf>.