# Biased reporting by the German media?

## An Analysis of Political News Coverage in Germany for the 2017 Bundestag Election

Franziska Löw[*]

October 2021

**Abstract**

The dynamics of online news and policy outcomes have been of great interest in several research areas in recent years. This paper provides a new method to estimate media bias using a structural topic model and cosine similarity to test slanting toward different political actors. For the empirical analysis, the content of German online newspapers and press releases of German parties during the election campaign before the federal election in 2017 is analyzed. Following the assumption that a) potential media bias is demand-driven and b) election results can be used as a proxy for reader beliefs, the results show that news articles of most newspapers slant towards AfD topics. Furthermore, we find evidence for the hypothesis that the election day results in changes in news coverage since newspapers can observe the true beliefs of readers.

[*]Institut für Industrieökonomik, Helmut Schmidt Universität. Email: franzi@localyzeapp.com

# Contents

# List of Tables

# List of Figures

# 1 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 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 [jandinter_wahlnachlese_]. After the 2017 federal elections in Germany, for example, the media were accused of contributing to the success of the right-wing populist AfD[1] by increasingly including the party's content and using the same language in their articles as the AfD. On the other hand, supporters of the AfD accuse the media of not covering their topics to a sufficient extent. Representatives of these media houses strongly opposed both accusations, claiming balanced reporting. The purpose of this study is to examine whether there is evidence that supports the allegation of biased media reporting in either direction, especially during election campaigns.

Economic literature examines both the supply and demand sides as factors driving media bias. In the former case, bias reflects the preferences of editors, owners (Besley and Prat 2006), or journalists (Baron 2006). On the other hand, bias may be driven by the demand side reflecting the profit-maximizing decision of news providers to satisfy consumer preferences. Advertising-financed media like online news, which offer their content to a large extent free of charge and generate revenue through advertising space, compete for readers' attention. Readers pay a non-monetary price providing their attention, which the media platform bundles and sells to advertising customers. This business model corresponds to that of a platform market. News outlets act as platforms that connect the advertising market with the reader market to exploit the indirect network effects between them (Dewenter and Rösch 2014). Therefore, a profit-maximizing publisher directs its economic decisions according to what will attract the most attention. In the traditional conception of the demand for news, where readers value the accuracy of the information, the market forces news outlets to deliver more accurate information. M. Gentzkow and Shapiro (2006) showed that increased competition among newspapers can reduce bias. Their setup assumes that newspapers want to build a reputation as providers of accurate information and that bayesian consumers base their beliefs about information quality on past reports. As a result, low-quality firms are incentivized to ignore signals that contradict prior common expectations. Although this information is valuable to readers, it also reveals that its sources are low quality.[2]

This logic of a rational reader that simply values the accuracy of information differs from noneconomic media studies. Instead, communication literature suggests that readers prefer news consistent with their beliefs (Graber 1984). These beliefs might come from different sources, like education, previous news, or views of politicians or political parties they trust. Especially during election campaigns, competing political actors attempt to generate support by presenting their viewpoints and defining the issue-based criteria on which voters will evaluate them (Eberl, Boomgaarden, and Wagner 2017). Parties instrumentalize their public

---

[1]Alternative für Deutschland (AfD) is a right-wing populist political party in Germany established in April 2013 *The Economist* (n.d.)

[2]See Prat and Strömberg (2013) for a survey about economic literature on the topic of mass media.

relations to highlight issues they perceive as competent on, that they "own", and are essential to their voters (Kepplinger and Maurer 2004). News outlets will try to attract the same audiences by adjusting their news content if the political actor can generate enough interest. Following that explanation of confirmation bias, Mullainathan and Shleifer (2005) show that heterogeneous reader beliefs incentivize news outlets in a competitive market to segment the market and slant towards extreme positions, generating, in the aggregate, an unbiased media landscape. Conversely, on topics where readers share common beliefs, competition among news outlets results in slanting towards reader biases, enforcing a biased media landscape. They use a standard Hotelling model with quadratic transportation costs, where the transportation cost is interpreted as the ideological distance between a reader and a newspaper. If reader beliefs are homogeneous, the monopoly and the duopoly result in the same bias. In the case of heterogeneous preferences, competition lead to market segmentation through extremely biased news because market participants want to avoid price competition. The underlying assumption for this model is that the payoff for readers depends on both the quality of information and how well the information corresponds to their prior beliefs.

The assumption that the distribution of bias in the population is the primary driver of bias is consistent with the concept of framing from the communication literature. The central argument is that newspapers tend to select frames people like to hear. Another important concept in communication studies to explain the emergence of bias is the entertainment factor of news (Takens et al. 2013). The underlying thesis is that political news content produces news values and narrative techniques that media use to attract audiences, i.e., the factors that turn an event into news worth reporting like conflict, drama, negativity, surprise, or proximity Blassnig et al. (2019). 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. 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. Blassnig et al. (2019) show that populist key messages by political and media actors in news articles provoke more reader comments. Therefore, new outlets competing for readers' attention have an incentive to pick up on the key messages of these parties.

This paper uses the content of German online newspapers and press releases of major German parties to analyze whether online news equally addresses the topics covered in these press releases during the election campaign for the federal elections in Germany in 2017. Furthermore, we analyze the effect of the election results on the news content. The interpretation of the results is based on the assumptions that a) media bias is demand-driven and b) that election results can be used as a proxy for reader beliefs. The results show that news articles of most newspapers slant towards AfD topics during the election campaign. Based on the studies discussed above, the cause for this bias could emerge from homogeneous reader beliefs or the fact that the entertainment factor of topics covered by AfD press releases is higher compared to other parties or both. Next, we test the hypothesis that the election day results in changes in news coverage since newspapers can observe the true beliefs of readers. Here, the results indicate that some newspapers adjust their content towards the election results. Although this paper does not estimate the cause for media bias, it provides a new

method to test slanting towards specific topics. Using the election day as an external event allows to understand the effect of the election results on the news content.

To answer these and other media-related questions in the political context, quantifying media content is a prerequisite. One of the critical challenges is determining the features used to describe media content - audio, video, or text content. Studies that rely on quantifying media content for their analyses use, for example, visibility (how often political actors appear in the media (Lengauer and Johann 2013)) or tonality (how they are evaluated (Eberl, Boomgaarden, and Wagner 2017)). Other studies examine the topics discussed or the language used in the media to identify whether political actors can place their 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 (M. A. 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 parties' press releases in the German "Bundestag" to measure the "slant" of these newspapers towards a political party. The structural topic model (STM) developed by M. E. Roberts, Stewart, and Airoldi (2016) is applied to discover the latent topics in the corpus of text data (see 3.2 Structural topic model). This probabilistic text model results in a probability distribution for each document across all topics, which is then aggregated to calculate the degree of similarity between the news articles of different news providers and the parties' press releases[3] (see 3.3 Similarity measure). This similarity measure is then used to examine the above research questions using a regression model in 3.4 Model estimations. Prior to a more detailed explanation and implementation of this empirical strategy in chapter 3 Empirical analysis, the following section provides an overview of the political situation surrounding the 2017 federal election.

This paper adds to the academic debate about media bias from economic and communication literature. Although the empirical approach does not estimate the cause for media bias, it provides a new method to test slanting towards specific topics using natural language processing tools. This method can be easily extended to similar use cases and data sets, allowing a new way of measuring media bias without the need of manual classification of text data.

---

[3]For the sake of simplicity, both news articles and press releases will be referred to as documents for the remaining of this paper.

# 2 The political situation in Germany

The articles analyzed in this paper cover a period from June 1, 2017, to March 1, 2018, and thus cover both the most crucial 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)[4], 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 profile. In reporting on the federal elections, leading party members of the AfD and party supporters repeatedly accused the mass media of reporting unilaterally and intentionally presenting the AfD badly.

After the election, forming 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/GRÜNE (Bündnis 90/Die Grünen, green party)). The SPD initially rejected the grand coalition. However, 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/GRÜNE 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, the media discussed possible re-election or a minority government as alternatives before the SPD decided to hold coalition talks with the CDU/CSU. This step provoked significant resistance from the party base, which called for a party-internal referendum on a grand coalition. However, after the party members voted in favor of the grand coalition, CDU/CSU and SPD formed a government 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.[5] 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.

---

[4]CDU/CSU, Union and CDU are used as synonyms in this paper for simplicity.

[5]The graph shows the moving average within 15 days from major German research institutes. Since the institutions do not all publish new values on the same days, the overall temporal accuracy is higher than the weekly accuracy. The data is scraped from https://www.wahlrecht.de/.

Figure 1: Election polls during the period under review

# 3 Empirical analysis

The empirical strategy used in this paper leverages the structure of the topic model framework, specifically the Structural Topic Model (STM), to generate topic distributions for each document which are then used to measure similarity between documents. The diagram below outlines the approach in more detail.



Figure 2: High level overview

In 3.1 Text pre-processing the text data is processed resulting in a matrix that represents a multi-dimensional space, where each dimension corresponds to a word in the document. Subsequently, in 3.2 Structural topic model this so-called document-term matrix is used as

input to calculate each document's topic distribution applying a STM. This, in turn, leads to a reduction in dimensionality in that each document is now represented as a distribution over the topics. These document-topic vectors are then used to calculate the cosine similarity between two documents as described in 3.3 Similarity measure. In the final section 3.4 Model estimations, this similarity measure is utilized as the dependent variable in a regression model with various specifications.

## 3.1 Text pre-processing

The analysis performed in this paper is based on a sample of 18,757 online news articles from seven German online news providers[6] and press releases of the seven parties that have been in the Bundestag since the 2017 federal elections[7]. Both news articles and press releases are dated from June 1, 2017 to March 1, 2018. The scraping code for both the news articles and press releases was written in R by the author if this paper. News articles were sc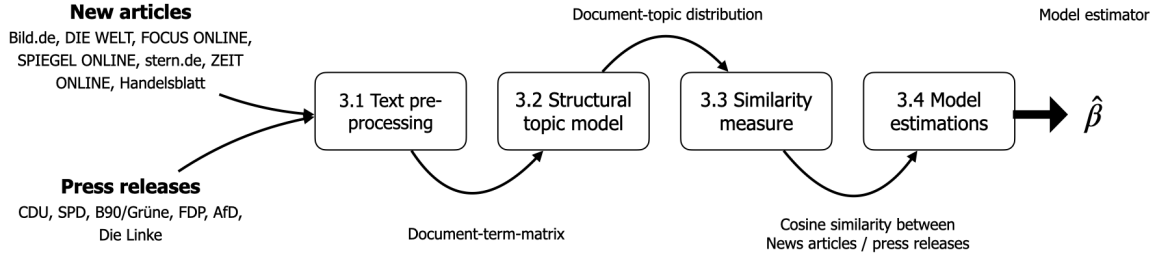raped from the Webhose.io API using the John Coene ([2018] 2019).[8] To consider only news about national politics, the articles were filtered based on their URL. The press releases were scraped from the public websites of the political parties and parliamentary groups.

As shown in Figure 3(a), except for Handelsblatt (position 53), these media outlets are among the top 30 German online news providers in the period under review in terms of visits.[9] The primary source of income for these privately managed media houses is digital advertising, even though paid content plays an increasingly important role. However, according to a survey on digital news by the Reuters Institute (N. 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[10] 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.[11] The results displayed in Figure 3(b) indicate that the media used for the analysis play a relevant role in their consumption.

Figure 4 shows the distribution of the number of articles by date and media outlet. There is a high peak around the federal elections on September 24 and another one shortly after the failure of the Jamaica coalition talks on November 19. The peak in July especially for stern.de is due to increased reporting about the G20 summit in Hamburg. Furthermore, Figure 4 shows that DIE WELT published the most articles on domestic policy, followed by stern.de, Handelsblatt and FOCUS ONLINE.

---

[6]Bild.de, DIE WELT, FOCUS ONLINE, SPIEGEL ONLINE, stern.de, ZEIT ONLINE, Handelsblatt

[7]CDU, SPD, B90/Grüne, FDP, AfD, Die Linke

[8]For more information see https://docs.webhose.io/reference#about-webhose.

[9]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.

[10]https://www.hans-bredow-institut.de/de/punctuationprojekte/reuters-institute-digital-news-survey

[11]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."

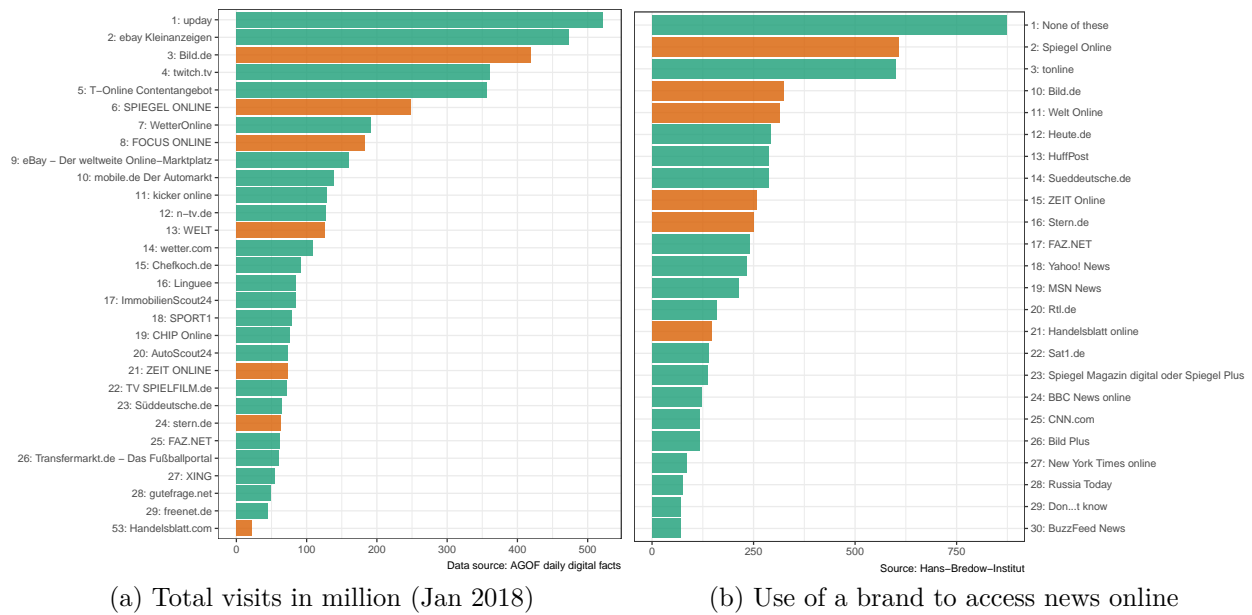(a) Total visits in million (Jan 2018)  (b) Use of a brand to access news online
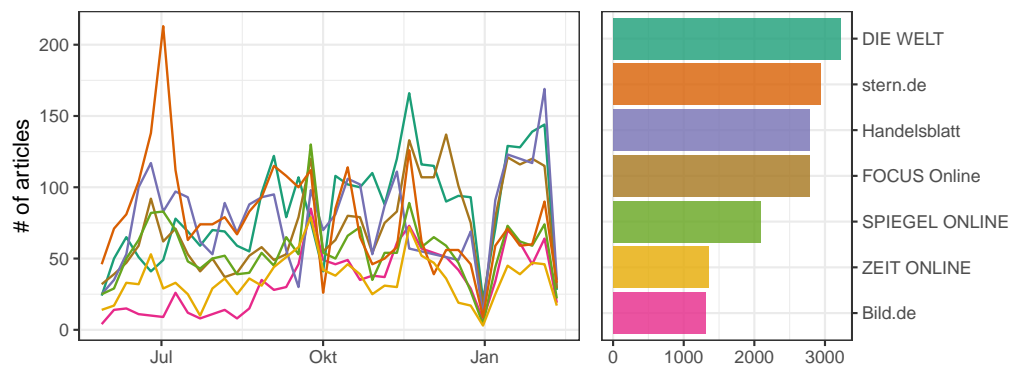
Figure 3: Selected german news brands



Figure 4: Distribution of news articles

Figure 5 shows that DIE LINKE published the most press releases in the period under review, followed by the AfD. Again, a peak can be discerned around the time of the G20 summit, especially in the press releases of DIE LINKE and SPD.



Figure 5: Distribution of press releases

Table 1 illustrates that, on average, news articles have a higher word count than the parties' press releases. While for news articles, the average is between 394 (FOCUS Online) and 590 (Handelsblatt), with press releases, the range is between 162 (FDP) and 275 (CDU). DIE WELT published the article with the most words (14.507) - the most extended press release has 1.048 words published by DIE LINKE.

Table 1: Summary statistics of word counts

| source | n | mean | sd | median | min | max |
|---|---|---|---|---|---|---|
| **News articles** | | | | | | |
| Bild.de | 1303 | 476.07 | 318.28 | 398.0 | 121 | 3710 |
| DIE WELT | 3222 | 509.57 | 612.06 | 380.0 | 121 | 14507 |
| FOCUS Online | 2780 | 393.89 | 317.05 | 297.5 | 121 | 5647 |
| Handelsblatt | 2785 | 589.51 | 495.82 | 488.0 | 121 | 6899 |
| SPIEGEL ONLINE | 2089 | 539.09 | 415.05 | 413.0 | 121 | 3466 |
| stern.de | 2943 | 514.66 | 616.55 | 373.0 | 121 | 9287 |
| ZEIT ONLINE | 1351 | 513.75 | 387.14 | 459.0 | 121 | 8015 |
| **Press releases** | | | | | | |
| AfD | 474 | 211.93 | 72.45 | 194.5 | 103 | 553 |
| B90/GRÜNE | 192 | 230.54 | 63.45 | 222.0 | 104 | 399 |
| CDU | 237 | 275.85 | 106.96 | 256.0 | 100 | 1030 |
| DIE LINKE | 631 | 200.36 | 70.66 | 190.0 | 101 | 1048 |
| FDP | 262 | 162.27 | 88.12 | 143.0 | 100 | 999 |
| SPD | 301 | 213.17 | 56.57 | 208.0 | 103 | 429 |

Several processing steps have to be performed to make the text quantifiable to use text as data input for statistical analyses. 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 (M. Gentzkow, Kelly, and Taddy (2017), Bholat et al. (2015)). Intuitively, the term frequency (tf) of a word measures how important that word may be for understanding the text. Word clouds are a commonly used visualization technique in text mining as they translate the tf into the size of the term in the cloud.

Words like "die," or "der" (eng. "the"), "and" (eng. "and"), and "ist" (eng. "is") are extremely common but unrelated to the quantity of interest. Often called stop words (M.

Gentzkow, Kelly, and Taddy 2017), these terms are essential to the grammatical structure but typically do not add any additional meaning and can be neglected. The predefined stop word list from the Snowball project[12] is used together with a customized, domain-specific list of words to identify and remove these distorting words. Additionally, punctuation characters (e.g. ., !, ?) and all numbers are removed from the data. The 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. We apply the widely used Porter-Stemmer algorithm based on a set of shortening rules applied to a word until it has a minimum number of syllables.[13]

As an example, the following word clouds represent the most frequent words of the pre-processed articles for Bild.de (Figure 6(a)) and press releases of AfD (Figure 6(b)). Thus, it becomes evident that these are texts discussing domestic policy issues. The SPD, in particular, seems to be highly frequent for Bild.de.



(a) Bild                    (b) AfD

Figure 6: Wordcloud after pre-processing

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. Next, for each document $d \in \{1, ..., D\}$ 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 $v \in \{1, ..., V\}$ and where $V$ is the number of unique terms. The $D$ x $V$ matrix $\boldsymbol{X}$ of all such counts is called the document-term matrix. Each row in this matrix represents a document, and each entry counts the occurrences of a unique term in that document. Table 2 provides

---

[12]http://snowball.tartarus.org/algorithms/german/stop.txt
[13]https://tartarus.org/martin/PorterStemmer/

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 (M. Gentzkow, Kelly, and Taddy 2017) since it disregards the words' order within a document.

Table 2: Document-term matrix - sample values

|  | wahlkreis | angelegenheit | einzutreten | eu | zunehmend | neuesten | widerspricht |
|---|---|---|---|---|---|---|---|
| 4149 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 679 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 7664 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 8389 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 8435 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
|  |  |  |  |  |  |  |  |
| 11293 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 9682 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5750 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| 17056 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 16762 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

## 3.2 Structural topic model

Next, a structural topic modeling (STM) developed by (M. E. Roberts, Stewart, and Airoldi 2016) is applied to discover the latent topics in the corpus of press releases and news articles. 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 is an extension of the standard topic modeling 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.[14]

The underlying idea for these models suggests that each topic $k$ potentially contains all of the unique terms within the vocabulary $V$ with a different probability. Therefore, each topic $k$ can be represented as a probability vector $\phi_k$ over all unique terms $V$. Simultaneously, each 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 its 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. At the same time, 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. $X$ gives the covariates for topic prevalence resulting in each document's

---

[14]See also Griffiths and Steyvers (2002), Griffiths and Steyvers (2004) and Hofmann (1999)

probability of a topic 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. Thus, the underlying data generating process to generate each word $w_{d,n}$ in 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 for topic models like LDA or STM is the number of topics ($K$) that occur over the entire corpus. Unfortunately, there is not a "right" answer to the number of appropriate topics 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). It is closely related to pointwise mutual information (D. Newman et al. 2010): it is maximized when the most probable words in a given topic are frequently used in a given topic co-occur together.

The function *searchK* from the *stm* package [stewart_bstewartstm_2021] supports the choice of the number of topics using several automated tests, including the average exclusivity and semantic coherence and 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, the author and bi-week dummies of a document are included as topical prevalence variables. In other words, we assume that the probability of a topic being included in a news article or a press release depends on the author and the publication date of that document. Therefore, we argue that these variables are best suited to capture temporal and publisher level variation in the documents.

In general, inference of mixed-membership models, such as the one applied in this paper, has been a thread of research in applied statistics Braun and McAuliffe (2010). However, topic models are usually imprecise as the function to be optimized has multiple modes so 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 is impossible, we compute various models and compare their posterior probability to evaluate how results vary for different model specifications (M. Roberts, Stewart, and Tingley 2016a). We then cross-checked some subset of assigned topic distributions to evaluate whether the estimates align with the concept of interest (M. Gentzkow, Kelly, and Taddy 2017).

### 3.2.1 Results of the STM

As mentioned in the previous section, the generative process of the STM results in a topic distribution $\theta_d$ for each document $d$ over all topics $k$ and a word distribution $\phi_k$ for each topic

over all terms in the vocabulary. Thus, the most probable words of each topic may help to understand the context of each topic.[15] However, since those most probable words are not necessarily the most exclusive words and only represent a small fraction of the probability distribution, interpretation should be made cautiously.

For the analysis, the topic distribution of each document is used to estimate the similarity of documents. Figure 7 illustrates such a topic distribution of two newspaper articles. The red numbers display the topic probability (for probabilities $>= 0.02$). News article 1[16] shows a definite distribution towards topic 36, for which terms like Bundeswehr, Soldaten (soldiers), Nato, Verteidigungsministerin (defense minister) are among the most probable words. News article 2[17] does not show such a clear tendency towards a single topic. However, for both topics with highest probability similar terms are among the top terms.

Similarly, Figure 8 illustrates the topic distribution for two press releases randomly chosen from the corpus. For press release 1[18], topic 24 is the most probable, containing terms about the G20 Summit, during which left-wing radicals caused considerable riots. Topic distribution of press article 2[19] shows peaks for topics 6, 21 and 35. The top terms of topic 6 contain the words trump, us, usa, deutschland (Germany), and präsident (president). Similarly, topic 35 seems to deal with German foreign policy since top terms include words like eu, deutschland (Germany), europa, and bundesregierung (Federal Government).
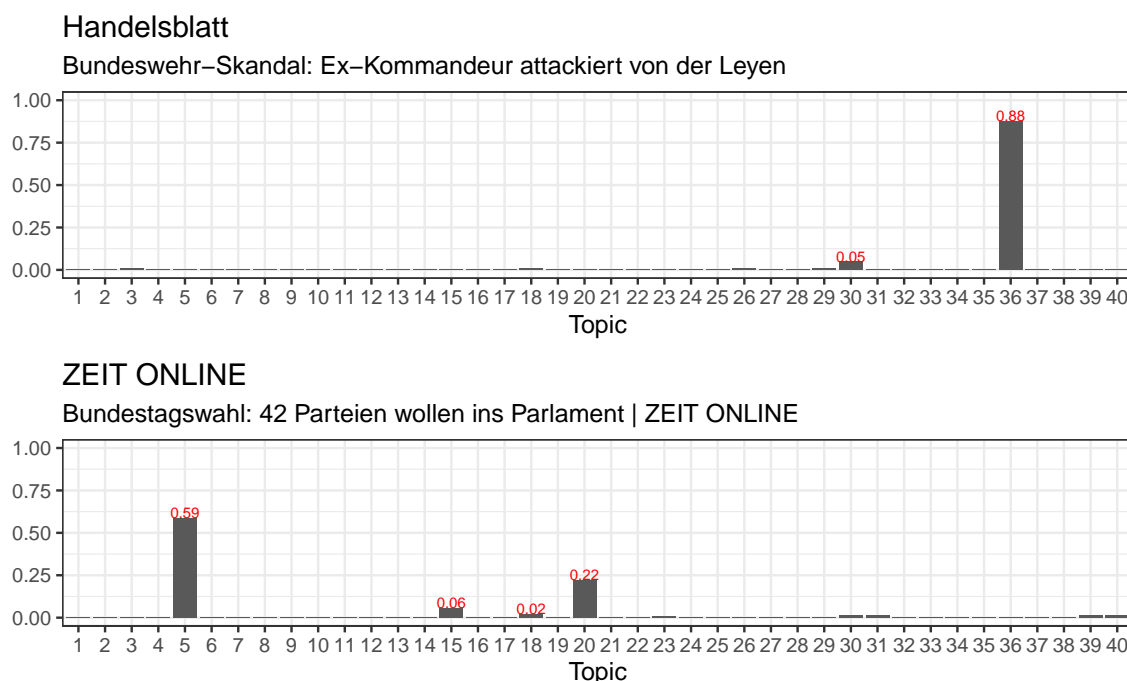


Figure 7: Topic probability of sample news articles

---

[15]Table 10 gives an overview of the most probable terms for each topic.

[16]Bundeswehr scandal: ex-commander attacks Von Der Leyen

[17]Bundestag elections: 42 parties want to be elected to parliament.

[18]Lars Herrmann: The danger for Germany and its Basic Law is also coming from the left

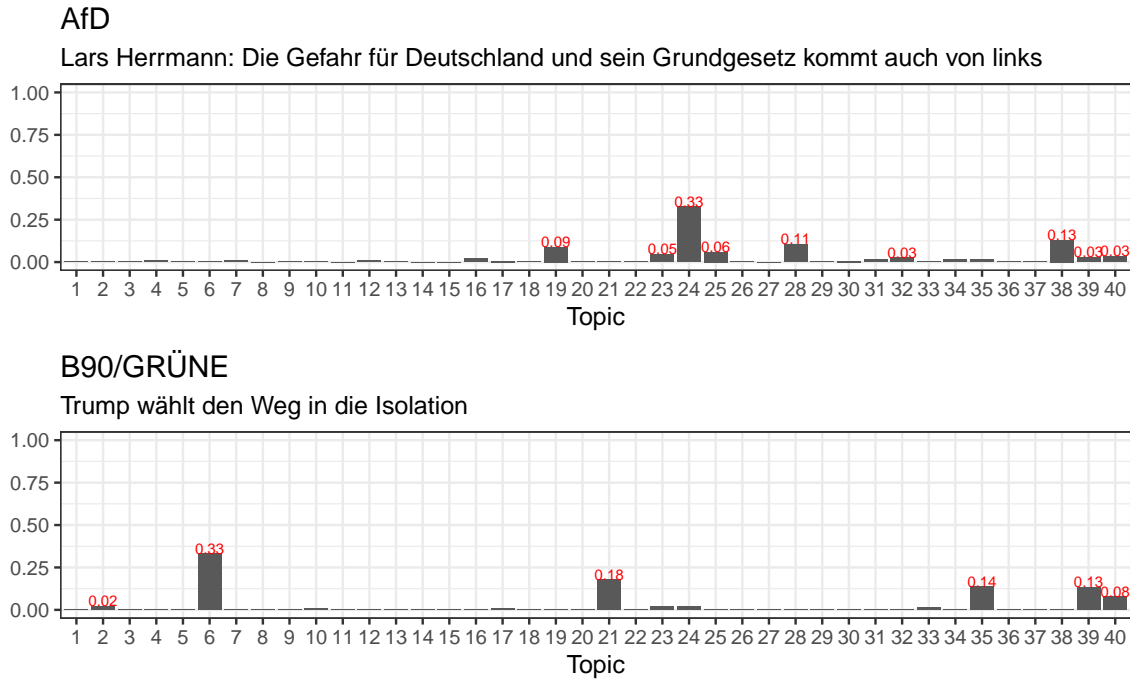[19]Trump chooses the path to isolation

Figure 8: Topic probability of sample press releases

Since each document's source and publication date are known, the probability of specific topics can be analyzed, aggregated by this metadata. The left chart of Figure 9 shows the 15 topics with the highest probability for press releases published by the AfD. The right side of the figure aggregates the probability by source and time (in weeks) for two sample topics, displaying how they change over time in the AfD press releases compared to two sample newspapers. It becomes clear that topic 9[20] is systematically more likely in the AfD's press releases compared to the two newspapers Bild.de and Handelsblatt. There is a noticeable increase in probability during the election campaign and ends in a peak on election day itself. For Handelsblatt and Bild.de, too, a slight increase of probability around election day is discernible.

The top words of topic 38 suggest that it addresses refugees - a topic for which the AfD has an absolute position. The probability of this topic increases in the AfD's press releases until about a month before the election and then levels off somewhat. A similar trend is discernible in the news articles from Bild.de. The curve from Handelsblatt is relatively flat and shows no apparent difference between before and after the election.

Figure 10 allows a similar analysis for the aggregated topic distribution in press releases of the FDP. The chart on the left illustrates that topic 39[21] has the highest probability in the FDP press releases. The two sample topics show clear temporal peaks: For topic 33[22], an increase can be seen in the FDP's press releases immediately after the election, when talks for a possible Jamaica coalition were taking place. However, for the two newspapers, the share

---

[20]translation: afd, gauland, weidel, alexander, alice, party, germany
[21]translation: germany, bund, states, federal government, education, states, municipalities
[22]translation: fdp, jamaika, coalition, lindner, union, re-elections, grünen

Figure 9: Comparison of topic probability - sample topics AfD

of this topic peaked around November 19, 2017, after the FDP announced its withdrawal from the negotiations. Topic 10[23] has a clear peak for both the newspapers and the FDP press releases around august 2017. There was a debate about whether and where driving bans for diesel cars would be introduced. After the states of Baden-Württemberg and North Rhine-Westphalia initially filed a lawsuit against this, the court proceedings that would decide whether driving bans are permissible began in mid-February 2018. The temporal curve of the FDP shows a further increase in topic probability at this time, which can also be detected at Handelsblatt. At Bild.de, however, the topic is only taken up once briefly in August 2017, as only a very low topic probability can be seen after that.



Figure 10: Comparison of topic probability - sample topics FDP

---

[23]translation: diesel, enterprises, germany, cars, german, industry, driving bans

## 3.3 Similarity measure

The topic distributions calculated by the STM are a vectorized representation of each document as represented by each row in the matrix in Table 3. Therefore, it is possible to calculate the similarity between two documents by estimating the cosine similarity between these vectors.[24]

Table 3: Document-topic distribution matrix

| doc_index | 1 | 2 | 3 | 4 | 5 | 6 | .. | 40 |
|---|---|---|---|---|---|---|---|---|
| 1 | 0.0006 | 0.0461 | 0.0008 | 0.0015 | 0.2259 | 0.0118 | ... | 0.0195 |
| 2 | 0.0045 | 0.0285 | 0.0001 | 0.0026 | 0.0005 | 0.1970 | ... | 0.0044 |
| 3 | 0.0044 | 0.0040 | 0.0017 | 0.0006 | 0.0046 | 0.0191 | ... | 0.0894 |
| 4 | 0.0005 | 0.0448 | 0.0006 | 0.0013 | 0.2575 | 0.0093 | ... | 0.0184 |
| 5 | 0.0003 | 0.0534 | 0.0004 | 0.0012 | 0.2859 | 0.0099 | ... | 0.0142 |

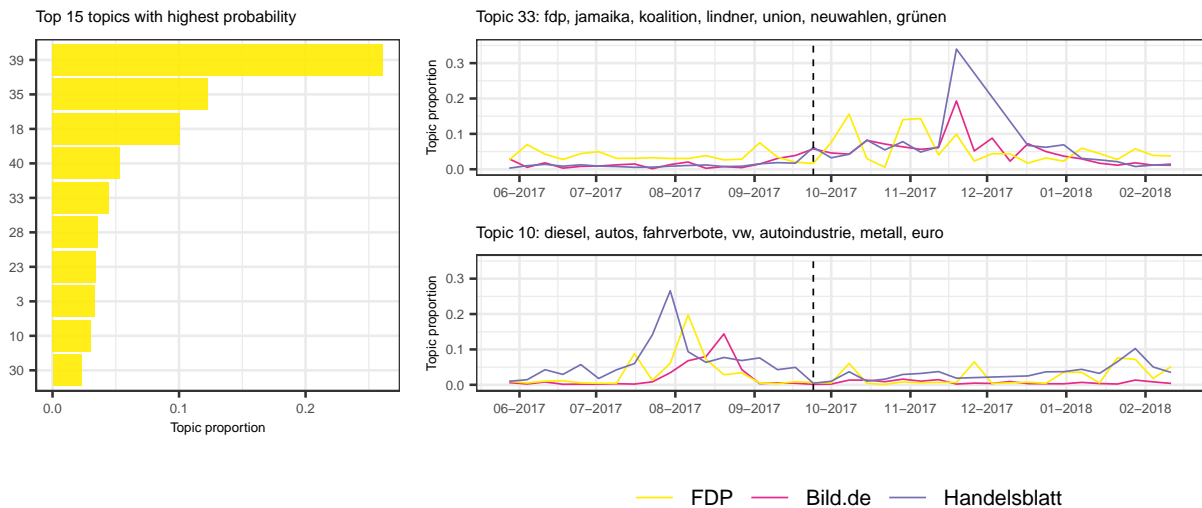The cosine similarity is the cosine of the angle between two vectors projected in a multi-dimensional space and is defined between zero and one; values towards 1 indicate similarity. For example, the cosine similarity (CS) between document 1 and 2 for $K = 40$ topics can be calculated as follows.

$$\text{CS} = \cos(\vec{doc_1}, \vec{doc_2}) = \frac{\vec{doc_1} * \vec{doc_2}}{||\vec{doc_1}|| ||\vec{doc_2}||} = \frac{\sum_{i=1}^{K} \vec{doc}_{1,i} \vec{doc}_{2,i}}{\sqrt{\sum_{i=1}^{K} \vec{doc}_{1,i}^2}, \sqrt{\sum_{i=1}^{K} \vec{doc}_{2,i}^2}}$$

For each newspaper, the cosine similarity between all topic-document distribution pairs between the newspapers articles and the press releases is calculated if that press release was published within seven days before the publication date of the news article. Thus, the topic distribution of news article 1 is compared to press releases 1, 2, 3, and so on for press releases published within seven days before the news article. Table 4 illustrates a sample subset of the data for DIE WELT.

Table 4: Dataset structure step 1 - DIE WELT

| title1 | title2 | cosine_sim | source1 | source2 | date1 | date2 |
|---|---|---|---|---|---|---|
| Wolfgang Schäuble... | Keine gemeinsame... | 0.87 | DIE WELT | FDP | 2018-01-18 | 2018-01-17 |
| Hartz IV: In Deut... | René Springer: Sc... | 0.05 | DIE WELT | AfD | 2017-11-27 | 2017-11-21 |
| SPD-Parteitag: So... | Holocaust nicht l... | 0.05 | DIE WELT | DIE LINKE | 2018-01-21 | 2018-01-18 |
| Merkel reagiert z... | Kampf gegen Mens... | 0.05 | DIE WELT | CDU | 2017-12-07 | 2017-12-01 |
| Behördenpanne: Zu... | Martin Hess: Fußf... | 0.05 | DIE WELT | AfD | 2017-11-18 | 2017-11-17 |

Next, the mean cosine similarity for each news article publication date (date1) and party (source2) is estimated to obtain the final data frame (see Table 5).

---

[24]For applications of cosine similarity to compare of topic model outcomes see e.g. Rehs (2020) and Ramage, Dumais, and Liebling (2010)

Table 5: Final dataset structure - DIE WELT

| date1 | source1 | source2 | cos_sim |
|---|---|---|---|
| 2017-10-12 | DIE WELT | B90/GRÜNE | 0.15 |
| 2017-06-04 | DIE WELT | DIE LINKE | 0.15 |
| 2018-01-16 | DIE WELT | FDP | 0.20 |
| 2017-11-24 | DIE WELT | CDU | 0.09 |
| 2017-11-25 | DIE WELT | FDP | 0.12 |

## 3.4 Model estimations

Finally, cosine similarity can be used as the independent variable in different model specifications to answer the research questions outlined previously. In 3.4.1 OLS dummy regression, an OLS model with party dummies is computed for the pre-election period to analyze whether online news equally addresses the topics covered in the press releases of different parties during the election campaign. The alternative hypothesis is that there is a significant difference between the topic similarity for different parties, indicating biased reporting of the individual newspapers. However, if different newspapers slant towards different parties, the overall landscape of political news would still be unbiased.

In 3.4.2. Regression discontinuity in time (RDiT) regression, a regression discontinuity is specified to test whether the election day affected the topic similarity overall (aggregated over all parties), respectively for different news/party combinations in detail (see 3.4.3 RDiT dummy regression).

### 3.4.1 OLS dummy regression

To measure whether there is a significant difference in the topic similarity for each party for a news publisher, a simple OLS regression is estimated, where the similarity score on day $t$ between the news articles and press releases is the dependent variable ($\mathrm{CS}_t$) and dummy-variables for different parties are the independent regressors.

$$\mathrm{CS}_t = \beta_0 + \beta_j D_{t,j} + \epsilon_t,$$

with $t = \mathrm{date}$[25] and $j = \{1, 2, \ldots, k-1\}$ for $k$ political parties[26].

The dummy-variable coefficients $\beta_j$ represent the mean difference between each of the other parties and the reference category $k$, conditional on any other predictors. The intercept is interpreted as the mean similarity score when the predictors are all 0. In the model estimated below, AfD is the reference group, i.e. the coefficients can be interpreted as the difference of topic similarity between any party and the AfD, whereas the intercept represent the mean topic similarity of AfD.

---

[25]date1 in Table 5
[26]source2 in Table 5

**OLS dummy results**   The columns in Table 6 report the results for each news publisher.[27] The F-statistic of each model indicates whether we can reject the null hypothesis that all regressor coefficients are equal to zero $H_0 : \beta_j = 0$. We can reject that hypothesis for all models except for Handelsblatt, meaning that the topic similarity can not be explained with the party dummies. Since all other newspaper models show a significant F-statistic, it can be concluded that topic similarity varies for different parties. Similarly, the p-values of the individual coefficients are significant at the 5% level, allowing to reject the null hypothesis. As stated above, the coefficients give the difference in intercepts compared to the base category AfD. Therefore, the coefficient for B90/GRÜNE in the first column - representing the model for Bild.de - indicates that the topic similarity between B90/GRÜNE and Bild.de is 0.055 points[28] lower than the topic similarity between AfD and Bild.de, holding everything else equal. Table 6 reveals that all coefficients are negative, meaning that the topic similarity is significantly lower between news articles and press releases when compared to AfD for all party/newspaper pairs (except for Handelsblatt).

Table 6: Results from the OLS dummy regression

| | *Dependent variable:* | | | | | | |
|---|---|---|---|---|---|---|---|
| | Cosine similarity of topic distribution | | | | | | |
| | Bild.de | DIE WELT | FOCUS Online | Handelsblatt | SPIEGEL ONLINE | stern.de | ZEIT ONLINE |
| | (1) | (2) | (3) | (4) | (5) | (6) | (7) |
| B90/GRÜNE | −0.057*** | −0.041*** | −0.052*** | 0.016* | −0.027*** | −0.034*** | −0.035*** |
| | (0.012) | (0.007) | (0.007) | (0.009) | (0.009) | (0.007) | (0.010) |
| CDU | −0.059*** | −0.044*** | −0.043*** | 0.019** | −0.033*** | −0.023*** | −0.041*** |
| | (0.012) | (0.007) | (0.007) | (0.009) | (0.009) | (0.007) | (0.010) |
| DIE LINKE | −0.040*** | −0.038*** | −0.039*** | 0.024*** | −0.026*** | −0.018** | −0.030*** |
| | (0.012) | (0.007) | (0.007) | (0.009) | (0.009) | (0.007) | (0.010) |
| FDP | −0.064*** | −0.051*** | −0.058*** | 0.009 | −0.050*** | −0.033*** | −0.040*** |
| | (0.012) | (0.007) | (0.007) | (0.009) | (0.009) | (0.007) | (0.010) |
| SPD | −0.068*** | −0.050*** | −0.056*** | 0.006 | −0.050*** | −0.032*** | −0.041*** |
| | (0.012) | (0.007) | (0.007) | (0.009) | (0.009) | (0.007) | (0.010) |
| Constant | 0.201*** | 0.185*** | 0.176*** | 0.164*** | 0.168*** | 0.148*** | 0.200*** |
| | (0.008) | (0.005) | (0.005) | (0.006) | (0.006) | (0.005) | (0.007) |
| Observations | 594 | 683 | 695 | 641 | 695 | 689 | 671 |
| R² | 0.077 | 0.097 | 0.118 | 0.017 | 0.064 | 0.048 | 0.039 |
| Adjusted R² | 0.069 | 0.090 | 0.112 | 0.010 | 0.057 | 0.041 | 0.031 |
| Residual Std. Error | 0.081 (df = 588) | 0.053 (df = 677) | 0.054 (df = 689) | 0.063 (df = 635) | 0.065 (df = 689) | 0.053 (df = 683) | 0.073 (df = 665) |
| F Statistic | 9.765*** (df = 5; 588) | 14.557*** (df = 5; 677) | 18.481*** (df = 5; 689) | 2.232** (df = 5; 635) | 9.409*** (df = 5; 689) | 6.847*** (df = 5; 683) | 5.352*** (df = 5; 665) |

*Note:* *p<0.1; **p<0.05; ***p<0.01

Figure 11 plots the sum of the intercept and the coefficients $(\beta_0 + \beta_j)$ for all models with a significant F-statistic to illustrate the overall magnitude of the effect. This shows that the topics in the news articles from Bild.de and ZEIT ONLINE are most similar to press releases of AfD (remember that the intercept represents the conditional mean topic similarity of the base category AfD), whereas the similarity is lowest in the case of stern.de and Handelsblatt. Furthermore, the figure visualizes that the topic similarity for all party/newspaper pairs is significantly smaller, when compared to the AfD. This difference is biggest for Bild.de, meaning that this newspaper has the strongest bias towards topics adressed in AfD press releases in the period under consideration.

The results also allow to compare between any other two parties, by taking the difference in their dummy-regressor coefficients. However, the differences are relatively small compared

---

[27]All regression output tables are created using Hlavac (2018)

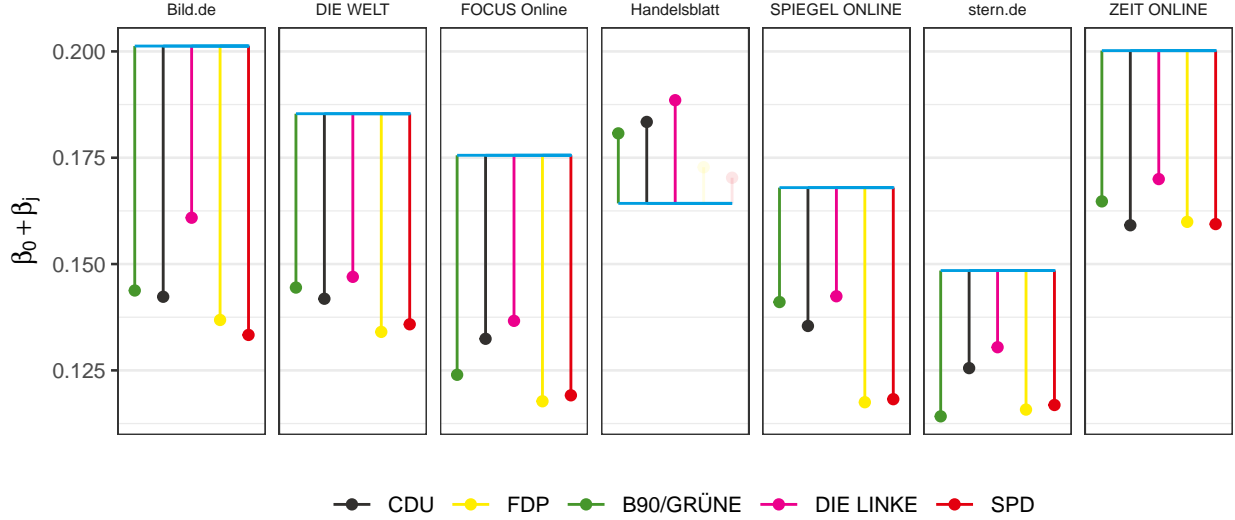[28]Remember that topic similarity is defined between 0 and 1.

Figure 11: Coefficients of OLS dummy regression

to the differences between the intercept (AfD) and the remaining coefficients. In summary, for all newspapers - except for Handelsblatt - the null hypothesis that the dummy regressors do not affect the topic similarity can be rejected. Additionally, the results show that topic similarity is significantly smaller for all parties compared to AfD. Thus, all newspapers under examination slant towards AfD topics during the election campaign resulting in a biased landscape for political news. However, it is worth noting that the model only considers the similarity of *which* topics are addressed and not *how* the topics are discussed.[29]

### 3.4.2 Regression discontinuity in time (RDiT) regression

As mentioned above, we assume that the election day affects the reporting since media outlets can observe the beliefs of potential readers. The underlying dynamic of this assumption coincides with the basic idea of regression discontinuity design (RDD). Therefore, an RDD is applied to identify the short-term effect of the election on the topic similarity between newspaper articles and press releases. The RDD was designed by Thistlethwaite and Campbell (1960) and formalized by Hahn, Todd, and Van der Klaauw (2001) to measure the effect of a treatment in a non-experimental setting, where the treatment is defined as a discontinuous function of a continuous, observed variable (the 'running' or 'forcing' variable). Like Thistlethwaite and Campbell (1960), who estimated the effect of receiving the National Merit Scholarship on future academic outcomes, early studies that rely on RD designs estimate the effects of certain thresholds of a running variable on educational outcomes (i.e., financial aid (van der Klaauw 2002) or class size (Angrist and Lavy 1999)). Following these early studies in education, the RDD has received attention in a broader range of the economic literature, including labor economics, political economy, health economics, and environmental economics. Compared to alternative quasi-experimental estimators like

---

[29]Nevertheless, it is assumed in communication literature that smaller, non-established parties benefit from placing their topics in the media to get them into the voters' heads. Here, the tendency of the reporting is irrelevant, but rather the quantity is decisive (see, e.g. Mazzoleni, Stewart, and Horsfield (2003)).

difference-in-difference and matching techniques, RDD is the estimator with the most significant internal validity (Lee and Lemieux 2010).

While RDD was applied initially in cross-sectional studies, an increasing number of studies, especially in environmental and energy economics, have adapted the framework to time series applications. In these studies, time is the running variable, and treatment begins at a particular threshold in time. A significant conceptual difference between regression discontinuity (RD) and regression discontinuity in time (RDiT) lies in the possible interpretation of the results. Since in RDiT, the running variable of time is not random eliminates the interpretation of local randomization. As noted by Jacob et al. (2012), although some researchers have focused on this interpretation of local randomization, in which the treatment status within a small neighborhood around the threshold can essentially be compared to a roll of the dice (Lee and Lemieux 2010), others have emphasized that RD is characterized by discontinuity at a threshold (Hahn, Todd, and Van der Klaauw 2001). Thus, to the extent that the RD framework is simply another quasi-experimental framework (one that uses discontinuity), RDiT is conceptually similar to RD.

In this paper, the date is the running variable, the election day is the treatment, and news publishers are the units that receive the treatment. A sharp regression design is used since the running variable (date) ultimately determines the treatment (election day). Thus, a news publisher's probability of receiving a treatment jumps from 0 to 1 at the cutoff. Specifically, the following equation is estimated:

$$\text{CS}_t = \beta_0 + \beta_1 T_t + f(W_t) + \epsilon_t$$

where

$$T_t = \begin{cases} 1, & \text{if date } \geq \text{election date} \\ 0, & \text{if date } < \text{election date} \end{cases}$$

The running variable $W_t$ is the time difference between date $i$ and the election date (in days), such that $\beta_1$ is the average treatment effect for observations with $W_t = 0$ (the election date). In other words, $\beta_1$ gives the average change of the similarity between news publisher content and press releases after the election day. Identification in the RD model comes from assuming that the underlying, potentially endogenous relationship between $\epsilon_t$ and the date is eliminated by the flexible function $f(.)$. In particular, the relationship between $\epsilon_t$ and the date must not change discontinuously on or near the election date.

Following Imbens and Lemieux (2008) we estimate a local linear regression model of the form:

$$\text{CS}_t = \beta_0 + \beta_1 T_t + \beta_2 W_t + \beta_3 W_t T_t + \epsilon_t$$

In this specification, the function $f(W_t)$ is specified as $\beta_2 W_t + \beta_3 W_t T_t$, where by $W_t T_t$ is assumed that in addition to the intercept (captured by the treatment effect $T_t$), the slope also changes after the election day. The interaction term, together with $W_t$, should absorb

any smooth relationship between the date and $\epsilon_t$ in the days surrounding the election day. Thus, if the RD assumption is valid (i.e., $\epsilon_t$ does not change discontinuously at the election day), the estimate of $\beta_1$, the coefficient of interest, will be unbiased even without further controls.

We specify a uniform kernel (Lee and Lemieux 2010) and use a bandwidth of 115 days on each side of the election day threshold. The election took place on September 24, 2017, so the sample includes dates between June 1, 2017, and January 17, 2018. Since the identification strategy only attempts to estimate $\beta$ at $W_t = 0$ (the election day), no additional dates beyond the 115-day bandwidth enter the sample. Alternative specifications with varying bandwidths led to similar results.

**RDiT Results** Since we are interested in the treatment effect at the cutoff point (remember that $W_t = 0$ for the election day) and since

$$\frac{\Delta Y}{\Delta T} = \beta_1 + \beta_3 W,$$

$\beta_1$ can be interpreted as the change in topic similarity with respect to the election day. The results in Table 7 show a significant F-statistic for all models, except for FOCUS ONLINE. We can reject the null hypothesis for all other newspapers that the regressors do not have a combined effect on the topic similarity. However, only the DIE WELT and SPIEGEL ONLINE models show a significant coefficient for $\beta_1$: For DIE WELT (-0.019) this effect is negative, indicating a drop in topic similarity overall. In the case of SPIEGEL ONLINE (0.015), the coefficient suggests an increase of topic similarity between the news articles and press releases.

Table 7: Results from the RDiT model

| | Bild.de | DIE WELT | FOCUS Online | Handelsblatt | SPIEGEL ONLINE | stern.de | ZEIT ONLINE |
|---|---|---|---|---|---|---|---|
| | \multicolumn{7}{c}{*Dependent variable:*} | | | | | | |
| | \multicolumn{7}{c}{Cosine similarity of topic distribution} | | | | | | |
| | (1) | (2) | (3) | (4) | (5) | (6) | (7) |
| T | −0.007 | −0.019*** | −0.002 | −0.005 | 0.015** | 0.004 | −0.008 |
| | (0.008) | (0.006) | (0.006) | (0.008) | (0.007) | (0.007) | (0.009) |
| W | −0.0002* | −0.0001** | −0.00003 | −0.0002*** | −0.0003*** | −0.0002** | −0.0003*** |
| | (0.0001) | (0.0001) | (0.0001) | (0.0001) | (0.0001) | (0.0001) | (0.0001) |
| TTRUE:W | 0.0004*** | 0.0003*** | −0.00000 | 0.001*** | 0.0003** | 0.0003*** | 0.0004*** |
| | (0.0001) | (0.0001) | (0.0001) | (0.0001) | (0.0001) | (0.0001) | (0.0001) |
| Constant | 0.144*** | 0.141*** | 0.132*** | 0.163*** | 0.119*** | 0.116*** | 0.150*** |
| | (0.006) | (0.004) | (0.004) | (0.006) | (0.005) | (0.004) | (0.006) |
| Observations | 1,161 | 1,212 | 1,268 | 1,046 | 1,264 | 1,218 | 1,244 |
| $R^2$ | 0.008 | 0.029 | 0.003 | 0.028 | 0.017 | 0.010 | 0.031 |
| Adjusted $R^2$ | 0.005 | 0.027 | 0.001 | 0.025 | 0.015 | 0.008 | 0.028 |
| Residual Std. Error | 0.073 (df = 1157) | 0.052 (df = 1208) | 0.053 (df = 1264) | 0.069 (df = 1042) | 0.062 (df = 1260) | 0.057 (df = 1214) | 0.079 (df = 1240) |
| F Statistic | 3.089** (df = 3; 1157) | 12.152*** (df = 3; 1208) | 1.422 (df = 3; 1264) | 10.090*** (df = 3; 1042) | 7.280*** (df = 3; 1260) | 4.072*** (df = 3; 1214) | 13.027*** (df = 3; 1240) |

*Note:* *p<0.1; **p<0.05; ***p<0.01

### 3.4.3 RDiT dummy regression

Since the model estimated in the previous section gives the effect of the election day on the overall topic similarity without differentiating per party, we now include dummy variables
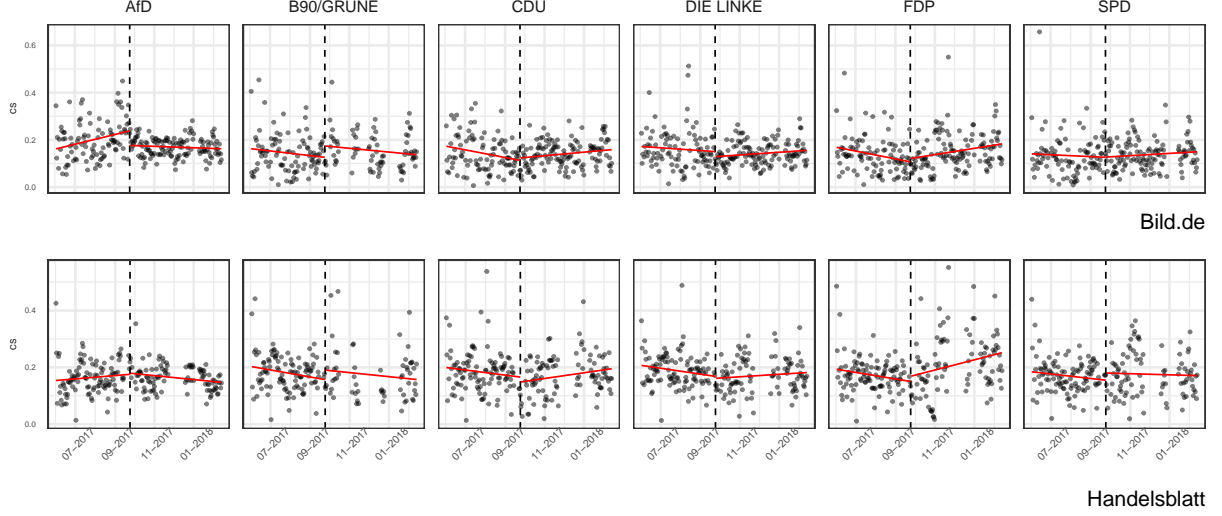
Figure 12: Mean cosine similarity between newspaper/press articles pairs - with cutoff value

for each party $k$ ($D_{t,k-1}$). In doing so, we can test whether the election day's effect on topic similarity differs for different parties and test our hypothesis that newspapers adjust their slant after observing the beliefs of potential readers.

Figure 12 visually captures that hypothesis for two sample news publishers, Bild.de and Handelsblatt. For the former, the illustration suggests a negative treatment effect for AfD and a positive effect for B90/GRÜNE. Similarly, in Handelsblatt's case, a negative effect for AfD and an adverse effect for CDU are observable. Since the figures (see Figure 14 for all news publishers) suggest that the slope changes after the election day for nearly all newspaper/party pairs, the interaction term $T_t D_{t,k-1}$ is included to capture this effect. Thus, $\gamma_1, ..., \gamma_{k-1}$ give the average treatment effect for each newspaper/party pair in the following equation.

$$\text{CS}_t = \beta_0 + \beta_1 T_t + \beta_2 W_t + \beta_3 W_t T_t + \beta_j D_{t,j} + \gamma_j T_t D_{t,j} + \epsilon_t$$

**RDiT dummy results**    Table 8 outputs the results for all newspaper models. The coefficients for the treatment variables (e.g., "TTRUE:FDP") show the effect of the election day on the topic similarity depending on the party for a given $W$. This effect can be illustrated using DIE WELT and FDP as an example and comparing the model equation for $D_{FDP} = 1$ and $D_{FDP} = 0$ for $W = 0$.

$$D_{FDP} = 1 : \hat{Y} = 0.178 + (-0.026)T + (0.025)T$$
$$D_{FDP} = 0 : \hat{Y} = 0.178 + (-0.026)T$$

In other words, when $D_{FDP}$ switches from 0 to 1, the treatment effect decreases by 0.025 compared to the base dummy group AfD, for which the treatment effect is $-0.025$.

Table 8: Results from the regression discontinuity model

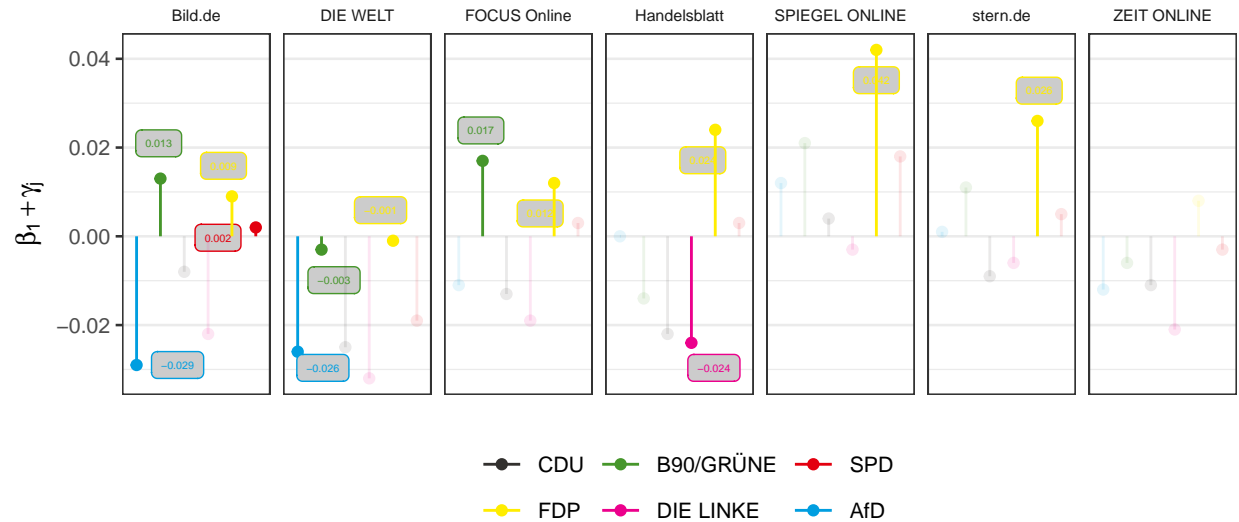| | | | | *Dependent variable:* | | | |
|---|---|---|---|---|---|---|---|
| | | | | Cosine similarity of topic distribution | | | |
| | Bild.de | DIE WELT | FOCUS Online | Handelsblatt | SPIEGEL ONLINE | stern.de | ZEIT ONLINE |
| | (1) | (2) | (3) | (4) | (5) | (6) | (7) |
| T | −0.029** | −0.026*** | −0.011 | −0.0001 | 0.012 | 0.001 | −0.012 |
| | (0.012) | (0.009) | (0.008) | (0.013) | (0.010) | (0.010) | (0.013) |
| W | −0.0002* | −0.0001** | −0.00003 | −0.0002*** | −0.0003*** | −0.0002*** | −0.0003*** |
| | (0.0001) | (0.0001) | (0.0001) | (0.0001) | (0.0001) | (0.0001) | (0.0001) |
| B90/GRÜNE | −0.057*** | −0.040*** | −0.051*** | 0.018* | −0.026*** | −0.033*** | −0.034*** |
| | (0.010) | (0.007) | (0.007) | (0.009) | (0.008) | (0.007) | (0.010) |
| CDU | −0.058*** | −0.043*** | −0.042*** | 0.020** | −0.031*** | −0.022*** | −0.039*** |
| | (0.010) | (0.007) | (0.007) | (0.009) | (0.008) | (0.007) | (0.010) |
| DIE LINKE | −0.040*** | −0.038*** | −0.038*** | 0.025*** | −0.025*** | −0.017** | −0.029*** |
| | (0.010) | (0.007) | (0.007) | (0.009) | (0.008) | (0.007) | (0.010) |
| FDP | −0.064*** | −0.051*** | −0.057*** | 0.009 | −0.050*** | −0.032*** | −0.039*** |
| | (0.010) | (0.007) | (0.007) | (0.009) | (0.008) | (0.007) | (0.010) |
| SPD | −0.068*** | −0.049*** | −0.056*** | 0.007 | −0.049*** | −0.031*** | −0.040*** |
| | (0.010) | (0.007) | (0.007) | (0.009) | (0.008) | (0.007) | (0.010) |
| TTRUE:W | 0.0003*** | 0.0003*** | −0.00001 | 0.001*** | 0.0003** | 0.0003*** | 0.0004*** |
| | (0.0001) | (0.0001) | (0.0001) | (0.0001) | (0.0001) | (0.0001) | (0.0001) |
| TTRUE:B90/GRÜNE | 0.042*** | 0.023** | 0.028*** | −0.014 | 0.009 | 0.010 | 0.006 |
| | (0.016) | (0.011) | (0.011) | (0.017) | (0.013) | (0.012) | (0.017) |
| TTRUE:CDU | 0.021 | 0.001 | −0.002 | −0.022 | −0.008 | −0.010 | 0.001 |
| | (0.014) | (0.010) | (0.009) | (0.014) | (0.011) | (0.011) | (0.015) |
| TTRUE:DIE LINKE | 0.007 | −0.006 | −0.008 | −0.024* | −0.015 | −0.007 | −0.009 |
| | (0.014) | (0.010) | (0.009) | (0.014) | (0.011) | (0.011) | (0.015) |
| TTRUE:FDP | 0.038*** | 0.025*** | 0.023** | 0.024* | 0.030*** | 0.025** | 0.020 |
| | (0.014) | (0.010) | (0.010) | (0.015) | (0.011) | (0.011) | (0.015) |
| TTRUE:SPD | 0.031** | 0.007 | 0.014 | 0.003 | 0.006 | 0.004 | 0.009 |
| | (0.014) | (0.010) | (0.010) | (0.015) | (0.012) | (0.011) | (0.015) |
| Constant | 0.192*** | 0.178*** | 0.173*** | 0.149*** | 0.149*** | 0.138*** | 0.181*** |
| | (0.009) | (0.006) | (0.006) | (0.008) | (0.007) | (0.006) | (0.009) |
| Observations | 1,161 | 1,212 | 1,268 | 1,046 | 1,264 | 1,218 | 1,244 |
| $R^2$ | 0.076 | 0.132 | 0.122 | 0.050 | 0.086 | 0.051 | 0.061 |
| Adjusted $R^2$ | 0.065 | 0.122 | 0.113 | 0.038 | 0.077 | 0.041 | 0.051 |
| Residual Std. Error | 0.071 (df = 1147) | 0.049 (df = 1198) | 0.050 (df = 1254) | 0.068 (df = 1032) | 0.060 (df = 1250) | 0.056 (df = 1204) | 0.078 (df = 1230) |
| F Statistic | 7.233*** (df = 13; 1147) | 13.990*** (df = 13; 1198) | 13.370*** (df = 13; 1254) | 4.216*** (df = 13; 1032) | 9.073*** (df = 13; 1250) | 4.988*** (df = 13; 1204) | 6.157*** (df = 13; 1230) |

*Note:* *p<0.1; **p<0.05; ***p<0.01



Figure 13: Coefficients of RDiT dummy regression

Figure 13 plots the sum of the treatment coefficients and the interaction terms $(\beta_1 + \gamma_j)$ for $W = 0$. Remember that the coefficient of the treatment effect $(\beta_1)$ shows the treatment effect for AfD since it is the base dummy group. While in the previous model without party dummies, the treatment effect for Bild.de was not significant, the present model gives a more differentiated analysis, showing a significant negative effect for the topic similarity for AfD ($-0.029$), as well as a significant positive effect for B90/GRÜNE ($0.013$), SPD ($0.002$) and FDP ($0.009$) holding the respective other variables equal. Besides Bild.de, in the case of DIE WELT, a significant negative effect of the election day on the topic similarity with AfD press releases is discernible ($-0.026$). The only other negative treatment effect - except a small negative effect for DIE WELT/B90/GRÜNE - exists for DIE LINKE in the case of Handelsblatt ($-0.024$). Positive effects mainly exist for FDP for nearly all newspapers (except DIE WELT and ZEIT ONLINE). In these cases, the election day positively affected the topic similarity between the press releases and the news articles. The same is true for B90/GRÜNE in the case of Bild.de and FOCUS Online. No significant effect of the election day can be detected in the case of ZEIT ONLINE on either of the model specifications.

# 4 Discussion and conclusion

In the run-up to the 2017 federal election, German media was accused of indirectly influencing the election through its political coverage. On the one hand, it was accused of providing a stage for the AfD through its choice of topics, which led to a rise in the party's popularity. But, on the other hand, the AfD accused the same media of devaluing the party through negative reporting. This paper investigates whether political reporting of German online newspapers was similar for the major political parties during the election campaign for the Bundestag 2017.

The STM (Structural Topic Modeling) as applied in this paper helps detect the latent topics discussed in news articles and press releases. The result from this unsupervised machine learning approach is a vectorized topic distribution for each document (either a news article or a press release), which allows to calculate the cosine similarity between documents. This approach helps analyze text content programmatically and put it into a format usable for a regression model. Using the topic similarity - measured as the cosine similarity between topic distributions - as the dependent variable, the applied regression models with different specifications allow estimating:

a) whether there is an observable difference in topic similarity between different newspaper/party pairs and

b) whether the election results had a significant effect.

Results show that the news articles of all newspapers (except for Handelsblatt) slant towards AfD topics during the election campaign resulting in a biased landscape for political news. Although no statement can be made about the tonality with which AfD-related issues are discussed, it can be assumed that the mere disproportionate mention of these topics in the media has brought the party more into the focus of voters.

Following Mullainathan and Shleifer (2005) - and assuming a competitive market - the results suggest that reader beliefs are homogeneous, causing newspapers uniformly to slant towards these beliefs. Furthermore, this explanatory approach, similar to the notion of framing, suggests that readers' beliefs are more consistent with topics discussed in AfD press releases than any other party, assuming that election results can be used as a proxy for reader beliefs. However, election polls during the research period show that the popularity of the AfD increases but stays below SPD and CDU. An alternative explanation for the results is the entertainment factor as a driver for bias. Newspapers competing for readers' attention are incentivized to cover topics with a high entertainment factor, i.e., topics containing conflict, drama, and negativity. Likewise, topics from populist parties like AfD often contain negative, emotionalized, or dramatized messages, thus utilizing similar mechanisms as the attention economy. Although the analysis conducted in this paper does not reveal the tonality of news articles, this mechanism can lead to the increased popularity of the AfD. In general, it is assumed that smaller, non-established parties benefit from placing their topics in the media to get them into the voters' heads. Here, the tendency of the reporting is irrelevant, but rather the quantity is decisive (Druckman and Parkin (2005), Eberl (2018)).

Results from the RDiT model allow analyzing the effect of the election day on news coverage. They indicate that the election day significantly affected the topic similarity for specific newspaper/party pairs. Particularly the decrease of the topic similarity with AfD (for Bild.de, DIE WELT, FOCUS ONLINE) and CDU (for Handelsblatt, stern.de) might indicate an adjustment of content based on the observed reader preferences. In the latter case, election results for CDU turned out to be worse than predicted, whereas, in the case of AfD, newspapers might realize an "over-reporting" that does not fit the true beliefs of readers. Similarly, the increase of topic similarity with FDP (Bild.de, FOCUS Online, SPIEGEL ONLINE, stern.de) could be interpreted as an upward adjustment based on the good election results of that party.

Again, it is essential to state that the interpretation of these results assumes that a) bias is driven by demand and b) that election results are a proxy of reader beliefs. Overall, the only evidence from these results is that the content of the newspapers was more similar to AfD press releases and that Election Day had a significant effect on this similarity for some newspapers. It is also necessary to remark that this research only had limited choice of newspapers. It would be interesting to reproduce the analysis for other German newspapers and extend it to other time frames. Since the empirical strategy used for this paper is a machine-based approach, it allows reproducibility and the possibility to adapt it to other datasets.

# Annex

Table 9: 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 10: 7 most probable terms per topic

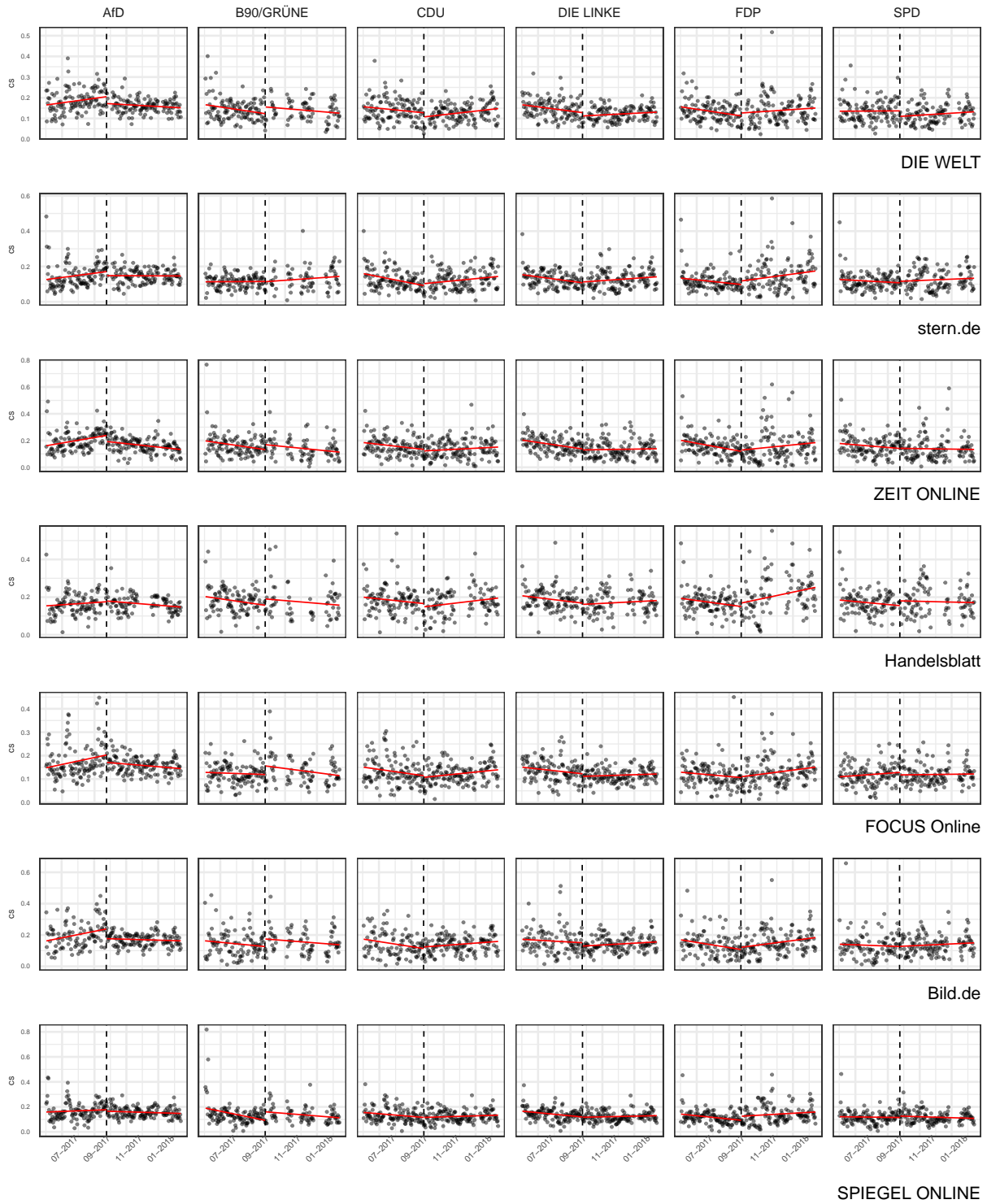| | Top Terms |
|---|---|
| 1 | a, the, s, of, b, u, to |
| 2 | merkel, angela, kanzlerin, bundeskanzlerin, cdu, merkels, wahlkampf |
| 3 | spd, union, koalitionsverhandlungen, koalitionsvertrag, groko, cdu, koalition |
| 4 | linke, linken, wagenknecht, rot, partei, linkspartei, bartsch |
| 5 | bundestagswahl, wahl, partei, afd, wähler, parteien, stimmen |
| 6 | trump, us, usa, deutschland, präsident, russland, donald |
| 7 | cdu, union, politiker, peter, zeitung, spahn, altmaier |
| 8 | csu, seehofer, parteitag, groko, söder, partei, horst |
| 9 | afd, gauland, weidel, alexander, alice, partei, deutschland |
| 10 | diesel, autos, fahrverbote, vw, autoindustrie, metall, euro |
| 11 | ge, ten, be, le, ver, te, li |
| 12 | gericht, staatsanwaltschaft, prozess, richter, urteil, verfahren, jahre |
| 13 | berlin, deutschen, jahre, tag, jahr, jahren, münchen |
| 14 | kohl, helmut, kohls, einheit, kanzler, tod, deutschen |
| 15 | august, cdu, spd, prozent, bundestagswahl, wahl, parteien |
| 16 | spd, nahles, andrea, partei, scholz, schwesig, stegner |
| 17 | csu, union, cdu, jamaika, seehofer, obergrenze, fdp |
| 18 | euro, milliarden, millionen, jahr, prozent, geld, kosten |
| 19 | polizei, stadt, polizisten, menschen, verletzt, täter, beamten |
| 20 | prozent, umfrage, spd, union, afd, cdu, fdp |
| 21 | grünen, özdemir, grüne, göring, eckardt, cem, partei |
| 22 | cdu, niedersachsen, spd, grünen, rot, landtag, fdp |
| 23 | welt, frage, lange, fragen, leute, lässt, wissen |
| 24 | g, hamburg, gipfel, polizei, hamburger, demonstranten, scholz |
| 25 | is, deutschland, verfassungsschutz, syrien, gefährder, islamisten, staat |
| 26 | steinmeier, bundespräsident, frank, walter, schmidt, spd, glyphosat |
| 27 | afd, petry, partei, fraktion, frauke, sachsen, meuthen |
| 28 | berliner, berlin, amri, maizière, innenminister, behörden, daten |
| 29 | gabriel, sigmar, außenminister, schröder, spd, amt, gerhard |
| 30 | bundestag, abgeordneten, abgeordnete, schäuble, spd, fraktion, parlament |
| 31 | frauen, kinder, deutschland, studie, eltern, muslime, antisemitismus |
| 32 | türkei, erdogan, türkischen, deutschland, bundesregierung, türkische, deutsche |
| 33 | fdp, jamaika, koalition, lindner, union, neuwahlen, grünen |
| 34 | facebook, twitter, maas, gesetz, heiko, netz, internet |
| 35 | eu, europa, deutschland, europäischen, staaten, europäische, kommission |
| 36 | bundeswehr, soldaten, leyen, nato, einsatz, ursula, verteidigungsministerin |
| 37 | spd, schulz, martin, union, kanzlerkandidat, partei, sozialdemokraten |
| 38 | flüchtlinge, deutschland, menschen, zahl, flüchtlingen, asylbewerber, jahr |
| 39 | deutschland, bund, länder, bundesregierung, bildung, ländern, kommunen |
| 40 | menschen, politik, land, deutschland, gesellschaft, politische, politischen |

Figure 14: Daily mean cosine similarity between newspaper/press articles pairs - with cutoff value

# References

Angrist, Joshua D., and Victor Lavy. 1999. "Using Maimonides' Rule to Estimate the Effect of Class Size on Scholastic Achievement." *The Quarterly Journal of Economics* 114 (2): 533–75. https://www.jstor.org/stable/2587016.

Baron, David P. 2006. "Persistent Media Bias." *Journal of Public Economics* 90 (1): 1–36. https://doi.org/10.1016/j.jpubeco.2004.10.006.

Besley, Timothy, and Andrea Prat. 2006. "Handcuffs for the Grabbing Hand? Media Capture and Government Accountability." *American Economic Review* 96 (3): 720–36. https://doi.org/10.1257/aer.96.3.720.

Bholat, David M., Stephen Hansen, Pedro M. Santos, and Cheryl Schonhardt-Bailey. 2015. "Text Mining for Central Banks." *SSRN Electronic Journal*, June. 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, 1–23. https://doi.org/10.1080/10584609.2019.1637980.

Blei, David M., Andrew Y Ng, and Michael I Jordan. 2003. "Latent Dirichlet Allocation." *Journal of Machine Learning Research* 3 (January): 993–1022.

Braun, Michael, and Jon McAuliffe. 2010. "Variational Inference for Large-Scale Models of Discrete Choice." *Journal of the American Statistical Association* 105 (489): 324–35. https://doi.org/10.1198/jasa.2009.tm08030.

Dewenter, Ralf, and Jürgen Rösch. 2014. *Einführung in die neue Ökonomie der Medienmärkte: Eine wettbewerbsökonomische Betrachtung aus Sicht der Theorie der zweiseitigen Märkte.* Springer-Verlag. https://books.google.com?id=7uXSBAAAQBAJ.

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.

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.

Gentzkow, Matthew, and Jesse M. Shapiro. 2006. "Media Bias and Reputation." *Journal of Political Economy* 114 (2): 280–316. https://doi.org/10.1086/499414.

Graber, Doris Appel. 1984. *Processing the News: How People Tame the Information Tide.* New York: Longman Press. https://books.google.com?id=pKTZAAAAMAAJ.

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 (April): 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.

Hlavac, Marek. 2018. *Stargazer: Well-Formatted Regression and Summary Statistics Tables. R Package Version 5.2.2.* https://CRAN.R-project.org/package=stargazer.

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.

Imbens, Guido W., and Thomas Lemieux. 2008. "Regression Discontinuity Designs: A Guide to Practice." *Journal of Econometrics*, The regression discontinuity design: Theory and applications, 142 (2): 615–35. https://doi.org/10.1016/j.jeconom.2007.05.001.

Jacob, Robin, Pei Zhu, Marie-Andrée Somers, and Howard Bloom. 2012. *A Practical Guide to Regression Discontinuity. MDRC.* MDRC. https://eric.ed.gov/?id=ED565862.

John Coene. (2018) 2019. *Webhoser: R Wrapper for the Webhose.io API.* news-r. https://github.com/news-r/webhoser.

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.

Klaauw, Wilbert van der. 2002. "Estimating the Effect of Financial Aid Offers on College Enrollment: A Regression-Discontinuity Approach." *International Economic Review* 43 (4): 1249–87. https://www.jstor.org/stable/826967.

Lee, David S., and Thomas Lemieux. 2010. "Regression Discontinuity Designs in Economics." *Journal of Economic Literature* 48 (2): 281–355. https://doi.org/10.1257/jel.48.2.281.

Lengauer, Günther, and David Johann. 2013. "Candidate and Party Bias in the News and Its Effects on Party Choice: Evidence from Austria." *Studies in Communication Sciences* 13 (1): 41–49. https://doi.org/10.1016/j.scoms.2013.04.011.

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.

Mazzoleni, G., J. Stewart, and Bruce Horsfield. 2003. "The Media and Neo-Populism : A Contemporary Comparative Analysis." *Undefined*. https://www.semanticscholar.org/paper/The-media-and-neo-populism-%3A-a-contemporary-Mazzoleni-Stewart/53e24419dbaddd6dbd4ca93aef60c472912fa9b0.

Mimno, David, Hanna M. Wallach, Edmund Talley, Miriam Leenders, and Andrew McCal-

lum. 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.

Mullainathan, Sendhil, and Andrei Shleifer. 2005. "The Market for News." *American Economic Review* 95 (4): 1031–53. https://doi.org/10.1257/0002828054825619.

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.

Prat, Andrea, and David Strömberg. 2013. "The Political Economy of Mass Media." In *Advances in Economics and Econometrics: Tenth World Congress: Volume 2: Applied Economics*, edited by Daron Acemoglu, Eddie Dekel, and Manuel Arellano, 2:135–87. Econometric Society Monographs. Cambridge: Cambridge University Press. https://doi.org/10.1017/CBO9781139060028.004.

Ramage, Daniel, Susan Dumais, and Daniel Liebling. 2010. *Characterizing Microblogs with Topic Models. ICWSM*.

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. Airoldi. 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.

*The Economist*. n.d. "Germany's Far-Right Party Will Make the Bundestag Much Noisier." Accessed June 12, 2022. https://www.economist.com/europe/2017/08/24/germanys-far-right-party-will-make-the-bundestag-much-noisier.

Thistlethwaite, Donald L., and Donald T. Campbell. 1960. "Regression-Discontinuity Anal-

ysis: 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.