

# Filter Bubbles, Echo Chambers, and Online News Consumption

Seth R. Flaxman                      Sharad Goel  
Carnegie Mellon University      Stanford University

Justin M. Rao  
Microsoft Research

## Abstract

Online publishing, social networks, and web search have dramatically lowered the costs to produce, distribute, and discover news articles. Some scholars argue that such technological changes increase exposure to diverse perspectives, while others worry they increase ideological segregation. We address the issue by examining web browsing histories for 50,000 U.S.-located users who regularly read online news. We find that social networks and search engines increase the mean ideological distance between individuals. However, somewhat counterintuitively, we also find these same channels increase an individual's exposure to material from his or her less preferred side of the political spectrum. Finally, we show that the vast majority of online news consumption is accounted for by individuals simply visiting the home pages of their favorite, typically mainstream, news outlets, tempering the consequences—both positive and negative—of recent technological changes. We thus uncover evidence for both sides of the debate, while also finding that the magnitude of the effects are relatively modest.

**WORD COUNT:** 5,762 words

The Internet has dramatically reduced the cost to produce, distribute, and access diverse political information and perspectives. Online publishing, for example, circumvents much of the costly equipment required to produce physical newspapers and magazines. With the rise of social media sites such as Facebook and Twitter, individuals can now readily share their favorite stories with hundreds of their contacts (Bakshy et al., 2012; Goel et al., 2012b). Moreover, search engines facilitate a diversity of voices by offering access to a range of opinions far broader than found in one’s local paper, greatly expanding the information available to citizens and their choices over news outlets.

What is the effect of such technological changes on ideological segregation? On the one hand, with more options individuals may choose only to consume content that accords with their previously held beliefs. Commentators such as Sunstein (2009) have thus predicted the rise of “echo chambers”, in which individuals are largely exposed to conforming opinions. Indeed, in controlled experiments, subjects tend to choose news articles from outlets aligned with their political opinions (Garrett, 2009; Iyengar and Hahn, 2009; Munson and Resnick, 2010). Additionally, search engines, news aggregators and social networks are increasingly personalizing content through machine learned algorithms (Agichtein et al., 2006; Das et al., 2007; Hannak et al., 2013), which can in principle create “filter bubbles” (Pariser, 2011) that amplify ideological segregation. Moreover, individuals are more likely to share information that conforms to opinions in their local social neighborhoods (Moscovici and Zavalloni, 1969; Myers and Bishop, 1970; Schkade et al., 2007; Spears et al., 1990). If realized, such information segregation is a serious concern, as it has long been thought that functioning democracies depend critically on voters who are exposed to and understand a variety of political views (Baron, 1994; Downs, 1957; Lassen, 2005).

On the other hand, Benkler (2006) and others have argued that increased choice and social networks lead to greater exposure to diverse ideas, breaking individuals free from insular consumption patterns (Goel et al., 2012a; Obendorf et al., 2007). Providing evidence for this view, Messing and Westwood (2012) show that social endorsements increase exposure to heterogeneous perspectives. Relatedly, Goel et al. (2010) show that a substantial fraction of ties in online social networks are between individuals on opposite sides of the political spectrum, opening up the possibility for diverse content discovery. Moreover, in the context of music consumption,

Hosanagar et al. (2013) find that personalized recommendation systems increase within-user diversity. Taken together, these results suggest technologies like web search and social networks reduce ideological segregation.

In short, there are compelling arguments on both sides of the debate. We investigate the issue by empirically examining news consumption patterns using the detailed web browsing records of 1.2 million anonymized U.S.-located Internet users. Our dataset records every web page viewed by these individuals over the three-month period between March and May of 2013—a total of 2.3 billion page views, organized as a time-series by user. The vast majority of these page views do not come from news sites, and even the majority of views on news sites concern topics for which ideological segregation is not particularly meaningful, such as sports and entertainment. To study the news consumption patterns that we are interested in, we must therefore identify substantively relevant articles (“hard news”); we must also quantify an outlet’s ideological leaning. For the first step, we apply machine learning algorithms to article text to identify hard news. We then further algorithmically separate out descriptive reporting from opinion pieces. For the second step, we use an audience-based approach (Gentzkow and Shapiro, 2011; Lawrence et al., 2010; Tewksbury, 2005) and estimate an outlet’s *conservative share*: the fraction of its readership that supported the Republican candidate in the most recent presidential election. Following past work, we then define (population-level) ideological segregation as the expected difference in the conservative shares of news outlets visited by two randomly selected individuals.

We find that segregation is slightly higher for descriptive news accessed via social media (0.12) than for articles read by directly visiting a news outlet’s home page (0.11). For opinion pieces, however, the effect is more substantial, moving from 0.13 for directly accessed articles to 0.17 for socially recommended pieces, to 0.20 for articles found via web search. To put these numbers in perspective, a difference of 0.20 corresponds to the ideological distance between the centrist *Yahoo News* and the left-leaning *Huffington Post* (or equivalently, *CNN* and the right-leaning *National Review*).

Our segregation measure is based on the distribution over the *mean* consumption for each individual. Consequently, the overall level of segregation we observe could be the result of two qualitatively different individual-level behaviors. A typical individual might regularly read a variety of liberal and conservative news outlets, but

still have a left- or right-leaning preference. Alternatively, individuals may choose to only read publications that are ideologically similar to one another, rarely reading opposing perspectives. We find strong evidence for the latter pattern. Specifically, users who predominately visit left-leaning news outlets only very rarely ( $< 5\%$  of the time overall) read substantive news articles from conservative sites, and vice versa for right-leaning readers, an effect that is even more pronounced for opinion articles. Interestingly, exposure to opposing perspectives is higher (more than double, though still low in absolute terms) for the channels associated with the highest segregation, search and social. Thus, counterintuitively, we find evidence that recent technological changes both increase *and* decrease various aspects of the partisan divide.

Finally, we note that directly accessed, descriptive reporting comprises 75% of traffic, primarily driven by mainstream news outlets. It accordingly appears that social networks and web search have not transformed news consumption to the degree many have hoped or feared. Why do these channels not dominate the circulation of news? One explanation is that social media platforms are used primarily for entertainment and interpersonal communication rather than political discussion. Indeed, we find that only about 1 in 300 outbound clicks from Facebook correspond to substantive news, with video and photo sharing sites far-and-away the most popular destinations. For web search, users may simply find it more convenient to visit their favorite news site rather than searching for a news topic. However, we find that for opinion stories—which account for 6% of hard-news consumption—about one-third come through social or search. So if opinion content has an outsized importance on citizen’s political views, these channels may still be substantively important. Moreover, the next generation of Internet users may increasingly rely on social media to obtain news and opinion, with corresponding implications for ideological segregation.

## 1 Data and Methods

Our primary analysis is based on web browsing records collected via the Bing Toolbar, a popular add-on application for the Internet Explorer web browser. Upon installing the toolbar, users can consent to sharing their data via an opt-in agree-

ment, and to protect privacy, all records are anonymized prior to our analysis. Each toolbar installation is assigned a unique identifier, giving the data a panel structure. We analyze the web browsing behavior of 1.2 million U.S.-located users for the three-month period between March and May of 2013, making this one of the largest studies of web content consumption to date. For each user, we have a time-stamped collection of URLs opened in the browser, along with the user’s geographic location, as inferred via the IP address. In total, our dataset consists of 2.3 billion distinct page views, with a median of 991 page views per individual.

As with nearly all observational studies of individual-level web browsing behavior, our study is restricted to individuals who voluntarily share their data, which likely creates selection issues. These users, for example, are presumably less likely to be concerned about privacy. Moreover, it is generally believed that Internet Explorer users are on average older than the Internet population at large. Instead of attempting to re-balance our sample using difficult-to-estimate and potentially incorrect weights, we acknowledge these shortcomings and note throughout where they might be a concern. When appropriate, we also replicate our analysis on different subsets of the full dataset, such as particularly heavy users.

## 1.1 Identifying News and Opinion Articles

We select an initial universe of news outlets (i.e., web domains) via the Open Directory Project (ODP, [dmoz.org](http://dmoz.org)), a collective of tens of thousands of editors who hand-label websites into a classification hierarchy. This gives 7,923 distinct domains labeled as: news, politics/news, politics/media, and regional/news. Since the vast majority of these news sites receive relatively little traffic, to simplify our analysis we restrict to the one hundred domains that attracted the largest number of unique visitors from our sample of toolbar users.<sup>1</sup> This list of popular news sites includes every major national news source, well-known blogs and many regional dailies, and collectively accounts for over 98% of consumption across all news sites. The complete list is given in the Appendix.

The bulk of the 4.1 million articles we consider do not fall into categories where political leaning has a meaningful interpretation, but rather relate to sports,

---

<sup>1</sup>This list has high overlap with the current Alexa rankings of news outlets (<http://www.alexa.com/topsites/category/Top/News>).

Table 1: Most predictive words for classifying articles as either news or non-news, and separately, for separating out descriptive news from opinion.

**Front-section news & opinion (+) vs. “non-news” (-)**

| Positive   | Negative  |
|--|---|
| contributed, democratic<br>economy, authorities,<br>leadership, read<br>republican, democrats<br>country’s, administration | film, today<br>pretty, probably<br>personal, learn<br>technology, mind<br>posted, isn’t |

**Opinion (+) vs. descriptive news (-)**

| Positive  | Negative  |
|---|---|
| stay, seem<br>important, seems<br>isn’t, fact<br>actually, reason<br>latest, simply | contributed, reporting<br>said, say<br>spokesman, experts<br>interview, expected<br>added, hers |

weather, lifestyle, entertainment, and other largely apolitical topics. We filter out these apolitical stories by training a binary classifier on the article text. The classifier identified 1.9 million stories (46%) as “front-section” news. Next, starting from this set of 1.9 million front-section news stories, we separate out descriptive news from opinion via a second classifier; 200,000 (11%) are ultimately found to be opinion stories. Details of the article classification, including performance benchmarks, are in the Appendix.

## 1.2 Measuring the Political Slant of Publishers

In the absence of human ratings, there are no existing methods to reliably assess *article* slant with both high recall and precision.<sup>2</sup> Since our sample has over 1.9 million articles classified as either front-section news or opinion, human labeling is not feasible. We thus follow the literature (Gentzkow and Shapiro, 2010, 2011; Groseclose and Milyo, 2005) and focus on the slant at the *outlet* level, ultimately

---

<sup>2</sup>High precision is possible by focusing on the use of highly polarizing phrases such as “death panel,” but the recall of this method tends to be very low, meaning most pieces of content are not rated. Even with human ratings, the wide variety of sites we investigate—ranging from relatively small blogs to national newspapers—exhibit correspondingly diverse norms of language usage, making any content-level assessment of political slant quite difficult.

Table 2: For the 20 most popular news outlets, each outlet’s estimated conservative share (i.e., the two-party fraction of its readership that voted for the Republican candidate in the last presidential election).

| Publication            | Cons. share | Publication            | Cons. share |
|------------------------|-------------|------------------------|-------------|
| BBC                    | 0.30        | L.A. Times             | 0.46        |
| New York Times         | 0.31        | Yahoo News             | 0.47        |
| Huffington Post        | 0.35        | USA Today              | 0.47        |
| Washington Post        | 0.37        | Daily Mail             | 0.47        |
| Wall Street Journal    | 0.39        | CNBC                   | 0.47        |
| U.S. News & World Rep. | 0.39        | Christian Sci. Monitor | 0.47        |
| Time Magazine          | 0.40        | ABC News               | 0.48        |
| Reuters                | 0.41        | NBC News               | 0.50        |
| CNN                    | 0.42        | Fox News               | 0.59        |
| CBS News               | 0.45        | Newsmax                | 0.61        |

assigning articles the polarity score of the outlet in which they were published. By doing so, we clearly lose some signal. For instance, we mislabel liberal op-eds on generally conservative news sites, and we mark neutral reporting of a breaking event as having the overall slant of the outlet. Nevertheless, such a compromise is common practice, and where possible, we attempt to mitigate any resulting biases.

Unfortunately, estimates from past work (Gentzkow and Shapiro, 2010; Groseclose and Milyo, 2005) cover less than half of the 100 outlets used in our main analysis. Our solution is to construct an audience-based measure of outlet slant. Specifically, we estimate the fraction of each news outlet’s readership that voted for the Republican candidate in the most recent presidential election, which we call the outlet’s *conservative share*. Thus, left-leaning, or “liberal,” outlets have conservative shares less than about 50%, and right-leaning, or “conservative,” outlets have conservative shares greater than about 50%. To estimate the political composition of a news outlet’s readership, we use the location of each webpage view as inferred from the IP address. We can then measure how the popularity of a news outlet varies across counties as a function of the counties’ political compositions, which in turn yields the estimates we desire. We detail our approach in the Appendix.

Table 2 lists estimated conservative shares for the 20 news outlets attracting the most number of unique visitors in our dataset, ranging from the *BBC* and *The New York Times* on the left to *Fox News* and *Newsmax* on the right. While our measure of conservative share is admittedly imperfect, the list does seem largely consistent

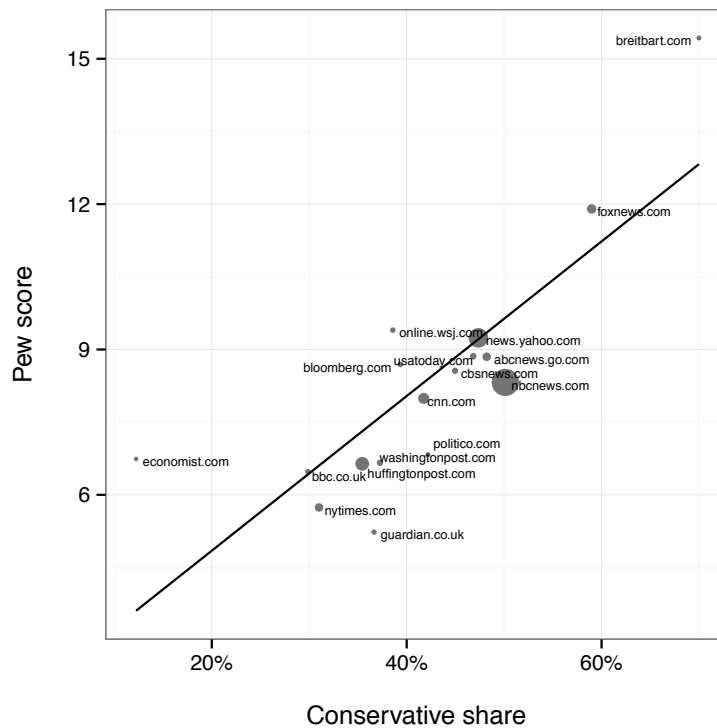


Figure 1: A comparison of our estimate of conservative share of an outlet’s audience to a Pew survey-based measure of audience ideology, where point sizes are proportional to popularity. For the 17 outlets for which both measures are available, the correlation between the two scores is 0.81.

with commonly held beliefs on the slant of particular outlets.<sup>3</sup> Furthermore, as shown in Figure 1, our ranking of news sites is highly correlated with the survey-based measure of audience ideology derived from the Pew (2014) study.<sup>4</sup> Among the 17 news sources on both lists, the correlation was 0.81.<sup>5</sup> Conservative shares for

<sup>3</sup>One exception is *The Wall Street Journal*, which we characterize as left-leaning even though it is generally thought to be politically conservative. We note, however, that the most common audience and content-based measures of slant also characterize the paper as relatively liberal (Gentzkow and Shapiro, 2011; Groseclose and Milyo, 2005). As a robustness check we repeated our analysis after omitting *The Wall Street Journal* from our dataset, and found that none of our substantive results changed.

<sup>4</sup>Pew Research Center, October 2014, “Political Polarization and Media Habits.” The report is accessible here: <http://www.journalism.org/interactives/media-polarization/>

<sup>5</sup>Comparing to the Gentzkow and Shapiro (2011) list based on 2008 audience data in which users’ party affiliations were explicitly collected, we find a correlation of 0.82 among the top 20 domains



our full list of 100 domains are given in the Appendix.

### 1.3 Inferring Consumption Channels

We define four channels through which an individual can discover a news story: direct, aggregator, social, and search. Direct discovery means a user directly and independently visits a top-level news domain such as `nytimes.com` (e.g., by typing the URL into the browser’s address bar, accessing it through a bookmark, or performing a “navigational search,” explained below), and then proceeds to read articles within that outlet. The aggregator channel refers to referrals from *Google News*—one of the last remaining popular news aggregators—which presents users with links to stories hosted on other news sites. We define the social channel to include referrals from Facebook, Twitter, and various web-based email services. Finally, the search category refers to news stories accessed as the result of web search queries on Google, Bing and Yahoo.

The time series of webpage views for an individual is not sufficient to perfectly determine discovery channel of a news article. We get around this problem with a “short” vs. “long” URL distinction in the following simple heuristic: define the “referrer” of a news article to be the most recently viewed URL that is a top-level domain such as `nytimes.com` or `facebook.com` (short URL), but not a specific story link, such as `nytimes.com/a-news-story` (long URL). We then use the referrer to classify the discovery channel. For example, if the referrer is a news domain, such as `foxnews.com`, then the channel is “direct navigation,” whereas the channel is “social” if the referrer is, for instance, `facebook.com`. The intuition behind this method is that a user is very unlikely to directly type in a specific long URL, so the visit must have a referrer, which can be inferred from the time series of URLs. Since users often use a search engine simply to navigate to a publisher’s front page (by searching for the publication’s name). This type of “navigational search” query is widely regarded as a convenient shortcut to typing in a web address (Broder, 2002) so we define it as direct navigation. The heuristic thus is based on two key assumptions: first, users do not typically type in the long, unwieldy web addresses assigned to individual articles, but rather are directed there via a previous visit to a top-level domain and a subsequent chain of clicks; and second, top-level domains are not typically shared or posted via email, social media or aggregators.

Even when referring pages can be perfectly inferred, there can still be genuine ambiguity in how to determine the consumption channel. For example, if an individual follows a Facebook link to a *New York Times* article and then proceeds to read three additional articles at that outlet, are all four articles “social” or just the first? Our solution is to take the middle ground: in this example, any subsequent article-to-article views (e.g., clicks on a “related story”) are classified as “social,” whereas an intermediate visit to the outlet’s front page results in subsequent views being classified as “direct.” Note that this is consistent with the simple procedure described above, since the site’s front page results in a short URL.

## 1.4 Limiting to Active News Consumers

As recent studies have noted, only a minority of individuals regularly read online news. For example, a 2012 survey by Pew Research showed that 39% of adults claimed to have read online news in the previous day.<sup>6</sup> Studies using actual browsing behavior tend to find that this number is quite a bit lower (Goel et al., 2012a; LaCour, 2013). Because our aim is to understand the preferences and choices of individuals who actively consume substantive news online, we limit to the subset of users who have read at least 10 substantive news articles and at least two opinion pieces in the three-month timeframe we consider. This first requirement reduces our initial sample of 1.2 million individuals to 173,450 (broadly consistent with past work); and the second requirement further reduces the sample to 50,383.

# 2 Results

## 2.1 Overall Segregation

Recall that the conservative share of a news outlet—which we also refer to as the outlet’s *polarity*—is the estimated fraction of the publication’s readership that voted for the Republican candidate in the most recent presidential election. We now define the polarity of an *individual* to be the typical polarity of the news outlet that he or she visits. We then define segregation to be the expected distance between the

---

<sup>6</sup><http://www.people-press.org/2012/09/27/in-changing-news-landscape-even-television-is-vulnerable/>

polarity scores of two randomly selected users. This definition of segregation, which is in line with past work (Dandekar et al., 2013), intuitively captures the idea that segregated populations are those in which pairs of individuals are, on average, far apart.

Due to sparsity in the data, however, our measure of segregation is not entirely straightforward to estimate. In particular, under a naive inference strategy, noisy estimates of user polarities would inflate the estimate of segregation. We thus estimate segregation via a hierarchical Bayesian model (Gelman and Hill, 2007), as described in more detail below. Finally, we note that throughout our analysis we consider the segregation associated with various subsets of consumption (e.g., views of opinions stories on social media sites). Intuitively, such measures correspond to first restricting to the relevant subset of consumption, and then computing the segregation effects; in practice, though, we simultaneously estimate the numbers in a single, random effects model.

We define the polarity score of an *article* to be the polarity score of the news outlet in which it was published.<sup>7</sup> Now, let  $X_{ij}$  be the polarity score of the  $j$ -th article read by user  $i$ . We model:

$$X_{ij} \sim N(\mu_i, \sigma_d^2) \tag{1}$$

where  $\mu_i$  is the latent polarity score for user  $i$ , and  $\sigma_d$  is a global dispersion parameter (to be estimated from the data). To mitigate data sparsity, we further assume the latent variables  $\mu_i$  are themselves drawn from a normal distribution. That is,

$$\mu_i \sim N(\mu_p, \sigma_p^2). \tag{2}$$

To complete the model specification, we assign weak priors to the hyperparameters  $\sigma_d$ ,  $\mu_p$  and  $\sigma_p$ . Ideally, we would perform a fully Bayesian analysis to obtain the posterior distribution of the parameters. However, for computational convenience, we use the approximate marginal maximum likelihood estimates obtained from the `lmer()` function in the R package `lme4` (Bates et al., 2013).

Though the distributional assumptions we make are standard in the litera-

---

<sup>7</sup>While this is standard practice, it ignores, for example, the possibility of a conservative outlet publishing liberal editorials. Ideally, the classification would be done at the article level, but there are no known methods for reliably doing so.

ture (Gelman and Hill, 2007), our modeling choices of course affect the estimates we obtain. As a robustness check, we note that a naive, model-free estimation procedure yields qualitatively similar, though ostensibly less precise, results. Moreover, in our analysis of Twitter in Section A.1—a setting where sparsity is not an issue—we estimate user polarity scores directly and find that they are indeed approximately normally distributed.

Having specified the model, we can now formally define segregation, which we do in terms of the expected squared distance between individuals’ polarity scores. Namely, we define segregation to be  $\sqrt{\mathbb{E}(\mu_i - \mu_j)^2}$ . After simple algebraic manipulation, our measure of segregation further reduces to  $\sqrt{2}\sigma_p$ . Higher values of this measure correspond to higher levels of segregation, with individuals more spread out across the ideological spectrum.

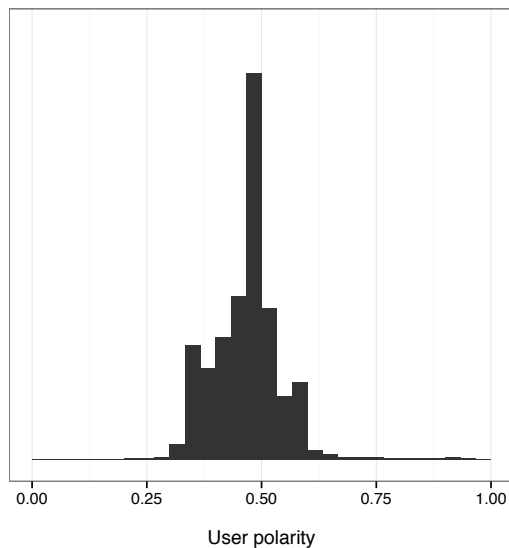


Figure 2: The distribution of individual-level polarity, where each individual’s polarity score is the (model-estimated) average conservative share of the news outlets he or she visits.

Figure 2 shows the distribution of polarity scores (i.e., the distribution of  $\mu_i$ ) for users in our sample. We find that most individuals are relatively centrist, with two-thirds of people having polarity scores between 0.41 and 0.54. Overall segregation is estimated to be 0.11, which means that for two randomly selected users, the ideological distance between the publications they typically read is on par with that

between the centrist *NBC News* and the left-leaning *Daily Kos* (or equivalently, *ABC News* and *Fox News*). Thus, though we certainly find a degree of ideological segregation, it would seem to be relatively moderate.

## 2.2 Segregation by Channel and Article Subjectivity

When measuring segregation across various distribution channels and levels of article subjectivity, the data sparsity issues we encountered above are exacerbated. For example, even among active news consumers, relatively few individuals regularly read news articles from both aggregators and social media sites. And when we further segment articles into opinion and descriptive news, it compounds the problem. However, the polarity of consumption for a user across channels should be correlated; for example, the opinion pieces one reads from Facebook are likely ideologically related to the articles one reads from Google News. There is thus opportunity to improve our estimates by “sharing strength” across channels and subjectivity levels, and accordingly to jointly estimate the segregation parameters of interest. Joint estimation with weak priors also mitigates channel selection issues.

The four consumption channels (aggregator, direct, web search and social media) and two subjectivity classes (descriptive reporting and opinion) give eight subjectivity-by-channel dimensions. Let  $X_{ijk}$  denote the polarity of the  $j$ -th article that user  $i$  reads in the  $k$ -th subjectivity-by-channel category, where we recall that the polarity of an article is defined to be the conservative share of the site on which it was published. Generalizing our hierarchical Bayesian framework, we model

$$X_{ijk} \sim N(\mu_i^k, \sigma_d^2) \quad (3)$$

where  $\mu_i^k$  is the  $k$ -th component in the latent 8-dimensional polarity vector  $\vec{\mu}_i$  for user  $i$ , and  $\sigma_d$  is a global dispersion parameter. As before, we deal with sparsity by further assuming a distribution on the latent variables  $\vec{\mu}_i$  themselves. In this case, we model each individual’s polarity vector as being drawn from a multivariate normal:

$$\vec{\mu}_i \sim N(\vec{\mu}_p, \Sigma_p) \quad (4)$$

where  $\vec{\mu}_p$  and  $\Sigma_p$  are global hyperparameters. The full Bayesian model is analyzed by assigning weak priors to the hyperparameters and computing posterior distribu-

| Consumption channel | Front-section news |            | Opinion |            |
|---------------------|--------------------|------------|---------|------------|
|                     | $\mu_p$            | $\sigma_p$ | $\mu_p$ | $\sigma_p$ |
| Aggregator          | 0.44               | 0.051      | 0.44    | 0.092      |
| Direct              | 0.47               | 0.076      | 0.47    | 0.094      |
| Social              | 0.46               | 0.087      | 0.47    | 0.12       |
| Search              | 0.46               | 0.087      | 0.46    | 0.14       |

Table 3: Bayesian model estimates of ideological consumption by channel and subjectivity type. The column  $\mu_p$  indicates the corresponding entry of  $\vec{\mu}_p$ , and the column  $\sigma_p$  indicates the corresponding diagonal entry of the model-estimated covariance matrix  $\Sigma_p$ .

tions of the latent variables, but in practice we simply fit the model with marginal maximum likelihood.

As with the analysis in Section 2.1, the diagonal entries of the covariance matrix  $\Sigma_p$  yield estimates of segregation for each of the eight subjectivity-by-channel categories. In particular, letting  $\sigma_k^2$  denote the  $k$ -th diagonal entry of  $\Sigma_p$ , segregation in the  $k$ -th category is  $\sqrt{2}\sigma_k$ . Table 3 lists these diagonal parameter estimates.<sup>8</sup> The off-diagonal entries of  $\Sigma_p$  measure the relationship between categories of one’s ideological exposure. For example, after normalizing  $\Sigma_p$  to generate the corresponding correlation matrix, we find the correlation between social media-driven descriptive news and opinion is 0.71. The full correlation matrix is included in the Appendix.

To help visualize these model estimates, Figure 3 plots segregation across the four consumption channels, for both opinion and descriptive news. The size of the markers is proportional to total consumption within the corresponding channel, normalized separately for opinion and descriptive news. To ground the scale of the  $y$ -axis, we note that among the top 20 most popular news outlets, conservative share ranges from 0.30 for the liberal *BBC* to 0.61 for the conservative *Newsmax*.

Figure 3 indicates that social media is indeed associated with higher segregation than direct browsing. For descriptive news this effect is subtle, with segregation increasing from 0.11 for direct browsing to 0.12 for articles linked to from social media. However for opinion pieces, the effect is more pronounced, rising from 0.13 to 0.17. It is unclear whether this increased segregation is the effect of algorithmic

---

<sup>8</sup>Given the large sample size, all estimates are statistically significant well beyond conventional levels.

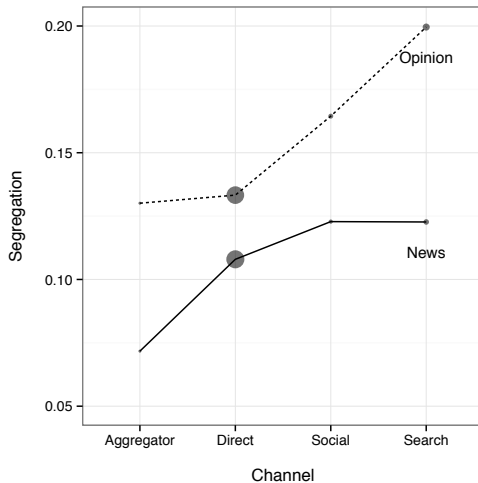


Figure 3: Estimates of ideological segregation across consumption channels. Point sizes indicate traffic fraction, normalized separately within the news and opinion lines.

filtering of the news stories appearing in one’s social feed (Pariser, 2011), the result of ideological similarity among one’s social contacts (Goel et al., 2010; McPherson et al., 2001), due to active selection by individuals of which stories in their feed to read (Garrett, 2009; Iyengar and Hahn, 2009; Munson and Resnick, 2010), or some combination of all three. In any case, however, our results are directionally consistent with worries that social media increase segregation.

We further find that search engines are associated with the highest levels of segregation among the four channels we investigate: 0.12 for descriptive news and 0.20 for opinion. Some authors have argued that web search personalization is a key driver of such effects (Pariser, 2011). There are two alternative explanations. The first is that users implicitly influence the ideological leanings of search results through their query formulation by, for example, issuing a query such as “obamacare” instead of “health care reform” (Borra and Weber, 2012). The second is that even when presented with the same search results, users are more likely to select outlets that share their own political ideology, especially for opinion content, has been found in laboratory studies (Garrett, 2009; Iyengar and Hahn, 2009; Munson and Resnick, 2010). While we cannot determine the relative importance of these factors, our findings do suggest that the relatively recent ability to instantly query large corpora of news articles—vastly expanding users’ choices sets—contributes to increased ideological

segregation, at least marginally for descriptive news and substantially for opinion stories.

At the other end of the spectrum, aggregators exhibit the lowest segregation. In particular, even though aggregators return personalized news results from a broad set of publications with disparate ideological leanings (Das et al., 2007), the overall effect is relatively low segregation. Though even for aggregators, segregation for opinion (0.13) is far higher than for descriptive news (0.07).

Given that our results are directionally consistent with filter bubble concerns, how is it that in Section 2.1 we found largely moderate overall levels of segregation? The answer is simply that only a relatively small fraction of consumption is of opinion pieces or from polarizing channels (social and search). Indeed, even after removing apolitical categories like sports and entertainment (which account for a substantial fraction of traffic), opinion still only constitutes 6% of consumption. Further, for both descriptive news and opinion, direct browsing is the dominant consumption channel (79% and 67%, respectively), dwarfing social media and search engines.

To help explain these results, we note that while sharing information is popular on social media, the dissemination of news is not its primary function. In fact, we find that only 1 in 300 clicks of links posted on Facebook lead to substantive news articles; rather, the vast majority of these clicks go to video and photo sharing sites. Moreover, we observe that even the most extreme segregation that we see (0.20 for opinion articles returned by search engines) is not, in our view, astronomically high. In particular, that level of segregation corresponds to the ideological distance between *Fox News* and *Daily Kos*, which represents meaningful differences in coverage (Baum and Groeling, 2008), but is within the mainstream political spectrum. Consequently, though the predicted filter bubble and echo chamber mechanisms do appear to increase online segregation, their overall effects at this time are somewhat limited.

### 2.3 Ideological Isolation

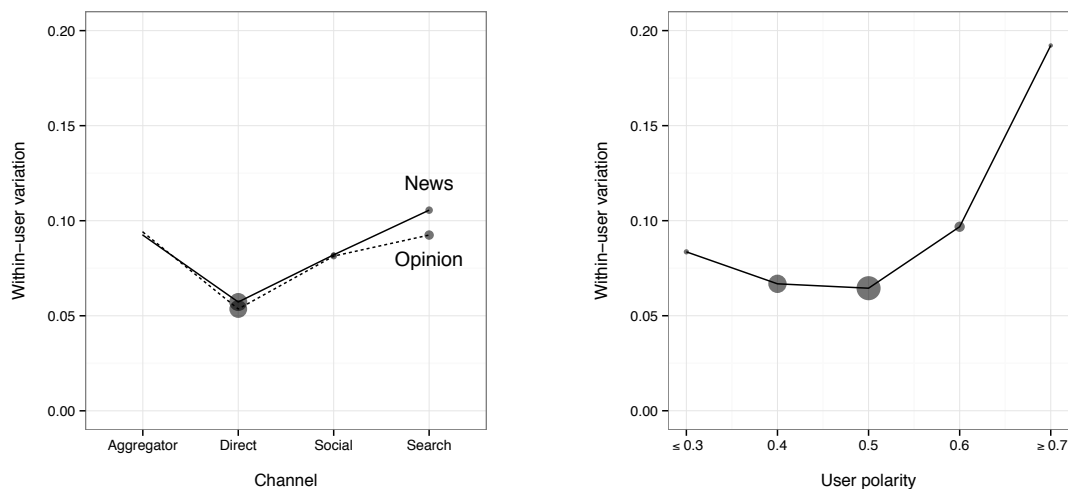
We have thus far examined segregation in terms of the distance between individuals' *mean* ideological positions. It could be the case, for example, that individuals typically consume content from a variety of ideological viewpoints, though ultimately skewing toward the left or right, leading to moderate overall segregation. Alterna-



tively, individuals might be tightly concentrated around their ideological centers, only rarely reading content from across the political spectrum. These two potential patterns have markedly different implications for the broader issues of political discussion and consensus formation (Benkler, 2006).

To investigate this question of within-user variation, we start by looking at the dispersion parameter  $\sigma_d$  in the overall consumption model described by Eqs. (1) and (2). We find that  $\sigma_d = 0.06$ , indicating that individuals typically read publications that are tightly concentrated ideologically.

This finding of within-user ideological concentration is driven in part by the fact that individuals often simply turn to a single news source for information: 78% of users get the majority of their news from a single publication, and 94% get a majority from at most two sources. As shown in the Appendix, however, this concentration effect holds even for those who visit multiple news outlets. Thus, even when individuals visit a variety of news outlets, they are, by and large, frequenting publications with similar ideological perspectives.



(a) Descriptive news (solid line) and opinion (dotted line). Point sizes indicate the relative fraction of traffic attributed to each source, normalized separately by category.

(b) Point sizes indicate the relative number of individuals in each polarity bin.

Figure 4: Within-user variation across consumption channel (a) and by mean polarity (b).

We now investigate ideological isolation across consumption channels and subjectivity categories. For each of the eight subjectivity-by-channel categories and for

each user, we first estimate the variance of the polarities of articles read by that user in that category.<sup>9</sup> For each category, we then average these individual-level estimates of variance (and take the square root of the result) to attain category-level estimates. Figure 4a plots these estimates of within-user variation by channel and subjectivity.

Across all four consumption channels, Figure 4a shows that descriptive and opinion articles are associated with similar levels of within-user variation. Social media, however, is associated with higher variation than direct browsing. Though this may at first seem surprising given that social media also has relatively high segregation, the explanation is clear in retrospect: when browsing directly, individuals typically visit only a handful of news sources, whereas social media sites expose users to more variety. Likewise, web search engines, while associated with high segregation, also have relatively high diversity. Finally, relatively high levels of within-user spread are observed for aggregators, as one might have expected.

We similarly examine within-user ideological variation as a function of user polarity (i.e., mean ideological preference). In this case, we first bin individuals by their polarity—as estimated in Section 2.1—and then compute the individual-level variation of article polarity, averaged over users in each group. As shown in Figure 4b, within-user variation is small and quite similar for users with polarity ranging from 0.3 to 0.6. Interestingly, however, the 2% of individuals with polarity of approximately 0.7 or more (significantly to the right of Fox News) exhibit a strikingly high within-user variation of 0.17.

This preceding result prompts a question: Does the high within-user variation we see among extreme right-leaning readers result from them reading articles from across the political divide, or are they simply reading a variety of right-leaning publications? More generally, across channels and subjectivity types, what is the relationship between within-user variation and exposure to ideologically divergent news stories? We conclude our analysis of ideological isolation by examining these questions.

We start by defining a news outlet as left-leaning (resp., right-leaning) if it is in the bottom (resp., top) third of the 100 outlets we consider; the full ranked list of publications is given in the Appendix. The left-leaning publications include

---

<sup>9</sup>For each category, we restrict to users who read at least two articles in that category.

newspapers from liberal areas, such as the *San Francisco Chronicle* and the *New York Times*, as well as blogs such as the *Huffington Post* and *Daily Kos*; the right-leaning set includes newspapers from historically conservative areas, such as the *Fort Worth Star-Telegram* and the *Salt Lake Tribune*, and online outlets such as *Newsmax* and *Breitbart*; and centrist publications (i.e., the middle third) include, for example, *Yahoo News* and *USA Today*. We refer to the combined collection of left- and right-leaning outlets as *partisan*.

For each user who reads at least two partisan articles, define his or her liberal exposure  $\ell_i$  to be the fraction of partisan articles read that are left-leaning. We define an individual’s *opposing partisan exposure*  $o_i = \min(\ell_i, 1 - \ell_i)$ . Thus, for individuals who predominantly read left-leaning articles,  $o_i$  is the proportion of partisan articles they read that are right-leaning, and vice-versa. We note  $o_i$  is only defined for the 82% of individuals in our sample that have read at least two partisan articles.

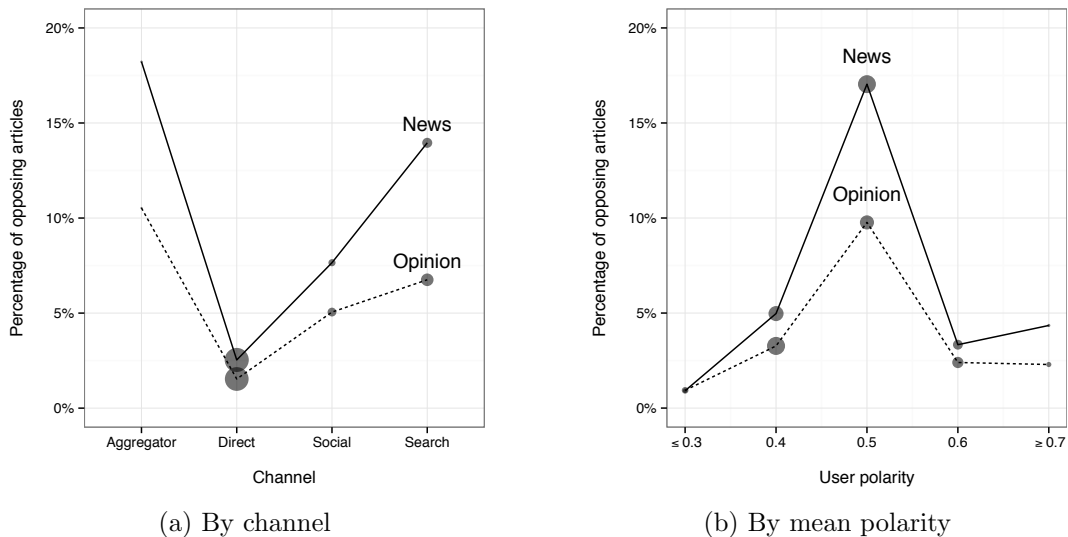


Figure 5: Opposing partisan exposure by channel (a) and polarity (b). Descriptive news (solid line) and opinion (dotted line). Point sizes indicate the relative fraction of traffic attributed to each source, normalized separately by article category.

Figure 5 shows average opposing partisan exposure, partitioned by article channel and subjectivity (Figure 5a), and by user polarity (Figure 5b).<sup>10</sup> For every

<sup>10</sup>To compute the estimates of average opposing partisan exposure shown in 5a,  $o_i$  is computed separately for each of the eight subjectivity-by-channel categories by restricting to the relevant articles, and limiting to users who read at least two partisan articles in that category.

subset we consider, only a small minority of articles—less than 20% in all cases, and less than 5% for all non-centrist users—comes from the opposite side of an individual’s preferred partisan perspective. Additionally, for every subset this opposing exposure is lower for opinion. Answering the question posed above, even extreme right-leaning readers have strikingly low opposing partisan exposure (3%); thus, their relatively high within-user variation is a product of reading a variety of centrist and right-leaning outlets, and not exposure to truly ideologically diverse content. In contrast, the relatively higher levels of within-user variation associated with social media and web search (Figure 4a) do translate to increased exposure to opposing viewpoints, though this effect is still small.

Summarizing our results on ideological isolation, we find that individuals generally read publications that are ideologically quite similar, and moreover, users that regularly read partisan articles are almost exclusively exposed to only one side of the political spectrum. In this sense, many, indeed nearly all, users exist in a so-called echo chambers. We note, however, two key differences between our findings and some previous discussions of this topic (Pariser, 2011; Sunstein, 2009). First, we show that while social media and search do contribute to segregation, the lack of within-user variation is primarily driven by direct browsing. Second, consistent with Gentzkow and Shapiro (2011), the outlets that dominate partisan news coverage are still relatively mainstream, ranging from *The New York Times* on the left to *Fox News* on the right; the more extreme ideological sites (e.g., *Breitbart*), which presumably benefited from the rise of online publishing, do not appear to qualitatively impact the dynamics of news consumption.

### 3 Discussion and Conclusion

Returning to our opening question—the effect of recent technological changes on ideological segregation—there are two competing theories. Some authors have argued that such changes would lead to “echo chambers” and “filter bubbles”, while others predicted these technologies would increase exposure to diverse perspectives. We addressed the issue directly by conducting one of the largest studies of online news consumption to date.

We showed that articles found via social media or web search engines are indeed

associated with higher ideological segregation than those an individual reads by directly visiting news sites. However, we also found, somewhat counterintuitively, that these channels are associated with greater exposure to opposing perspectives. Finally, we showed that the vast majority of online news consumption mimicked traditional offline reading habits, with individuals directly visiting the home pages of their favorite, typically mainstream, news outlets. We thus uncovered evidence for both sides of the debate, while also finding that the magnitude of the effects are relatively modest.

We conclude by noting some limitations of our study. First, as with past work (Gentzkow and Shapiro, 2010, 2011; Groseclose and Milyo, 2005), for methodological tractability we focus on the ideological slant of news outlets, as opposed to that of specific articles. As such, we would misinterpret, for example, the news preferences of an individual who primarily reads liberal articles from generally conservative sites. We suspect, however, that this type of behavior is relatively limited, in part because individuals typically visit ideologically similar news outlets, suggesting their own preferences are in line with those of the sites that they frequent. Second, we focus exclusively on news consumption itself, and not on the consequences such choices have on, for example, voting behavior or policy preferences.<sup>11</sup> Relatedly, social networks can impact political outcomes through means other than exposure to news, for instance by allowing users to broadcast their decision to vote (Gerber et al., 2008). Third, it is plausible the stories that have the greatest impact are disproportionately discovered via social networks or search engines, meaning the true impact of these channels is larger than the raw figures indicate. Fourth, and related to the previous point, as we have focused our study on the (natural) sub-population of active news consumers, it is unclear what impact recent technological changes have on the majority of individuals who have little exposure to the news, but who may get that limited amount largely from social media. Finally, we note that precisely defining causation in this setting is a difficult issue. For example, is the counterfactual thought experiment one in which social media or search engines do not exist? Or perhaps one imagines the change in experience of a single, prototypical individual who joins (or is prevented from joining) a social media site? Nevertheless, despite these limitations, we believe our findings provide an empiri-

---

<sup>11</sup>Establishing and measuring the causal effects of partisan news exposure is difficult, though not impossible (Prior, 2013).

cal starting point for understanding how novel means of news consumption affect ideological polarization.

## References

- Agichtein, E., Brill, E., and Dumais, S. (2006). Improving web search ranking by incorporating user behavior information. In *Proceedings of the 29th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 19–26. ACM.
- Bakshy, E., Rosenn, I., Marlow, C., and Adamic, L. (2012). The role of social networks in information diffusion. In *Proceedings of the 21st international conference on World Wide Web*, pages 519–528. ACM.
- Baron, D. P. (1994). Electoral competition with informed and uninformed voters. *American Political Science Review*, 88(01):33–47.
- Bates, D., Maechler, M., and Bolker, B. (2013). *lme4: Linear mixed-effects models using Eigen and Eigenfaces*. R package version 0.999999-2.
- Baum, M. A. and Groeling, T. (2008). New media and the polarization of american political discourse. *Political Communication*, 25(4):345–365.
- Benkler, Y. (2006). *The wealth of networks: How social production transforms markets and freedom*. Yale University Press.
- Borra, E. and Weber, I. (2012). Political insights: exploring partisanship in web search queries. *First Monday*, 17(7).
- Broder, A. (2002). A taxonomy of web search. In *ACM Sigir forum*, volume 36, pages 3–10. ACM.
- Dandekar, P., Goel, A., and Lee, D. T. (2013). Biased assimilation, homophily, and the dynamics of polarization. *Proceedings of the National Academy of Sciences*, 110(15):5791–5796.

- Das, A. S., Datar, M., Garg, A., and Rajaram, S. (2007). Google news personalization: scalable online collaborative filtering. In *Proceedings of the 16th international conference on World Wide Web*, pages 271–280. ACM.
- DellaVigna, S. and Kaplan, E. (2007). The Fox News effect: media bias and voting. *The Quarterly Journal of Economics*, 122(3):1187–1234.
- Downs, A. (1957). *An economic theory of democracy*. New York.
- Garrett, R. K. (2009). Echo chambers online?: Politically motivated selective exposure among internet news users. *Journal of Computer-Mediated Communication*, 14(2):265–285.
- Gelman, A. and Hill, J. (2007). *Data analysis using regression and multi-level/hierarchical models*. Cambridge University Press.
- Gentzkow, M. and Shapiro, J. M. (2006). Media bias and reputation. *Journal of Political Economy*, 114(2):280–316.
- Gentzkow, M. and Shapiro, J. M. (2010). What drives media slant? evidence from US daily newspapers. *Econometrica*, 78(1):35–71.
- Gentzkow, M. and Shapiro, J. M. (2011). Ideological segregation online and offline. *The Quarterly Journal of Economics*, 126(4):1799–1839.
- Gerber, A. S., Green, D. P., and Larimer, C. W. (2008). Social pressure and voter turnout: Evidence from a large-scale field experiment. *American Political Science Review*, 102(01):33–48.
- Glover, E. J., Flake, G. W., Lawrence, S., Birmingham, W. P., Kruger, A., Giles, C. L., and Pennock, D. M. (2001). Improving category specific web search by learning query modifications. In *Symposium on Applications and the Internet*, pages 23–32. IEEE.
- Goel, S., Hofman, J. M., and Sirer, M. I. (2012a). Who does what on the web: A large-scale study of browsing behavior. In *ICWSM*.
- Goel, S., Mason, W., and Watts, D. J. (2010). Real and perceived attitude agreement in social networks. *Journal of Personality and Social Psychology*, 99(4):611.

- Goel, S., Watts, D. J., and Goldstein, D. G. (2012b). The structure of online diffusion networks. In *Proceedings of the 13th ACM Conference on Electronic Commerce*, pages 623–638. ACM.
- Groseclose, T. and Milyo, J. (2005). A measure of media bias. *The Quarterly Journal of Economics*, 120(4):1191–1237.
- Hannak, A., Sapiezynski, P., Molavi Kakhki, A., Krishnamurthy, B., Lazer, D., Mislove, A., and Wilson, C. (2013). Measuring personalization of web search. In *Proceedings of the 22nd international conference on World Wide Web*, pages 527–538. International World Wide Web Conferences Steering Committee.
- Hosanagar, K., Fleder, D., Lee, D., and Buja, A. (2013). Will the global village fracture into tribes? recommender systems and their effects on consumer fragmentation. *Management Science*, 60(4):805–823.
- Iyengar, S. and Hahn, K. S. (2009). Red media, blue media: Evidence of ideological selectivity in media use. *Journal of Communication*, 59(1):19–39.
- Kwak, H., Lee, C., Park, H., and Moon, S. (2010). What is Twitter, a social network or a news media? In *WWW '10: Proceedings of the 19th international conference on World wide web*, pages 591–600, New York, NY, USA. ACM.
- LaCour, M. J. (2013). A balanced news diet, not selective exposure: Evidence from a direct measure of media exposure.
- Lassen, D. D. (2005). The effect of information on voter turnout: Evidence from a natural experiment. *American Journal of Political Science*, 49(1):103–118.
- Lawrence, E., Sides, J., and Farrell, H. (2010). Self-segregation or deliberation? blog readership, participation, and polarization in American politics. *Perspectives on Politics*, 8(01):141–157.
- Liu, D. C. and Nocedal, J. (1989). On the limited memory BFGS method for large scale optimization. *Mathematical programming*, 45(1-3):503–528.
- Manning, C. D. and Schütze, H. (1999). *Foundations of statistical natural language processing*, volume 1. MIT Press.



- McPherson, M., Smith-Lovin, L., and Cook, J. M. (2001). Birds of a feather: Homophily in social networks. *Annual review of sociology*, pages 415–444.
- Messing, S. and Westwood, S. J. (2012). Selective exposure in the age of social media: Endorsements trump partisan source affiliation when selecting news online. *Communication Research*, page 0093650212466406.
- Moscovici, S. and Zavalloni, M. (1969). The group as a polarizer of attitudes. *Journal of personality and social psychology*, 12(2):125.
- Mullainathan, S. and Shleifer, A. (2005). The market for news. *American Economic Review*, pages 1031–1053.
- Munson, S. A. and Resnick, P. (2010). Presenting diverse political opinions: how and how much. In *Proceedings of the SIGCHI conference on human factors in computing systems*, pages 1457–1466. ACM.
- Myers, D. G. and Bishop, G. D. (1970). Discussion effects on racial attitudes. *Science*.
- Obendorf, H., Weinreich, H., Herder, E., and Mayer, M. (2007). Web page revisitation revisited: implications of a long-term click-stream study of browser usage. In *Proceedings of the SIGCHI conference on Human factors in computing systems*, pages 597–606. ACM.
- Pariser, E. (2011). *The filter bubble: What the Internet is hiding from you*. Penguin UK.
- Pennebaker, J. W., Chung, C. K., Ireland, M., Gonzales, A., and Booth, R. J. (2007). The development and psychometric properties of liwc2007. *Austin, TX, LIWC. Net*.
- Pennebaker, J. W., Francis, M. E., and Booth, R. J. (2001). Linguistic inquiry and word count: Liwc 2001. *Mahway: Lawrence Erlbaum Associates*, page 71.
- Prior, M. (2013). Media and political polarization. *Annual Review of Political Science*, 16:101–127.

Schkade, D., Sunstein, C. R., and Hastie, R. (2007). What happened on deliberation day? *California Law Review*, pages 915–940.

Spears, R., Lea, M., and Lee, S. (1990). De-individuation and group polarization in computer-mediated communication. *British Journal of Social Psychology*, 29(2):121–134.

Sunstein, C. R. (2009). *Republic.com 2.0*. Princeton University Press.

Tewksbury, D. (2005). The seeds of audience fragmentation: Specialization in the use of online news sites. *Journal of broadcasting & electronic media*, 49(3):332–348.

# A For Online Publication

## A.1 Ideological Segregation on Twitter

Our main analysis investigated a variety of channels through which individuals read the news, but it was limited to a particular opt-in sample of individuals. In this supplementary section, we augment our analysis by examining the news consumption habits of a nearly complete set of users on one specific social information channel, Twitter, one of the largest online social networks, and arguable the largest designed primarily for information discovery and dissemination, as exemplified by their instructions to users to “simply find the accounts you find most compelling and follow the conversations.”<sup>12</sup>

The Twitter and toolbar datasets differ on two additional substantively important dimensions. First, Internet Explorer and Twitter users are demographically quite different. For example, whereas Internet Explorer users are believed to be, on average, older than those in the general Internet population, Twitter users skew younger. In particular, 27% of 18–29 year-olds use Twitter, compared to 10% of those aged 50–64 (Pew Research, 2013). Second, because of differing levels of information in the two datasets, in the toolbar analysis we examine the articles that an individual *viewed*, whereas with Twitter we look at the articles that were merely *shared* with that individual, regardless of whether or not he or she read the story. Thus, given these differences, to the extent that our results extend to this setting, we can be further assured of the robustness of our findings.

To generate the Twitter dataset, we start with the nearly complete set of U.S.-located individuals who posted a tweet during the two-month period March–April, 2013.<sup>13</sup> We focus on accounts maintained and used by an individual (as opposed to corporate accounts), and so further restrict to those that receive content from (“follow”) between 10 and 1,000 users on the network. This process yields approximately 7.5 million individuals. Finally, similar to our restriction in the toolbar analysis, we limit to active news consumers, who received (i.e., followed individuals who posted) at least 10 front-section news articles and at least 2 opinion pieces.<sup>14</sup> In total, 1.5

---

<sup>12</sup>Twitter positions itself as a fully-customizable information portal, this quote comes from [www.twitter.com/about](http://www.twitter.com/about).

<sup>13</sup>Twitter offers the option of “protected accounts,” which are not publicly accessible. These accounts are rare and are not part of our study.

<sup>14</sup>As with the toolbar analysis, articles were classified as front-section news and opinion according

million users meet all of these restrictions.

We begin our analysis by estimating the distribution of user polarity. In this setting, user polarity is the typical polarity of the articles to which a user is exposed (i.e., articles that are posted by an account the user follows), where we recall that the polarity of an article is the conservative share of the outlet in which it was published. Since users on Twitter often receive news by following the accounts of major news outlets rather than accounts of actual individuals (Kwak et al., 2010), and since these news outlets typically post hundreds of articles per day, individuals in our sample are generally exposed to large numbers of news articles—4,008 on average during the two-month time frame we study. As a consequence, data sparsity is not a serious concern, which in turn significantly simplifies our estimation procedure. Specifically, for each Twitter user, we estimate polarity by simply averaging the polarities of the articles to which he or she is exposed.

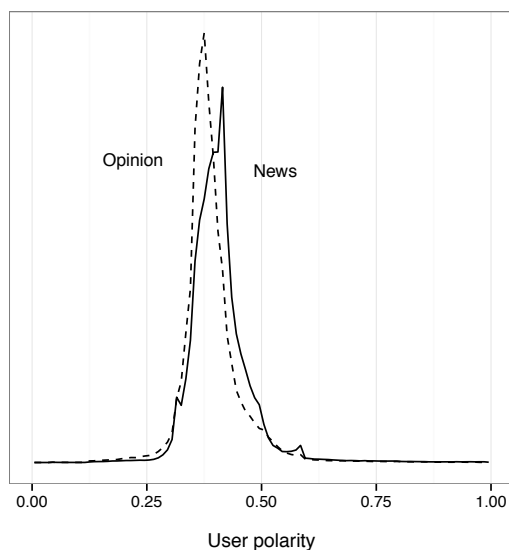


Figure 6: Distribution of individual-level polarity for Twitter users, where an individual’s polarity score is the average conservative share of news outlets to which he or she is exposed, computed separately for descriptive news articles (solid line) and opinion pieces (dashed line).

Figure 6 shows the resulting distribution of user polarity, where we separately plot the user polarity distribution computed for descriptive news articles (solid line) to the methods described in Section 1.

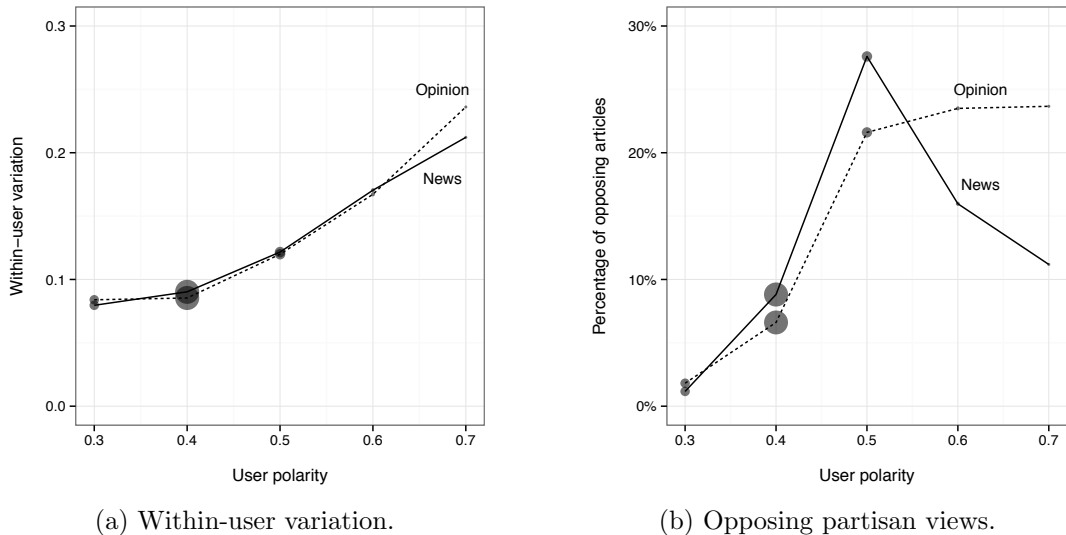


Figure 7: Within-user variation (a) and opposing partisan views (b) on Twitter, as a function of individual-level polarity. The sizes of the points indicate the relative number of individuals in each polarity bin, normalized separately for front-section news (solid line) and opinion (dashed line).

and opinion stories (dashed line). This plot illustrates two points. First, despite a slight leftward ideological skew relative to toolbar users, the bulk of Twitter users exhibit quite moderate news preferences. For example, 70% of Twitter users have polarity scores between 0.35 and 0.45, ranging from *The Huffington Post* to *CBS*. Second, segregation is correspondingly moderate, 0.10, and remarkably similar to our estimate from the toolbar data (0.11). Thus, despite the relative ease with which individuals may elect to follow politically extreme news publishers, and despite the worry that algorithmic recommendations of whom to follow could spur segregation, ideological segregation on Twitter looks very much like what we observe in direct navigation web browsing.

We investigate the exposure distribution further with two individual-level metrics: (1) within-user variation, defined as the standard deviation of the polarities of articles to which an individual is exposed; and (2) opposing partisan views, defined as the fraction of partisan articles from an individual’s less preferred ideological perspective. The results are plotted in Figure 7, as a function of user polarity.

As indicated by Figure 7a, average within-user variation—averaged over all individuals in our sample of Twitter users—is 0.10, significantly higher than the 0.05

we observed for direct web browsing, but comparable to the 0.09 we found for articles obtained through aggregators (Figure 4a), consistent with the general view of Twitter as a custom aggregator. Further, as we saw before, within-user variation increases substantially as we move to the conservative end of the spectrum; that is, individuals who on average consume more conservative content also tend to consume content from a wider variety of ideological viewpoints.

We plot opposing partisan exposure in Figure 7b, restricting to individuals who are exposed to at least two partisan articles (as we required in the toolbar analysis). Average opposing partisan exposure is 11%, very close to the 10% we observe in the toolbar dataset—the vast majority of an individual’s partisan views come from their preferred political side. However, a notable difference between the two datasets is that whereas in the toolbar data both left- and right-leaning individuals have little exposure to opposing views, on Twitter, right-leaning individuals have considerably more exposure to opposing views than left-leaning users. Though it is not entirely clear what is driving this effect, it is likely in part due to the overall leftward skew of Twitter, where it is thus harder for right-leaning individuals to isolate themselves from the majority view.

## A.2 News and opinion classifier

To train the news and opinion classifiers, we require datasets consisting of a representative set of articles known to be front-section news, and another known not to be (i.e., a sampling of articles from the categories we wish to filter out, hereafter referred to as “non-news”); we likewise require labeled examples of descriptive versus opinion articles. To generate these sets we make use of the fact that many popular publishers indicate an article’s classification in its URL (web address). For example, a prototypical story on *USA Today* (in this case, about U.S. embassy security) has the address <http://www.usatoday.com/story/news/world/2013/08/01/us-embassies-sunday-security/2609863/>, where “news/world” in the URL indicates the article’s category. Identifying these URL patterns for 21 news websites, we are able to produce 70,406 examples of front-section news and opinion, and 73,535 examples of non-news. We use the same approach (looking for URLs with the word “opinion”) to generate a separate training dataset to distinguish between opinion pieces and descriptive news articles.

Given these training datasets, we next build a natural language model. We first compute the 1,000 most frequently occurring words in our corpus of articles, excluding so-called stop words, such as “and”, “the”, and “of”.<sup>15</sup> We augment this list with a set of 39 first and third person pronouns (Pennebaker et al., 2007, 2001), since opinion pieces—unlike descriptive articles—are often written in the first person, and including such pronouns has been shown to improve performance (Glover et al., 2001). Each article is subsequently represented as a 1,039-dimensional vector, where the  $i$ -th component indicates the number of times the  $i$ -th word in our list appears in the article, normalized by the total number of words in the article. Using fractional scores rather than raw frequencies is a standard approach in natural language classification tasks for dealing with differences in article length (Manning and Schütze, 1999). To retain the predictive power of the pronouns, quotations are removed from the articles before representing them as vectors of relative word frequencies.

Having defined the predictors (i.e., the relative frequencies of various popular words), and having generated a set of labeled articles, we now use logistic regression to build the classifiers. Given the scale of the data, we fit the models with the L-BFGS algorithm (Liu and Nocedal, 1989), as implemented in the open-source machine learning package Vowpal Wabbit. Applying the fitted model to the entire collection of 4,127,140 articles in our corpus, we obtained 2,226,170 stories (46%) classified as front-section news or opinion, and of these 11% are classified as opinion. Note that as mentioned in the text, we use the classifier even for outlets that indicate the article category in the URL, which guards against differing editorial policies biasing the results.

The accuracy of our classifiers is quite high. When tested on a 10% hold-out sample of articles whose categories can be inferred from their URLs, the front-section news classifier obtains 96% accuracy, with 95% precision and 97% recall (where the positive class is news and the negative class is non-news). We also achieve good performance on a hand-labeled set of 100 randomly selected articles from the full corpus: 81% accuracy, 84% precision, and 79% recall. Accuracy for the opinion classifier is high as well: 92% on a hold-out set of URL-labeled articles,

---

<sup>15</sup>me, ive, myself, weve, we, wed, i, ill, were, well, mine, us, ourselves, lets, im, ours, our, my, id,shes, shed, himself, theyll, her, hes, theyve, them, hed, their, his, she, they, theyd, hers, shell, themselves, herself, him, and he.

with 96% precision and 76% recall (where the positive class is opinion and the negative class is descriptive news). On a randomly selected hand-labeled subset of 100 news articles accuracy is 88%, with 80% precision and 57% recall. Table 1 lists words with the highest positive and negative weights for both classifiers—the words accord with common intuition. While the overall performance of our classifiers is quite good, it is by no means perfect. However, we note that in many cases there is genuine ambiguity, even among human judges, as to what constitutes, for example, descriptive news versus opinion.

### A.3 Measuring the Political Slant of Publishers

Approaches for measuring the political slant of news outlets broadly fall into one of two categories: content-based and audience-based. Content-based approaches compare the entire body of published textual content from a source (rather than individual articles) to sources with known political slants. For example, Groseclose and Milyo (2005) use the co-citation matrix of newspapers and members of Congress referencing political think tanks. Similarly, Gentzkow and Shapiro (2010) use congressional speeches to identify words and phrases associated with a stance on a particular issue, and then tabulate the frequencies of such phrases in newspapers. Audience-based approaches, on the other hand, use the political preferences of a publication’s readership base to measure political slant (Gentzkow and Shapiro, 2011; Tewksbury, 2005). Empirical evidence suggests that audience and content-based measures of slant are closely related (Gentzkow and Shapiro, 2006; Mullainathan and Shleifer, 2005). In particular, Iyengar and Hahn (2009) show that individuals select media outlets based on the match between the outlet’s and their own political positions, and moreover, it has been shown that outlets tailor their coverage to match the preferences of their base (Baum and Groeling, 2008; DellaVigna and Kaplan, 2007; Gentzkow and Shapiro, 2010).

Here we use an audience-based measure of news outlet slant. Specifically, we estimate the fraction of each news outlet’s readership that voted for the Republican candidate in the most recent presidential election (among those who voted for one of the two major-party candidates), which we call the outlet’s *conservative share*. Thus, liberal outlets have conservative shares less than about 50%, and conservative outlets have conservative shares greater than about 50%, in line with the usual left-



to-right ideology spectrum. To estimate the political composition of a news outlet’s readership, we make use of geographical information in our dataset. Specifically, each webpage view includes the county in which the user resides, as inferred by his or her IP address. With this information, we then measure how the popularity of a news outlet varies across counties as a function of the counties’ political compositions, which in turn yields the estimate we desire.

More formally, as a first approximation we start by assuming that the probability any user  $u_i$  views a particular news site  $s$  is solely a function of his or her party affiliation. Namely, for a fixed news site  $s$ , we assume Democrats view the site with probability  $p_d$  and Republicans view the site with probability  $p_r$ .<sup>16</sup> Reparameterizing so that  $\beta_0 = p_d$  and  $\beta_1 = p_r - p_d$ , we have

$$\mathbb{P}(u_i \text{ views } s) = \beta_0 + \beta_1 \delta_r(u_i) \tag{5}$$

where  $\delta_r(u_i)$  indicates whether user  $u_i$  is a Republican. Though our ultimate goal is to estimate  $\beta_0$  and  $\beta_1$ , we cannot observe an individual’s party affiliation. To circumvent this problem, for each county  $C_k$  we average (5) over all users in the county, yielding

$$\frac{1}{N_k} \sum_{u_i \in C_k} \mathbb{P}(u_i \text{ views } s) = \beta_0 + \beta_1 \frac{1}{N_k} \sum_{u_i \in C_k} \delta_r(u_i) \tag{6}$$

where  $N_k$  is the number of individuals in our sample who reside in county  $C_k$ .

While the left-hand side of (6) is observable—or at least is well approximated by the fraction of users in our sample that visit the news site—we cannot directly measure the fraction of Republicans in our sample (i.e., the sum on the right-hand side of (6) is not directly observable). To address this issue, we make a further assumption that our sample of users is representative of the county’s voting population, a population for which we can estimate party composition via the 2012 election returns. We thus have the following model:

$$P_k = \beta_0 + \beta_1 R_k \tag{7}$$

---

<sup>16</sup>As discussed later, by “Democrats” we in fact mean those who voted for the Democratic candidate in the last presidential election, and similarly for “Republicans.”

where  $P_k$  is the fraction of toolbar users in county  $C_k$  that visit the particular news outlet  $s$ , and  $R_k$  is the fraction of voters in county  $C_k$  that supported the Republican candidate, Mitt Romney, in the 2012 U.S. presidential election. To estimate the parameters  $\beta_0$  and  $\beta_1$  in (7), we fit a weighted least squares regression over the 2,654 counties for which we have at least one toolbar user in our sample, weighting each observation by  $N_k$  (i.e., the number of people in our dataset in county  $C_k$ ).

Clearly, (7) is only an approximation of actual behavior, with our specification ruling out the possibility that a generally liberal outlet is disproportionately popular in conservative counties. In particular, our model ignores the impact of local news coverage, with individuals living in the outlet’s county of publication visiting the site regardless of its political slant. Addressing this local effect, we modify our generative model to include an additional term. Namely, outside a news outlet’s local geographic region, we continue to assume that Democrats visit the site with probability  $p_d$ , and Republican’s visit the site with probability  $p_r$ , and we use (7)—fit on all non-local counties—to estimate  $p_r$  and  $p_d$ . Inside the local region we assume individuals visit the site with probability  $p_\ell$ , irrespective of their political affiliation, and we estimate  $p_\ell$  to be the empirically observed fraction of local toolbar users who visited the news outlet.

Finally, we approximate the conservative share  $p(s)$  of a news outlet  $s$  as the estimated fraction of Republicans that visit the site normalized by the total number of Democratic and Republican visitors. Specifically,

$$p(s) = \left[ N_\ell r_\ell p_\ell + p_r \sum_{k: C_k \text{ non-local}} N_k r_k \right] / \left[ N_\ell p_\ell + \sum_{k: C_k \text{ non-local}} N_k (r_k p_r + (1 - r_k) p_d) \right]$$

where  $N_k$  is the number of people in our dataset in county  $C_k$ ,  $p_d = \beta_0$ ,  $p_r = \beta_0 + \beta_1$ ,  $r_k$  is the two-party Romney vote share in county  $C_k$  (i.e., the number of Romney supporters divided by the total number of Romney and Obama supporters, excluding third party candidates), and parameters subscripted with  $\ell$  indicate values for the outlet’s local county of publication. This entire process is repeated for each of the 100 news outlets in our dataset.

## A.4 Additional Tables and Figures

Table 4: Conservative shares for the top 100 news outlets, ranked by share.

|    | <b>Domain</b>               | <b>Publication Name</b>      | <b>Conservative Share</b> |
|----|-----------------------------|------------------------------|---------------------------|
| 1  | timesofindia.indiatimes.com | Times of India               | 0.04                      |
| 2  | economist.com               | The Economist                | 0.12                      |
| 3  | northjersey.com             | North Jersey.com             | 0.14                      |
| 4  | ocregister.com              | Orange Country Register      | 0.15                      |
| 5  | mercurynews.com             | San Jose Mercury News        | 0.17                      |
| 6  | nj.com                      | NewJersey.com†               | 0.17                      |
| 7  | sfgate.com                  | San Francisco Chronicle      | 0.19                      |
| 8  | baltimoresun.com            | Baltimore Sun                | 0.19                      |
| 9  | courant.com                 | Hartford Courant             | 0.22                      |
| 10 | jpost.com                   | Jerusalem Post (EN-Israel)   | 0.25                      |
| 11 | prnewswire.com              | PR Newswire                  | 0.27                      |
| 12 | sun-sentinel.com            | South Florida Sun Sentinel   | 0.27                      |
| 13 | nationalpost.com            | National Post (CA)           | 0.28                      |
| 14 | thestar.com                 | Tornoto Star                 | 0.28                      |
| 15 | bbc.co.uk                   | BBC (UK)                     | 0.30                      |
| 16 | wickedlocal.com             | Wicked Local (Boston)        | 0.30                      |
| 17 | nytimes.com                 | New York Times               | 0.31                      |
| 18 | independent.co.uk           | The Independent              | 0.32                      |
| 19 | philly.com                  | Philadelphia Herald          | 0.32                      |
| 20 | hollywoodreporter.com       | Hollywood Reporter           | 0.33                      |
| 21 | miamiherald.com             | Miami Herald                 | 0.35                      |
| 22 | huffingtonpost.com          | Huffington Post              | 0.35                      |
| 23 | guardian.co.uk              | The Guardian                 | 0.37                      |
| 24 | washingtonpost.com          | Washington Post              | 0.37                      |
| 25 | online.wsj.com              | Wall Street Journal          | 0.39                      |
| 26 | news.com.au                 | News.com (AU)                | 0.39                      |
| 27 | dailykos.com                | Daily Kos                    | 0.39                      |
| 28 | bloomberg.com               | Bloomberg                    | 0.39                      |
| 29 | dailyfinance.com            | Daily Finance                | 0.39                      |
| 30 | syracuse.com                | Syracuse Gazette             | 0.39                      |
| 31 | usnews.com                  | US News and World Report     | 0.39                      |
| 32 | timesunion.com              | Times Union (Albany)         | 0.40                      |
| 33 | time.com                    | Time Magazine                | 0.40                      |
| 34 | reuters.com                 | Reuters                      | 0.41                      |
| 35 | telegraph.co.uk             | Daily Telegraph (UK)         | 0.41                      |
| 36 | businessweek.com            | Business Week                | 0.42                      |
| 37 | cnn.com                     | CNN                          | 0.42                      |
| 38 | politico.com                | Politico                     | 0.42                      |
| 39 | theatlantic.com             | The Atlantic                 | 0.42                      |
| 40 | nationaljournal.com         | National Journal             | 0.43                      |
| 41 | altnet.org                  | Altnet                       | 0.43                      |
| 42 | ajc.com                     | Atlanta Journal Constitution | 0.44                      |
| 43 | forbes.com                  | Forbes                       | 0.44                      |
| 44 | seattletimes.com            | Seattle Times                | 0.44                      |
| 45 | rawstory.com                | The Raw Story                | 0.44                      |
| 46 | newsday.com                 | News Day                     | 0.44                      |
| 47 | cbsnews.com                 | CBS                          | 0.45                      |
| 48 | rt.com                      | Russia Today                 | 0.45                      |
| 49 | theepochtimes.com           | The Epoch Times              | 0.46                      |
| 50 | latimes.com                 | Los Angeles Times            | 0.47                      |

|     | <b>Domain</b>          | <b>Publication Name</b>                               | <b>Conservative Share</b> |
|-----|------------------------|---|---------------------------|
| 51  | csmonitor.com          | Christian Science Monitor                             | 0.47                      |
| 52  | realclearpolitics.com  | Real Clear Politics                                   | 0.47                      |
| 53  | usatoday.com           | USA Today   | 0.47                      |
| 54  | cnbc.com               | CNBC  | 0.47                      |
| 55  | dailymail.co.uk        | The Daily Mail (UK)                                   | 0.47                      |
| 56  | mirror.co.uk           | Daily Mirror (UK)                                     | 0.47                      |
| 57  | news.yahoo.com         | Yahoo! News   | 0.47                      |
| 58  | abcnews.go.com         | ABC News  | 0.48                      |
| 59  | upi.com                | United Press International                            | 0.48                      |
| 60  | chicagotribune.com     | Chicago Tribune                                       | 0.49                      |
| 61  | ap.org                 | Associated Press                                      | 0.50                      |
| 62  | nbcnews.com            | NBC News  | 0.50                      |
| 63  | suntimes.com           | Chicago Sun-Times                                     | 0.51                      |
| 64  | freep.com              | Detroit Free Press                                    | 0.52                      |
| 65  | azcentral.com          | Arizona Republics                                     | 0.53                      |
| 66  | tampabay.com           | Tampa Bay Times                                       | 0.54                      |
| 67  | orlandosentinel.com    | Orlando Sentinel                                      | 0.54                      |
| 68  | thehill.com            | The Hill  | 0.57                      |
| 69  | nationalreview.com     | The National Review                                   | 0.57                      |
| 70  | news.sky.com           | SKY   | 0.58                      |
| 71  | detroitnews.com        | Detroit News  | 0.59                      |
| 72  | express.co.uk          | The Daily Express (UK)                                | 0.59                      |
| 73  | weeklystandard.com     | The Weekly Standard                                   | 0.59                      |
| 74  | foxnews.com            | Fox News  | 0.59                      |
| 75  | washingtontimes.com    | Washington Times                                      | 0.59                      |
| 76  | jsonline.com           | Milwaukee Journal Sentinel                            | 0.61                      |
| 77  | newsmax.com            | Newsmax   | 0.61                      |
| 78  | factcheck.org          | factcheck.org   | 0.62                      |
| 79  | reason.com             | Reason Magazine                                       | 0.63                      |
| 80  | washingtonexaminer.com | Washington Examiner                                   | 0.63                      |
| 81  | ecanadanow.com         | E Canada Now  | 0.63                      |
| 82  | americanthinker.com    | American Thinker                                      | 0.65                      |
| 83  | twincities.com         | St. Paul Pioneer Press                                | 0.67                      |
| 84  | jacksonville.com       | Florida Times Union                                   | 0.67                      |
| 85  | opposingviews.com      | Opposing Views  | 0.67                      |
| 86  | chron.com              | Houston Chronicle                                     | 0.67                      |
| 87  | startribune.com        | Minneapolis Star Tribune                              | 0.68                      |
| 88  | breitbart.com          | Breitbart   | 0.70                      |
| 89  | star-telegram.com      | Ft. Worth Star-Telegram                               | 0.74                      |
| 90  | stltoday.com           | St. Louis Post-Dispatch                               | 0.75                      |
| 91  | mysanantonio.com       | San Antonio Express News                              | 0.77                      |
| 92  | denverpost.com         | Denver Post   | 0.80                      |
| 93  | triblive.com           | Pittsburg Tribune-Review                              | 0.85                      |
| 94  | sltrib.com             | Salt Lake Tribune                                     | 0.85                      |
| 95  | dallasnews.com         | Dallas Morning News                                   | 0.86                      |
| 96  | kansascity.com         | Kansas City Star                                      | 0.93                      |
| 97  | deseretnews.com        | Deseret News (Salt Lake City)                         | 0.94                      |
| 98  | topix.com              | Topix   | 0.96                      |
| 99  | knoxnews.com           | Knoxville News Sentinel                               | 0.96                      |
| 100 | al.com                 | Huntsville News/Mobile Press Register/Birmingham News | 1.00                      |

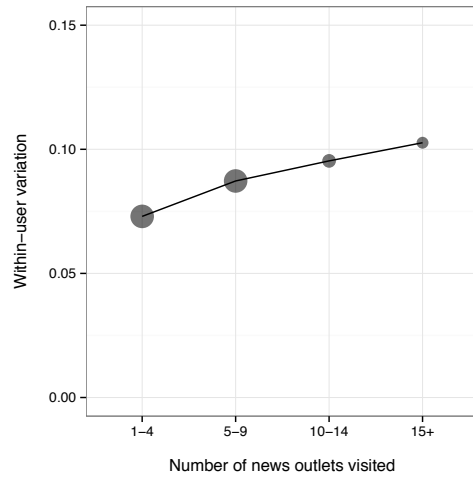


Figure 8: For a typical individual, within-user variation (i.e., standard deviation) of the conservative share of news outlets he or she visits, as a function of the number of outlets visited.

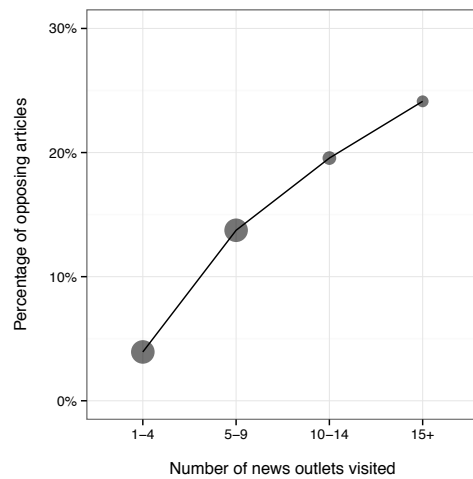


Figure 9: For a typical individual, fraction of partisan articles that are on the opposite side of the ideological spectrum from those he or she generally reads, as a function of the number of news outlets visited.

|         |            | News       |        |        |        | Opinion    |        |        |        |
|---------|------------|------------|--------|--------|--------|------------|--------|--------|--------|
|         |            | aggregator | direct | search | social | aggregator | direct | search | social |
| News    | aggregator | 0.0026     |        |        |        |            |        |        |        |
|         | direct     | 0.0007     | 0.0058 |        |        |            |        |        |        |
|         | search     | 0.0008     | 0.0033 | 0.0075 |        |            |        |        |        |
|         | social     | 0.0010     | 0.0043 | 0.0042 | 0.0075 |            |        |        |        |
| Opinion | aggregator | 0.0018     | 0.0013 | 0.0011 | 0.0010 | 0.0085     |        |        |        |
|         | direct     | 0.0007     | 0.0064 | 0.0039 | 0.0050 | 0.0024     | 0.0089 |        |        |
|         | search     | 0.0011     | 0.0038 | 0.0068 | 0.0048 | 0.0030     | 0.0057 | 0.0199 |        |
|         | social     | 0.0008     | 0.0043 | 0.0048 | 0.0072 | 0.0030     | 0.0064 | 0.0089 | 0.0135 |

Table 5: Variance-covariance matrix for the model used to estimate ideological consumption by channel and subjectivity type, as described in Eqs. (3) and (4).

|         |            | News       |        |        |        | Opinion    |        |        |        |
|---------|------------|------------|--------|--------|--------|------------|--------|--------|--------|
|         |            | aggregator | direct | search | social | aggregator | direct | search | social |
| News    | aggregator |            |        |        |        |            |        |        |        |
|         | direct     | 0.17       |        |        |        |            |        |        |        |
|         | search     | 0.18       | 0.51   |        |        |            |        |        |        |
|         | social     | 0.23       | 0.65   | 0.56   |        |            |        |        |        |
| Opinion | aggregator | 0.39       | 0.18   | 0.14   | 0.12   |            |        |        |        |
|         | direct     | 0.15       | 0.89   | 0.48   | 0.61   | 0.28       |        |        |        |
|         | search     | 0.16       | 0.35   | 0.56   | 0.4    | 0.23       | 0.43   |        |        |
|         | social     | 0.13       | 0.49   | 0.47   | 0.71   | 0.28       | 0.58   | 0.54   |        |

Table 6: Correlation matrix for the model used to estimate ideological consumption by channel and subjectivity type, as described in Eqs. (3) and (4).