

Government Websites As Data: A methodological pipeline with application to the websites of municipalities in the United States

Markus Neumann*

Fridolin Linder[†]

Bruce Desmarais[‡]

September 28, 2021

Abstract

The content of a government’s website is an important source of information about policy priorities, procedures, and services. Existing research on government websites has relied on manual methods of website content collection and processing, which imposes cost limitations on the scale of website data collection. In this research note, we propose that the automated collection of website content from large samples of government websites can offer relief from the costs of manual collection, and enable contributions through large-scale comparative analyses. We also provide software to ease the use of this data collection method. In an illustrative application, we collect textual content from the websites of over two hundred municipal governments in the United States, and study how website content is associated with mayoral partisanship. Using statistical topic modeling, we find that the partisanship of the mayor predicts differences in the contents of city websites that align with differences in the platforms of Democrats and Republicans. The application **illustrates** the utility of website content data extracted via our methodological pipeline.

1 Introduction

Government websites convey voluminous information about all aspects of government policymaking, policy implementation, and public deliberation. A substantial body of research has focused on the contents of government websites (e.g., Grimmelikhuijsen 2010; Wang et al. 2005; Osman et al. 2014; Eschenfelder et al. 1997). The conventional approach to data collection in projects focused on government websites involves manual content extraction from each website in

*Department of Political Science, The Pennsylvania State University, University Park, PA 16802, USA. Email: mvn5218@psu.edu. Corresponding author.

[†]Department of Political Science, Social Media and Political Participation Lab, New York University, New York, NY 10012, USA. Email: fridolin.linder@nyu.edu

[‡]Department of Political Science, The Pennsylvania State University, University Park, PA 16802, USA. Email: bdesmarais@psu.edu. This work was supported by the National Science Foundation [1320219, 1637089, 1641047].

the dataset. The manual content analysis of government websites is only feasible when researchers have access to substantial resources or focus on a small sample of websites. Consider the example of Feeney and Brown (2017). They analyzed the websites of 500 U.S. municipalities, and recorded information about the sites' coverage of, e.g., the city's social media accounts, council agenda, department descriptions, and police reports. They employed two researchers to code each of the 500 websites. If each researcher could complete 5–10 websites per hour, between hours spent on training and data collection, the budget to complete such a project would need to be on the order of thousands of dollars.

Our primary research objective is to introduce a methodological pipeline that will enable researchers to gather and analyze government website content without incurring the substantial costs associated with manual data collection. Since it is possible to use computer-automated methods to (1) access websites, (2) parse content, and (3) analyze and categorize content, we propose that automated methods can be combined to dramatically lower the cost and increase the scale of projects involving the analysis of government website content. Our pipeline ingests a set of urls that point to government websites, and outputs a corpus of plain text that is extracted from the files, of all types (e.g., HTML, DOC, PDF, and TXT) available at each URL.

Our secondary objective is to illustrate the application and use of our pipeline through a study of the relationship between the political party of a (U.S.) city's mayor, and the textual contents of the city's website. We gather and analyze a dataset that covers the textual contents of websites from over two hundred municipal governments in the United States. By studying the covariation of topical contents on these websites with the partisanship of the city mayors, we validate the utility of both the pipeline, and this specific dataset.

2 The study of Government Website Content

Existing research that analyzes government websites *manually* provides an indication of the potential benefits of automated analysis of web contents. Research on ‘e-governance’ evaluates government websites in terms of accessibility, ease-of-use, mobile accessibility, and overall function (e.g., Urban 2002; Tolbert et al. 2008; McNutt 2010; Armstrong 2011; Feeney and Brown 2017; Mossey et al. 2019). As an example, Grimmelikhuijsen and Welch (2012) study local government websites of Dutch municipalities to measure government transparency regarding air quality in the municipalities. In most research on website content, researchers manually visit each website and record whether a feature is present or not (Kaylor et al. 2001; Urban 2002; Jeffres and Lin 2006; Armstrong 2011; Dolson and Young 2012; Feeney and Brown 2017)¹. This approach is also used by the Rutgers E-Governance Institute, whose survey instrument has influenced the e-government literature at large. Holzer and Manoharan (2016) have conducted seven studies, between 2003 and 2016, of the largest city’s website in each of the world’s 100 nations with the most internet users. In each case, evaluators visit each website and manually score it according to 104 items. However, this level of thoroughness comes at a price—Holzer and Manoharan (2016) name a total of 127 website evaluators in their list of acknowledgements.

The manual approach is expensive, and is also subject to human error, judgement calls, and inter-coder reliability problems. Furthermore, researchers looking only for specific pieces of information are bound to miss a lot of content that doesn’t fall within their pre-defined objectives. Moreover, city websites often contain thousands of pages (Urban 2002)², so even projects with substantial human resources might struggle to provide a complete understanding of individual websites. **It is now standard practice to combine human coding with machine-learning based content**

¹A related literature concerns the websites of politicians and their parties. By and large, researchers in this field also rely on hand-coding (Norris 2003; Druckman et al. 2009, 2010; Esterling et al. 2011), albeit with some exceptions, who do use targeted scrapers (Therriault 2010; Cryer 2019).

²Urban (2002) relies on a webcrawler to measure how many pages each city website is comprised of, which is also the first step of `wget`. In this way, his research is a precursor to our own, albeit without the actual analysis of each page.

classification (e.g., Sebők and Kacsuk 2021; Fowler et al. 2021). If researchers had access to machine-readable content associated with government websites, a smaller quantity of human coding could be used to develop machine classifiers for attributes of interest, which could in turn be used to classify much larger sets of websites than is possible to do with human-only coding methods. Unfortunately, scrapers conventionally used in the social sciences (e.g., Cryer 2019) rely on extracting text through specific HTML tags, which makes it difficult to apply to a large number of websites, because the process needs to be adjusted for each site. By contrast, the only input our pipeline requires is the site URL.

Though our main focus in the current article is on the application to government websites, it is important to note that the pipeline we develop will be useful in the context of many other research designs that involve the collection of website content. Previous work in other domains has involved computational content analysis of websites. For example, Coyle and Nguyen (2020) analyze the content of company websites to find manufacturing companies that use contract manufacturing arrangements in their business models. Klüver and Mahoney (2015) use the textual content presented on interest group websites to classify the groups' issue agendas and constituencies.

3 Application: Mayoral Partisanship and U.S. Municipal Government Website Text

Though government websites serve largely instrumental service-delivery purposes, they also offer officials a prime venue via which to communicate policy goals and accomplishments, which inevitably reflect officials' politics. In the current paper, we focus on the running example of the reflection of mayoral partisanship on municipal government websites. A substantial body of research has found that the partisanship of the mayor affects city governance along multiple dimensions of spending and policy attention (Gerber and Hopkins 2011; de Benedictis-Kessner and Warshaw 2016; Einstein and Glick 2016; Marion and Oliver 2013). Official city websites allow

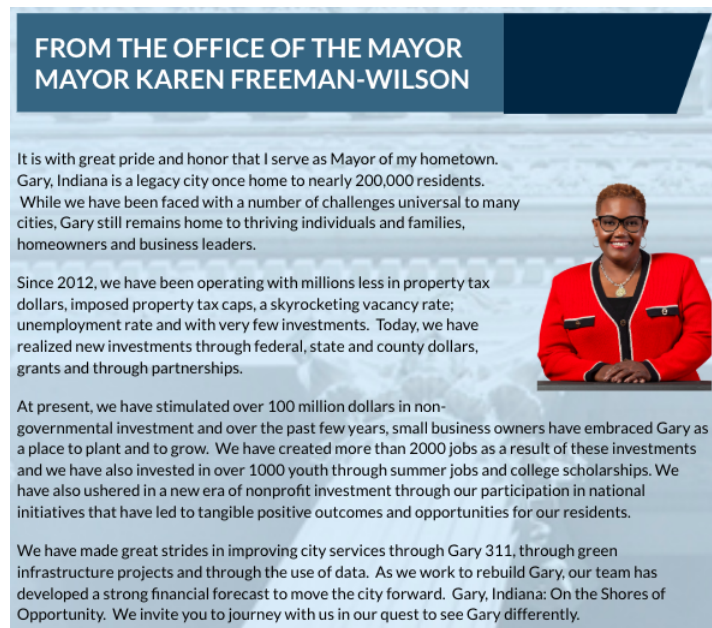


Figure 1: Screenshot from the homepage at <https://garyin.us/>, accessed on 05/22/2019. Image depicts Democratic mayor of Gary, IN, Karen Freeman-Wilson.

mayors to present their views and policy priorities to the public. In local politics, where campaign funds are low, this lends incumbents a crucial advantage in becoming more well-known among their constituencies (Stanyer 2008). Local government websites are frequently visited by the public (Thomas and Streib 2003). City websites can be used to communicate the stance of a mayor on social or economic programs. Consider the example of the Gary, Indiana homepage, depicted in Figure 1. This screenshot provides a clear example of the utility of a city website for communicating the mayor's policy priorities and accomplishments.

For data availability reasons, we focus our analysis of municipal websites on six states—Indiana, Louisiana, New York, Washington, California, and Texas. The websites were scraped in March 2018. The selection of states and cities is largely dictated by the presence of partisan mayors and availability of the relevant data. Municipal elections in Indiana and Louisiana are partisan across the board, so our sample is primarily focused on these two states. For Indiana and

Louisiana, all cities with a website are included, resulting in a considerably larger sample than for the other four states. New York and Washington do not have nominally partisan elections, but for a subset of cities, partisanship can be determined through contribution data (see appendix for more detail). California and Texas contain a number of large cities whose mayors are sufficiently well-known for their partisanship to be available. Our sample is well-balanced on a number of theoretically important dimensions. One, each of the four Census regions are represented with at least one state. Two, we have a fairly well-balanced sample with respect to the urban/rural cleavage. Furthermore, the sample is politically balanced—we have three blue states, and three red states. The partisan breakdown of city websites by state is provided in the appendix. This dataset of city website contents represents a contribution in the growing area of cross-municipality datasets covering local governments (e.g., Marschall and Shah 2013; Sumner et al. 2020). Details on the sources and methods of raw data collection can be found in the online appendix.

4 The Web-to-Text Pipeline

In this section, we describe our methodological pipeline, with which we take an archive of website files, and output a corpus of formatted plain text documents. We address three methodological challenges. First, though they contain significant amounts of text, websites are not comprised of clean plain text files. The first step is aimed at extracting clean plain text from this heterogeneous file base. The second step in our pipeline is to process the text to remove language that is effective at differentiating one website from another but is substantively uninformative. Finally, these tools need to work consistently across all of the websites in our corpus, in spite of the fact that relevant information is stored and structured in different ways. We make a software recommendation for each of these steps and gather most of them in our R package, `gov2text`. Some of the steps we take in this processing pipeline are universally applicable in the analysis of textual data, and some of them are most appropriate for the particular type of text analysis that we apply

to this data—statistical topic modeling. We will clarify this distinction as we describe steps in our pipeline.

4.1 Site to Text Conversion

4.1.1 File Type Detection

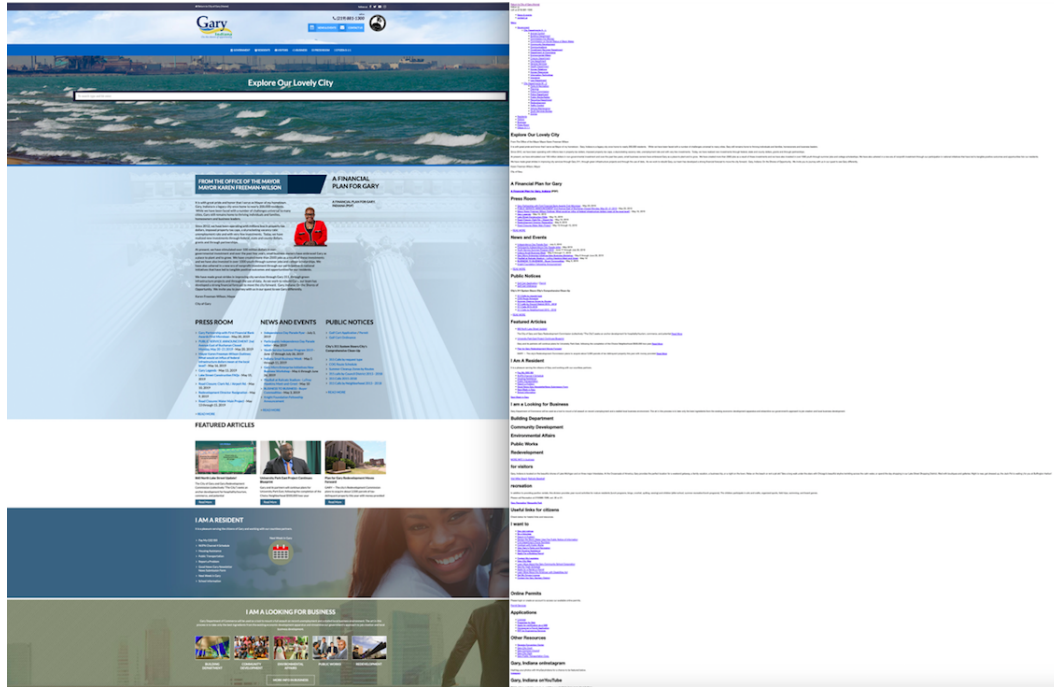
The format of a file impacts how text can be extracted from a document. For the most part, the file type of a document can be correctly determined by its extension. However, there are exceptions to this. For example, we found thousands of documents that ended in .html, when they were actually PDFs. A more accurate test for file type relies on the use of magic numbers, a short sequence of bytes at the start (and sometimes end) of files that is unique for each file type. We implement this method using the R package `wand` (Rudis et al. 2016).

4.1.2 Extracting Text from HTML

The HTML files that websites are comprised of contain a large amount of useful information, but also completely irrelevant text such as menus, navigational elements and other “boilerplate”. The side-by-side screenshots presented in Figure 2 convey the challenges presented by extracting content for text analysis for websites. The textual content that is substantive and unique to the Gary, IN homepage is the Mayor’s message depicted in Figure 1. The top row of Figure 2 presents the complete homepage, along with all of the text that can be naively extracted from the site. The Mayor’s message represents a relatively small fraction of the total text on the page.

To remove text that does not represent substantive content we rely on the `boilerpipe` classifier described in Kohlschütter et al. (2010), which is implemented through the R package `boilerpipeR`. The `boilerpipe` algorithm has been widely used in the natural language processing literature, but to our knowledge has not been previously used in the social sciences. The complete text extracted from the Gary, IN homepage using `boilerpipe` is depicted in the screenshot

(a) Naive Parsing



(b) Boilerpipe

From The Office of the Mayor Mayor Karen Freeman-Wilson

It is with great pride and honor that I serve as Mayor of my hometown. Gary, Indiana is a legacy city once home to nearly 200,000 residents. While we have been faced with a number of challenges universal to many cities, Gary still remains home to thriving individuals and families, homeowners and business leaders.

Since 2012, we have been operating with millions less in property tax dollars, imposed property tax caps, a skyrocketing vacancy rate; unemployment rate and with very few investments. Today, we have realized new investments through federal, state and county dollars, grants and through partnerships.

At present, we have stimulated over 100 million dollars in non-governmental investment and over the past few years, small business owners have embraced Gary as a place to plant and to grow. We have created more than 2000 jobs as a result of these investments and we have also invested in over 1000 youth through summer jobs and college scholarships. We have also ushered in a new era of nonprofit investment through our participation in national initiatives that have led to tangible positive outcomes and opportunities for our residents.

We have made great strides in improving city services through Gary 311, through green infrastructure projects and through the use of data. As we work to rebuild Gary, our team has developed a strong financial forecast to move the city forward. Gary, Indiana: On the Shores of Opportunity. We invite you to journey with us in our quest to see Gary differently.

Karen Freeman-Wilson, Mayor

Figure 2: The top image provides a side-by-side depiction of the entire homepage of <https://garyin.us/>, accessed on 05/22/2019, and complete/naive extraction of all of the text on the site. Bottom image provides the result of running <https://garyin.us/> through boilerpipe.

in the bottom row of Figure 2. We see that only the Mayor’s message is extracted, leaving the rest of the text as boilerplate.³

4.1.3 Extracting Text from PDF, DOC, DOCX and TXT

The extraction of information from other text-based file formats is more straightforward. To this end, we rely on `readtext` R package (Benoit and Obeng 2019), which is a wrapper for a set of parsers.⁴ The breakdown of all files by type is given in the online appendix. The most frequent file type besides HTML is PDF, from which we are able to extract a substantial amount of usable text. Files of type DOC, TXT, and DOCX, also occur regularly in our corpus and offer a considerable volume of textual data.

4.2 Preprocessing

Preprocessing is an important part of text-as-data research (Denny and Spirling 2018). Our advice given in this section, more than in any other, is specific to the problem of extracting meaningful textual information from municipal government websites, with the end goal of its use in a bag-of-words-based model. The challenge in conducting preprocessing for a comparative analysis of websites lies in the considerable, yet substantively irrelevant, differences between websites. For example, names of city officials, streets, and the city itself feature frequently in city documents. Since individual names recur at a much higher rate within a city than across the entire corpus, this would cause a topic model to cluster its topics by city.

To process the text we extract, we turn to a common method in natural language processing—part-of-speech (POS) tagging and named entity recognition (NER) to remove most names.⁵ For

³In the online appendix we show that without using `boilerpipe`, some of the most partisan ‘topics’ are simply website boilerplate text.

⁴`readtext` determines a document’s type solely through its ending—so the conversion described above is necessary.

⁵We retain laws, nationalities or religious or political groups, and works of art.

applications where the names of political actors might be of interest, this step is not recommended. Furthermore, we select words on the basis of their POS-tags, retaining only nouns, verbs, and adjectives. We keep proper nouns that also occur as nouns—this removes names, but retains titles such as “Police Chief” which can appear as proper nouns if they are followed by a name. Finally, we also conduct lemmatization to reduce words to their basic form.⁶ POS-tagging, NER and lemmatization are all implemented through `spacyr`. To deal with any leftover issues, we remove words with less than three characters, stopwords and non-English words (using the R package `Hunspell`). We also remove duplicate documents, which occur very frequently on websites.

After preprocessing, our corpus consists of 356,911 documents. In Table 1 we summarize all of the steps we take in gathering and processing our data. The summary includes a brief description of the step, the software packages used, and an indicator of whether the method is implemented in our R package, `gov2text`. The package is hosted on GitHub at (*blinded for anonymity*). We have also created an interactive tutorial for `gov2text` hosted on Google Colab (Bisong 2019).

Process	Software dependency	in <code>gov2text</code>
1. Assemble url list.	Selenium	no
2. Collect website files.	wget	no
3. Correct file extensions.	wand (Rudis et al. 2016)	yes
4. Discard website boilerplate.	boilerpipeR (Annau et al. 2015)	yes
5. Convert non-HTML files to text.	readtext (Benoit and Obeng 2019)	yes
6. Lemmatize text.	spacyr (Benoit and Matsuo 2018)	yes
7. Remove names.	spacyr	yes
8. Retain nouns, verbs, adjectives.	spacyr	yes
9. Stopword/number removal.	quanteda (Benoit et al. 2018)	yes
10. Retain only English words.	Hunspell (Ooms 2018)	yes
11. Removal of duplicate documents.	gov2text	yes

Table 1: Data collection and processing pipeline. Steps to collect and prepare text for topic modeling.

The biggest limitation in our pipeline, and an open area for future research, is the reliance on

⁶Lemmatization is similar to stemming, but works differently by taking grammar and surrounding words into account to identify the dictionary form of a word.

wget to gather website files. By using wget, we miss content that is displayed dynamically on websites using JavaScript. We investigated whether the presence of JavaScript was related to the amount of text we gathered from the website. We calculated the correlation between the number of `<script>` HTML tags on a city's website, which indicate the use of JavaScript on a site, and the number of text tokens we scrape from the site. This correlation is -0.059, which indicates a very weak relationship between the use of JavaScript and the amount of text scraped from the site.

5 Partisan Language on Municipal Websites

The policy priorities of U.S. mayors are shaped in large part by their partisan affiliations. For example, Einstein and Glick (2018) finds that Democratic mayors are much more likely than Republicans to support municipal-level redistributive policies, and Einstein et al. (2021) finds that Democratic mayors are more likely to support city-level policies aimed at addressing racial inequality. City mayors use government websites to present their policy priorities to the public. Consider an example; Formicola et al. (2003, p.55) document a significant website content change during a transition in mayoral administrations in the city of Indianapolis. Under the Republican mayor Stephen Goldsmith, voluminous content was added to the city website in connection with a faith-based initiative to create administrative partnerships with religious organizations. Faith-based initiatives represent a type of public-private partnership that is popular with Republicans (Saperstein 2003). When Democrat Bart Peterson took office in 2000, the material related to the faith-based initiative was removed from the website.

We illustrate the analysis of municipal website content by studying how differences in website content correlate with the partisanship of the city's mayor. As we reviewed above, the partisanship of the mayor has been found in past research to affect several features of city governance. However, Gerber and Hopkins (2011) note that, due to the constraints of state and national policies, municipalities lack discretion in many domains of governance. These constraints do not apply to

website contents. City governments have great discretion in composing their websites, and modifying website content is low cost.

To study content differences between government websites based on mayoral partisanship, we draw upon a recently-developed class model for text, the structural topic model (STM), developed by Roberts et al. (2014). Building on the conception of “topics” in conventional topic models (Valdez et al. 2018), in the STM a topic is a multinomial distribution defined on the word types in the corpus dictionary. The log-odds of the topic probabilities in each document-specific multinomial distribution over topics are drawn from a multivariate normal distribution in which the topic-specific means are determined by a linear regression function that associates document-attributed covariates with topics. For our primary empirical investigation, the STM provides a tool to estimate the relationship between the party of the city’s mayor and the prevalence of each topic. We also include the municipality’s population and median income as covariates. Further details on and results from our STM specification can be found in the online appendix.

5.0.1 Structural topic model results

The results are shown in Table 2. First, it is notable that the 95% credible interval includes zero in only seven of the sixty topics, indicating that the topics discussed on city websites varies systematically with the partisanship of the mayor. Many of the topics associated with Democrats fit with what we understand to be national party priorities. Topic **52**, on affordable housing, clearly resonates with the Democratic party’s appeal to low-income voters. Topic **6** (‘race’, ‘islander’, ‘census’, ‘female’) covers racial and gender identity issues. Similarly, employee rights and benefits are represented in topics **10** and **29**. Democrats also exhibit a strong preference for words related to public finances, such as Topic **58** (‘budget’, ‘revenue’, ‘expenditure’), Topic **45** (‘asset’, ‘actuarial’, ‘liability’, ‘financial’), Topic **35** (‘bond’, ‘obligation’, ‘proceeds’) as well as Topic **55** (‘taxable’, ‘deed’, ‘value’). We suspect that the association of Democratic mayors with finance-related terms

is indicative of a greater emphasis on the city's efforts to raise and spend money, and take credit for those efforts (e.g., the Gary, IN example in Figure 1). This finding is consistent with Einstein and Kogan (2015), who show that Democratic mayors tend to favor greater spending. A second, consistent Democratic focus appears to be law enforcement: The most Democratic topic, **59** ('burglary', 'robbery', 'theft', 'homicide') is clearly focused on crime. On the one hand, Democratic partisans have a more negative perception of the police, rating it considerably more negatively on the appropriate use of force and the equal treatment of minorities (Brown Jr 2017). On the other hand, the literature has also shown that cities with a higher Democratic vote share spend more on law enforcement, even after controlling for crime (Einstein and Kogan 2015).

City websites with Republican mayors, meanwhile, exhibit a pronounced focus on the essential functions of government. Basic utilities such as energy (Topic **20**), fire protection (Topic **51**), vaccination (Topic **2**), and sanitation (Topic **47**), are prevalent among cities with Republican mayors. These basic service topics cannot be found among topics prevalent in cities with Democratic mayors. Similarly, zoning issues figure prominently in the set of Republican topics (Topic **19**), which fits with the findings of Sorens (2018) that Republicans are more supportive of restrictive residential zoning rules. The Democratic topics also include one that is somewhat focused on zoning, Topic **39** ('downtown', 'mixed', 'density'), but emphasizes mixed-use zoning—a loosening of conventional single-use zoning rules.

In Figure 3 we illustrate the difference between Democrat and Republican mayor city documents in terms of the proportion of content predicted to be allocated to the respective topics. We illustrate the topics that are discussed above. Looking at the Democrat-leaning topics, we see that documents from Democratic mayor cities are approximately 1/2 of a percent more likely to include the affordable housing topic than are documents from Republican mayor cities. This may seem like a small effect, but 1/2 of a percent represents 30% of the average proportion of content represented by each of the sixty topics in this STM. Looking at the Republican portion of the figure, the

proportion of content appearing in documents from Republican mayor cities that is allocated to the topic on sanitation is approximately 2/3 of a percent higher than in Democratic mayor cities. This represents 40% of the average proportion of content allocated to a given topic.

6 Conclusion

We have developed a methodological pipeline for automatically gathering and preparing government websites for comparative content analysis. The key implication of our results is that researchers can use the methods we introduce to scale up, and lower the cost of, data collection for projects focused on the comparative analysis of website content. Through an application to the analysis of municipal websites in six different states, we find that government website contents are associated with the partisanship of the mayor in ways that would be expected based on the parties' national priorities and past research on the effects of mayoral partisanship on city governments.

We have highlighted projects focused on government and politician websites. However, the application of the methods we propose is not limited to government and politician websites. Our pipeline can be used in the context of any project that would benefit from the analysis of textual content from a large population of websites. There is a broad body of research programs that could make use of the data collection and processing pipeline we present.

References

- Annau, M., C. Kohlschuetter, and A. Clark (2015). *boilerpipeR: Interface to the Boilerpipe Java Library*. R package version 1.3.
- Armstrong, C. L. (2011). Providing a clearer view: An examination of transparency on local government websites. *Government Information Quarterly* 28(1), 11–16.

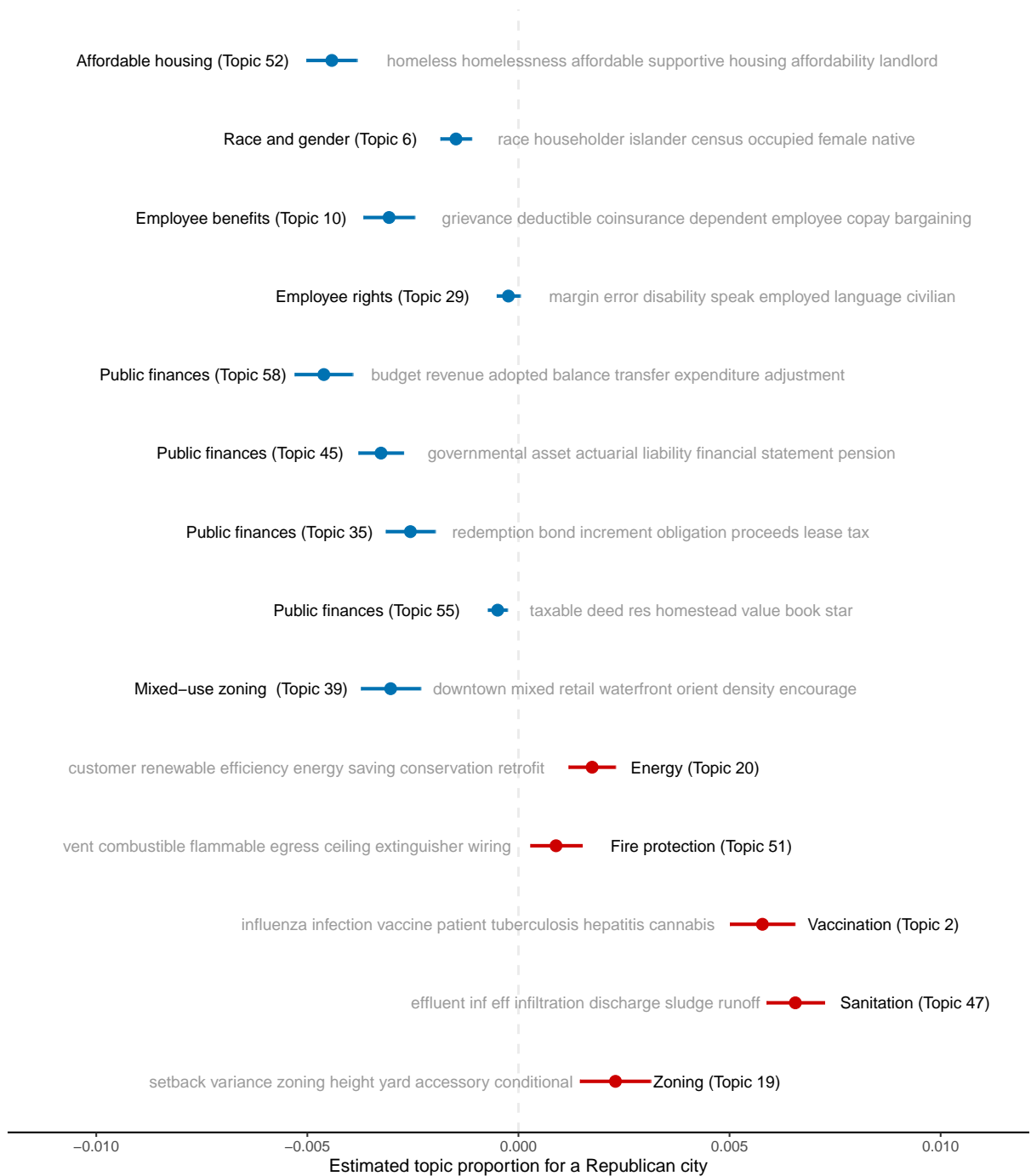


Figure 3: Topic differences in proportion to prevalence between Democrat and Republican mayor city websites. Depicting topics highlighted in discussion. Manually assigned labels are in black, top words for each topic as determined by the STM are in gray.

- Benoit, K. and A. Matsuo (2018). *spacyr: Wrapper to the 'spaCy' 'NLP' Library*. R package version 1.0.
- Benoit, K. and A. Obeng (2019). *readtext: Import and Handling for Plain and Formatted Text Files*. R package version 0.74.
- Benoit, K., K. Watanabe, H. Wang, P. Nulty, A. Obeng, S. Müller, and A. Matsuo (2018). quanteda: An r package for the quantitative analysis of textual data. *Journal of Open Source Software* 3(30), 774.
- Bisong, E. (2019). Google colab. In *Building Machine Learning and Deep Learning Models on Google Cloud Platform*, pp. 59–64. Springer.
- Brown Jr, L. T. (2017). Different lyrics, same song: Watts, ferguson, and the stagnating effect of the politics of law and order. *Harv. CR-CLL Rev.* 52, 305.
- Coyle, D. and D. Nguyen (2020). No plant, no problem? factoryless manufacturing, economic measurement and national manufacturing policies. *Review of International Political Economy*, 1–21.
- Cryer, J. (2019). Navigating identity in campaign messaging: The influence of race & gender on strategy in us congressional elections. In *2019 National Conference of Black Political Scientists (NCOBPS) Annual Meeting*.
- de Benedictis-Kessner, J. and C. Warshaw (2016). Mayoral partisanship and municipal fiscal policy. *The Journal of Politics* 78(4), 1124–1138.
- Denny, M. J. and A. Spirling (2018). Text preprocessing for unsupervised learning: why it matters, when it misleads, and what to do about it. *Political Analysis* 26(2), 168–189.

- Dolson, J. and R. Young (2012). Explaining variation in the e-government features of municipal websites: An analysis of e-content, e-participation, and social media features in canadian municipal websites. *Canadian Journal of Urban Research* 21(2), 1–24.
- Druckman, J. N., C. L. Hennessy, M. J. Kifer, and M. Parkin (2010). Issue Engagement on Congressional Candidate Web Sites, 2002—2006. *Social Science Computer Review* 28(1), 3–23.
- Druckman, J. N., M. Kifer, and M. Parkin (2009). Campaign Communications in U.S. Congressional Elections. *American Political Science Review* 103(03), 343–366.
- Einstein, K. L. and D. M. Glick (2016). Mayors, partisanship, and redistribution: Evidence directly from us mayors. *Urban Affairs Review*, 1078087416674829.
- Einstein, K. L. and D. M. Glick (2018). Mayors, partisanship, and redistribution: Evidence directly from us mayors. *Urban Affairs Review* 54(1), 74–106.
- Einstein, K. L., L. Godinez Puig, and S. Piston (2021). The pictures in their heads: How us mayors think about racial inequality. *Urban Affairs Review* 57(3), 611–642.
- Einstein, K. L. and V. Kogan (2015). Pushing the City Limits: Policy Responsiveness in Municipal Government. *Urban Affairs Review*, 1–30.
- Eschenfelder, K. R., J. C. Beachboard, C. R. McClure, and S. K. Wyman (1997). Assessing U.S. federal government websites. *Government Information Quarterly* 14(2), 173–189.
- Esterling, K. M., D. M. Lazer, and M. A. Neblo (2011). Representative communication: Web site interactivity and distributional path dependence in the us congress. *Political Communication* 28(4), 409–439.
- Feeney, M. K. and A. Brown (2017). Are small cities online? Content, ranking, and variation of U.S. municipal websites. *Government Information Quarterly* 34(1), 62–74.

- Formicola, J., M. Segers, and P. Weber (2003). *Faith-based Initiatives and the Bush Administration: The Good, the Bad, and the Ugly*. Rowman & Littlefield.
- Fowler, E. F., M. M. Franz, G. J. Martin, Z. Peskowitz, and T. N. Ridout (2021). Political advertising online and offline. *American Political Science Review* 115(1), 130–149.
- Gerber, E. R. and D. J. Hopkins (2011). When mayors matter: estimating the impact of mayoral partisanship on city policy. *American Journal of Political Science* 55(2), 326–339.
- Grimmelikhuijsen, S. G. (2010). Transparency of public decision-making: Towards trust in local government? *Policy & Internet* 2(1), 5–35.
- Grimmelikhuijsen, S. G. and E. W. Welch (2012). Developing and testing a theoretical framework for computer-mediated transparency of local governments. *Public administration review* 72(4), 562–571.
- Holzer, M. and A. Manoharan (2016). *Digital governance in municipalities worldwide (2015-2016): Seventh global e-governance survey: A longitudinal assessment of municipal websites throughout the world*. E-Governance Institute, National Center for Public Performance, Rutgers University.
- Jeffres, L. W. and C. A. Lin (2006). Metropolitan websites as urban communication. *Journal of Computer-Mediated Communication* 11(4), 957–980.
- Kaylor, C., R. Deshazo, and D. Van Eck (2001). Gauging e-government: A report on implementing services among american cities. *Government Information Quarterly* 18(4), 293–307.
- Klüver, H. and C. Mahoney (2015). Measuring interest group framing strategies in public policy debates. *Journal of Public Policy*, 223–244.

- Kohlschütter, C., P. Fankhauser, and W. Nejdil (2010). Boilerplate Detection using Shallow Text Features. In *Web Search and Data Mining*.
- Marion, N. E. and W. M. Oliver (2013). When the mayor speaks... mayoral crime control rhetoric in the top us cities: Symbolic or tangible? *Criminal justice policy review* 24(4), 473–491.
- Marschall, M. and P. Shah (2013). Local elections in america project. *Center for Local Elections in American Politics. Kinder Institute for Urban Research, Rice University.(Database)*.
- McNutt, K. (2010). Virtual policy networks: Where all roads lead to rome. *Canadian Journal of Political Science/Revue canadienne de science politique* 43(4), 915–935.
- Mossey, S., D. Bromberg, and A. P. Manoharan (2019). Harnessing the power of mobile technology to bridge the digital divide: a look at us cities’ mobile government capability. *Journal of Information Technology & Politics* 16(1), 52–65.
- Norris, P. (2003). Preaching to the Converted?: Pluralism, Participation and Party Websites. *Party Politics* 9(1), 21–45.
- Ooms, J. (2018). *hunspell: High-Performance Stemmer, Tokenizer, and Spell Checker*. R package version 3.0.
- Osman, I. H., A. L. Anouze, Z. Irani, B. Al-Ayoubi, H. Lee, A. Balci, T. D. Medeni, and V. Weerakkody (2014). Cobra framework to evaluate e-government services: A citizen-centric perspective. *Government Information Quarterly* 31(2), 243–256.
- Roberts, M. E., B. M. Stewart, D. Tingley, C. Lucas, J. Leder-Luis, S. K. Gadarian, B. Albertson, and D. G. Rand (2014). Structural topic models for open-ended survey responses. *American Journal of Political Science* 58(4), 1064–1082.

- Rudis, B., C. Zoulas, M. Rullgard, and J. Ong (2016). *wand: Retrieve 'Magic' Attributes from Files and Directories*. R package version 0.2.0.
- Saperstein, D. (2003). Public accountability and faith-based organizations: A problem best avoided. *Harvard Law Review* 116(5), 1353–1396.
- Sebők, M. and Z. Kacsuk (2021). The multiclass classification of newspaper articles with machine learning: The hybrid binary snowball approach. *Political Analysis* 29(2), 236–249.
- Sorens, J. (2018). The effects of housing supply restrictions on partisan geography. *Political Geography* 66, 44–56.
- Stanyer, J. (2008). Elected representatives, online self-presentation and the personal vote: Party, personality and webstyles in the united states and united kingdom. *Information, Community & Society* 11(3), 414–432.
- Sumner, J. L., E. M. Farris, and M. R. Holman (2020). Crowdsourcing reliable local data. *Political Analysis* 28(2), 244–262.
- Therriault, A. (2010). Taking Campaign Strategy Online: Using Candidate Websites to Advance the Study of Issue Emphases. *APSA 2010 Annual Meeting Paper, Available at SSRN: <https://ssrn.com/abstract=1643374>*, 1–23.
- Thomas, J. C. and G. Streib (2003). The new face of government: citizen-initiated contacts in the era of e-government. *Journal of public administration research and theory* 13(1), 83–102.
- Tolbert, C. J., K. Mossberger, and R. McNeal (2008). Institutions, policy innovation, and e-government in the american states. *Public administration review* 68(3), 549–563.
- Urban, F. (2002). Small town, big website? Cities and their representation on the internet. *Cities* 19(1), 49–59.

Valdez, D., A. C. Pickett, and P. Goodson (2018). Topic modeling: Latent semantic analysis for the social sciences. *Social Science Quarterly* 99(5), 1665–1679.

Wang, L., S. Bretschneider, and J. Gant (2005). Evaluating web-based e-government services with a citizen-centric approach. In *System Sciences, 2005. HICSS'05. Proceedings of the 38th Annual Hawaii International Conference on*, pp. 129b–129b. Ieee.

#	Top Word 1	Top Word 2	Top Word 3	Top Word 4	Top Word 5	Top Word 6	Tokens assigned	
49	artist	fun	music	beginner	player	prize	4565	<div></div>
46	chair	subcommittee	speaker	agenda	committee	commission	446	<div></div>
16	motion	second	adjourn	carry	unanimous	chairman	419	<div></div>
47	effluent	inf	eff	infiltration	discharge	sludge	751	<div></div>
21	everybody	think	something	thing	try	want	2609	<div></div>
2	influenza	infection	vaccine	patient	tuberculosis	hepatitis	2980	<div></div>
27	article	subsection	shall	franchisee	paragraph	meaning	658	<div></div>
30	subcontractor	bid	bidder	proposer	subcontract	bidding	512	<div></div>
12	craftsman	architecture	brick	distinctive	revival	storefront	1731	<div></div>
24	mail	fax	application	click	applicant	copy	367	<div></div>
34	playground	recreation	picnic	park	restroom	zoo	546	<div></div>
19	setback	variance	zoning	height	yard	accessory	453	<div></div>
26	mesa	canyon	via	odd	unidentified	paradise	1886	<div></div>
23	bag	recyclable	recyclables	reusable	vegetable	bait	2254	<div></div>
20	customer	renewable	efficiency	energy	saving	conservation	652	<div></div>
31	student	teacher	preschool	academic	kindergarten	youth	855	<div></div>
28	garland	assoc	association	firefighter	duke	xerox	480	<div></div>
50	trench	manhole	ductile	excavation	pipe	grout	1436	<div></div>
32	canceled	dwelling	suite	ave	tad	alteration	491	<div></div>
51	vent	combustible	flammable	egress	ceiling	extinguisher	1160	<div></div>
44	findings	tank	string	carcinogen	lust	sic	255	<div></div>
17	portfolio	micron	maturity	treasury	yield	investment	538	<div></div>
48	contributor	filer	officeholder	political	rouge	payee	293	<div></div>
5	draft	comment	review	revision	clarify	process	356	<div></div>
37	endorsed	endorse	rescue	assistant	analyst	technician	355	<div></div>
9	trust	revocable	planned	mfr	apportionment	exhibit	361	<div></div>
8	imp	assessor	taxpayer	petition	preliminary	determination	91	<div></div>
40	amt	invoice	acct	exp	unencumbered	encumbrance	116	<div></div>
57	councilman	introduced	alderman	whereas	resolved	councilwoman	615	<div></div>
11	obesity	sugary	epidemic	drink	calorie	sensible	96	<div></div>
15	credit	docket	app	post	download	month	61	<div></div>
3	wetland	specie	species	vernal	ecological	riparian	2293	<div></div>
29	margin	error	disability	speak	employed	language	180	<div></div>
43	medicare	payroll	blanket	contractual	undistributed	dept	322	<div></div>
42	incumbent	prep	batch	qualifier	analytical	examination	1091	<div></div>
55	taxable	deed	res	homestead	value	book	87	<div></div>
22	allocation	subtotal	admin	cost	yon	allocate	190	<div></div>
25	mitigation	impact	significant	adverse	environmental	measure	217	<div></div>
56	savings	neighborhood	village	excise	ltd	matrix	131	<div></div>
33	thence	east	south	corner	west	avenue	340	<div></div>
7	fugitive	bio	emission	coal	unmitigated	exhaust	773	<div></div>
18	perm	queue	delay	peak	adj	flt	187	<div></div>
54	license	licensee	citation	tow	fee	taxicab	710	<div></div>
6	race	householder	islander	census	occupied	female	160	<div></div>
60	bicycle	bike	pedestrian	route	sidewalk	bicyclist	561	<div></div>
14	accomplishment	grantee	narrative	outcome	grant	recipient	255	<div></div>
53	applied	col	dist	occupancy	monoxide	valuation	128	<div></div>
4	audit	auditor	procedure	timely	implemented	oversight	472	<div></div>
35	redemption	bond	increment	obligation	proceeds	lease	339	<div></div>
39	downtown	mixed	retail	waterfront	orient	density	419	<div></div>
10	grievance	deductible	coinsurance	dependent	employee	copay	583	<div></div>
38	para	persona	horas	bud	contracted	ante	1334	<div></div>
36	respondent	compare	figure	trend	appendix	satisfied	696	<div></div>
45	governmental	asset	actuarial	liability	financial	statement	235	<div></div>
41	complainant	allegation	defendant	offender	commander	complaint	1695	<div></div>
52	homeless	homelessness	affordable	supportive	housing	affordability	394	<div></div>
58	budget	revenue	adopted	balance	transfer	expenditure	176	<div></div>
13	initiative	outreach	strategy	leadership	engagement	focus	502	<div></div>
1	absent	preside	authorize	ordained	int	tag	377	<div></div>
59	burglary	robbery	theft	homicide	murder	gunshot	945	<div></div>

Table 2: Top words from a structural topic model with 60 topics and FREX scoring. Colors depict partisanship based on coefficient size. White cells are non-significant topics.