# Nationalized Elections, Localized Campaigns[*]
## Evidence from Televised U.S. Debates

Derek Holliday

Department of Political Science, UCLA

April 2020

### Abstract

Within states, partisan vote shares of presidential and gubernatorial candidates are becoming increasingly correlated, raising an important question for representation: Are election results increasingly similar because voters genuinely prefer candidates of a particular type or because voters see state races as referenda on national politics? Recent research focuses on explanations derived from voters, elites, and media. Largely ignored, however, is where these explanations intersect: campaigns and campaign rhetoric. Because national issues are efficient and proven motivators, state candidates may emphasize them to influence voters. I evaluate this possibility using text analysis and machine learning techniques on a corpus of presidential and gubernatorial electoral debates from 2000 to 2018. I find little evidence of nationalized campaign appeals; candidates tend to emphasize issues germane to their jurisdiction. These findings suggest similarities in outcomes independent of direct candidate appeals, necessitating renewed focus on the voter decision-making process in state and local politics.

# 1 Introduction

State-level election results have "nationalized" in recent years, with gubernatorial election results more closely resembling presidential results within states (Abramowitz and Webster 2016; Hopkins 2018; Sievert and McKee 2019). Extant literature has focused on voter- and elite-level mechanisms for nationalization such as the rise of negative partisanship amongst the electorate, dwindling access to local news and state political information, and declining heterogeneity in state party platforms (Abramowitz and Webster 2016; Hayes and Lawless 2018; Hopkins 2018; Martin and McCrain 2019). Less scholarly attention, however, has been given to potential campaign-level mechanisms, where voter and elite behaviors interact. Given the national media often characterizes *campaigns* as also being nationalized rather than just the *results*, this is an important omission from our understanding of nationalization. Do gubernatorial candidates increasingly emphasize positions on national-level political issues known to motivate votes during the campaign at the expense of state-level issues that appeal to constituents? If so, is campaign nationalization correlated with and/or partly responsible for the nationalization of electoral results?

In this study, I investigate the potential nationalization of gubernatorial campaigns using both supervised and unsupervised machine learning and quantitative text analysis techniques. Utilizing two original corpuses of state of the union/state of the state addresses and televised election debate transcripts from 2000 to 2018, I find limited evidence gubernatorial candidates mention national issues either at the expense of state issues or at a higher frequency than in the past. Generally, gubernatorial candidates talk about state issues in debates about five times as often as national issues. Additionally, I show the level of campaign nationalization is only mildly correlated with the nationalization of election results. These results have implications for our understanding of state campaign dynamics, the relative importance of voter- and elite-level mechanism of nationalization (such as voter information, issue prioritization, or party platforms), and the quality of state representation.

# 2 Why Nationalize Campaigns?

While scholarly work has largely focused on nationalized election *results*, media coverage of elections has used "nationalization" to mean the use of national issues and/or figures in state-level *campaigns*. For example, national media outlets characterized the 2019 gubernatorial races in Kentucky, Louisiana, and Mississippi as being nationalized due to an emphasis on impeachment of Donald Trump as a campaign issue and Trump's personal involvement in some of the races (Martin 2019; Manchester 2019; Rojas and Alford 2019). In 2018, the Virginia democratic gubernatorial primary was characterized as "a test of whether the boiling rage toward the new president [Trump] can be harnessed to win a state campaign" (Martin 2017). Other outlets gave similar appraisals of other races, including Washington in 2016 ("Inslee... was happy to nationalize the governor's race, sounding at many events like he was running against Trump"), West Virginia in 2011 (where the Republican Governor's Association spend $3.5 million in ad buys in an attempt to link the democratic candidate to Obamacare), and Texas in 2010 ("Mr. Perry turned the race into a referendum on federal spending") (Brunner 2016; Catanese 2011; McKinley Jr. 2010).

Of course, national media outlets calling certain gubernatorial elections nationalized may simply be a case of such outlets highlighting a small facet of the race that is of interest to a broader readership outside the state. Anecdotal evidence from the 2019 Kentucky gubernatorial debates, however, suggests nationalized appeals may indeed be a major part of candidates' campaign strategies. The following are excerpts from the October 15, 2019 debate between Kentucky gubernatorial candidates Matt Bevin (R-Incumbent) and Andy Beshear (D-State Attorney General):

> "National issues are also Kentucky issues." — Bevin, on whether this race is a referendum on national or state issues.

> "This governor is an extremist. He supports a total ban even for victims of rape and incest, something the President doesn't support, nor does his Attorney General candidate support." — Beshear

Both candidates invoke nationalized appeals here, both explicitly (in Bevin's case in the first comment) and implicitly by referencing President Trump's stance on abortion.

Given media portrayals of nationalized gubernatorial campaigns and some supporting anecdotal evidence, what incentives might candidates have to nationalize their campaigns?

First, gubernatorial candidates may find position-taking on national issues to be a stronger signal of "type" to the electorate. Research on political campaigns suggest candidates identify and campaign on salient political issue positions with majority support amongst the electorate (Carsey 2001; Vavreck 2009). Perhaps the simplest dividing line to campaign on is party identification. Voters are increasingly sorted along partisan lines on most national issues, strengthening in-group partisan identity and out-group partisan antipathy (Abramowitz and Webster 2016; Mason 2015). It is less clear if voters are as clearly divided along partisan lines for state-level political issues; Jensen et al. (2019) suggests fairly little partisan sorting on issues related to local economic development, but what counts as a state versus national political issue changes over time (Kousser 2014). It is plausible, however, that the easiest division for campaigns to exploit relates to national party identification and the issue positions the national parties espouse.

Campaigning on national issues may also come with the benefit of more crystallized issue positions in the electorate, preventing opposing campaigns changing voter opinions on what might otherwise be a winning campaign issue. Tesler (2015) suggests more crystallized beliefs (such as those pertaining to religious convictions or deeply held group-identities) are more likely to be primed than changed in elections. Experimental work by Mullinix (2016) and Zhang (2019) also give evidence that issue positions seen as more important and/or held with deeper conviction are more resistant to change. While voters may have deeply held beliefs on state-level issues, the rising importance of national party identification as a social identity and its convergence with other (religious, racial/ethnic) identities suggests national issue positions are more likely to be deeply held (Huddy and Bankert 2017). Furthermore, additional experimental work in the state context by Broockman and Butler (2017) suggests

voters are fairly easily persuaded to change opinions on state issues when contacted by their state representative.

While there are many plausible reasons to nationalize a gubernatorial campaign, there are equally plausible reasons to keep the campaign localized. The most obvious reason is voters may recognize a candidate running on nationalized appeals has no jurisdiction over the issues being discussed. Current research is divided on the extent to which voters hold politicians accountable for conditions under their jurisdiction; Arceneaux (2006) finds survey respondents tend to attribute credit/blame to offices which they (fairly accurately) assign functional responsibility to, whereas Brown (2010) finds partisanship moderates the attribution of functional responsibility and subsequent credit/blame. Benedictis-Kessner and Warshaw (N.d.) find some evidence for both sides using times series, cross-sectional models; voters routinely hold the president's party responsible for local economic conditions, but also hold governors accountable for such conditions. Therefore, the incentives for candidates to nationalize gubernatorial campaigns seem mixed.

# 3  Data and Methods

In this paper, I consider one avenue through which candidates may act on incentives to make nationalized issue appeals in campaigns: televised election debates. Specifically, I analyze the text of an original corpus of 405 electoral debates (93 presidential and 312 gubernatorial) between 2000 and 2018 retrieved from closed-captioned transcripts from the C-SPAN video archives to assess the proportion of national versus state rhetoric over time.[1] These debates are typically broadcast either directly on C-SPAN or through local public affiliates.

Research on gubernatorial debates is rare, but the few studies that have been conducted conclude candidates largely focus on policy positions rather than character (Benoit, Brazeal, and Airne 2007) and viewers of debates are often able to correctly identify the eventual

---

[1]Transcripts were retrieved using a combination of headless web browsing and scraping. Transcripts for non-closed-captioned videos are not available.

winner of the contest (Benjamin and Shapiro 2009). Research on the effects of presidential debates largely conclude such events have some short-term effect on candidate preference (Hillygus and Jackman 2003) and issue knowledge/salience (Benoit, Hansen, and Verser 2003). Therefore, debates provide a comparable campaign event across states and levels of government to examine the potential existence of nationalized issue appeals, especially in environments where baseline voter knowledge is low. Additionally, the debate context helps control for candidate-level confounders such as ideology, campaign resources, or campaign activity level that may bias results in a different context (such as television advertisements or social media).

I analyze the state and national content of these debates over time using a structural topic model (hereafter STM). A full technical description of STM is beyond the scope of this paper, so I will briefly cover its main functionality here.[2] STMs treat texts as "bags of words," disregarding sentence structure. Like the traditional latent Dirichlet allocation (LDA) approach, STMs assign words to topics and topics to documents probabilistically (Blei, Ng, and Jordan 2003). The output gives the word probabilities associated with each topic and the topic proportions for each document. STM builds upon LDA by allowing topic probabilities to vary according to researcher-specified covariates such as time or party ID. STMs offer a unique opportunity to analyze large quantities of text data in a systematic way, and have recently gained popularity in political science (Cryer 2019; Das et al. 2019; Parthasarathy, Rao, and Palaniswamy 2019).

One common pitfall of topic modeling is that researchers must subjectively apply meaning or interpretation to topic-word probabilities. For example, one can imagine a topic model applied to a corpus of Gospel and country music songs. One topic resulting from this model may include words like "God," "prayer," and "blessed," while another might include words like "truck," "horse," and "farm." A researcher may justifiably associate these topics with Gospel and country songs (respectively) and call the first topic "religiosity" and the second

---

[2]but see Roberts, Stewart, and Airoldi (2016)

"rural." The problem is that such intervention on the part of the researcher is not a result of the topic model but simply a subjective interpretation of the words in each topic. Indeed, it is rare that topics contain *only* words with such clear meanings. In many ways, topic modeling resembles principle component analysis in that it reduces the corpus to a series of frequent and meaningful words, but interpreting those words is still subject to researcher bias.

My solution is to remove myself from the interpretation of topical content by training a STM on a corpus known to contain state and national political content, extracting the topics associated with the state or national origin of each document, and applying the trained model to a test corpus (the debate transcripts) to measure the prevalence of the topics associated with state and national content over time. At no point do I intervene as a researcher to determine whether a topic has state or national content; this is done purely by the trained model.

For this strategy to succeed, I must collect a training corpus of documents where we know the topics being discussed are of state or national origin. For documents of national topical origin, I gather every public presidential address made between 1998 and 2019, relying heavily on State of the Union and inaugural addresses. For documents of state origin, I gather every available State of the State (or Commonwealth) Address or State Budget Address given by governors from 2000 to 2019.[3] These addresses are often mandated by state constitutions and deal specifically with the political issues that face each state, making them ideal for training a topic model. This yields an original corpus of 984 speeches (116 presidential, 868 gubernatorial). Following best practices for text data, for both the training and test corpuses I remove the first 750 characters (as they often involve only thanking specific members of the audience), stem, lemmatize, and stopword (Parthasarathy, Rao, and Palaniswamy 2019). Because speeches often contain direct references to the state or occasion each speech is from, it is plausible such words may end up in topics that become significant predictors of state or

---

[3]Speeches from 2000-2010 where gathered from a preexisting repository, but speeches from 2011 to 2019 were gathered mostly by hand via local news transcripts and archives of gubernatorial websites.

national content without referencing actual political issues. Therefore, I create a custom list of stopwords that include state names, names and nicknames for residents of certain states (e.g. Hoosiers), names of elected positions, and words common in transcriptions of public speeches (such as "applause," "laughter," and "crosstalk").

One remaining issue with the training corpus is its class imbalance; there are many more documents from governors than there are from presidents. This is not due to undersampling, as the corpus contains almost the universe of cases for each level of government, but rather because there are simply more gubernatorial speeches than there are presidential ones. This presents a problem because the resulting topic distribution will be skewed heavily toward state content, making future categorization of rare events (national content) more difficult. To solve this, I implement a Synthetic Minority Over-sampling Technique (SMOTE) to create more equal numbers of synthetic documents originating from presidents and governors (Chawla et al. 2002). The resulting training corpus contains 443 synthetic documents: 220 presidential and 223 gubernatorial speeches.[4]

## 4  Results

The STM model I apply to the synthetic training corpus contains covariates for year (pooled in two-year intervals) and level of government (state or national). I determine the optimal number of topics (70) using cross validation techniques recommended by Roberts, Stewart, and Airoldi (2016), although getting the "correct" number of topics here is less important than getting a set of topics distinctly related to state and national content. Figure 1 shows the words that appear with the highest probability for four topics highly associated with state or national origin. The nationally-associated topics include references to America, the nation as a whole, national institutions, and security, while state-associated topics have words associated with education, business, and budgeting. Because I do not analyze these word proportions closely to avoid researcher interference, this exercise is more to show that the

---

[4]See appendix A1 for additional details on the SMOTE process I apply to the training data.

topic model is returning topics with policy/political content rather than procedural words or other vocabulary that might be indicative of the origin of the document but not its content.
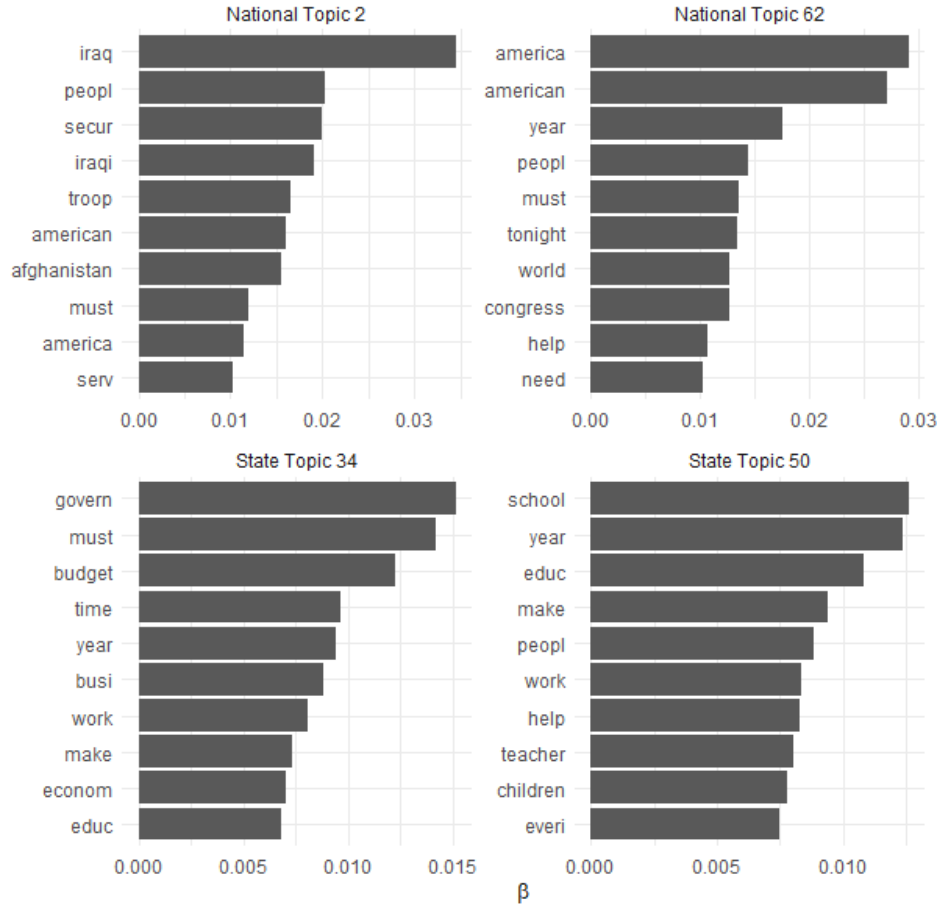


Figure 1: Highest Probability Words for Topics with most State/National Content

To estimate the extent to which the STM topics are indicative of state or national content, I include the 70 STM topic proportions for each synthetic (post-SMOTE) document in an L1 regularized regression (lasso) model predicting the state or national origin of the document. Lasso is a type of linear regression that avoids overfitting by using a data-driven process to remove less meaningful predictors. Lasso is optimal here since topic proportions less informative of state or national origin will have their model coefficients shrunk to zero, so we avoid making classifications based on topics not indicative of document origin. Then, I fit the STM model to the raw (pre-SMOTE) training documents to estimate the 70 topic

proportions for each document. I am then able to predict the origin of each raw document in the training corpus with the fitted lasso model to check whether the topic proportions do indeed help distinguish between state and national content. Figure 2 shows the resulting confusion matrix. These results are promising; the model accurately predicted the training document origin 98% of the time using only the fitted STM topic proportions, indicating the estimated topics are useful in distinguishing between state and national content.
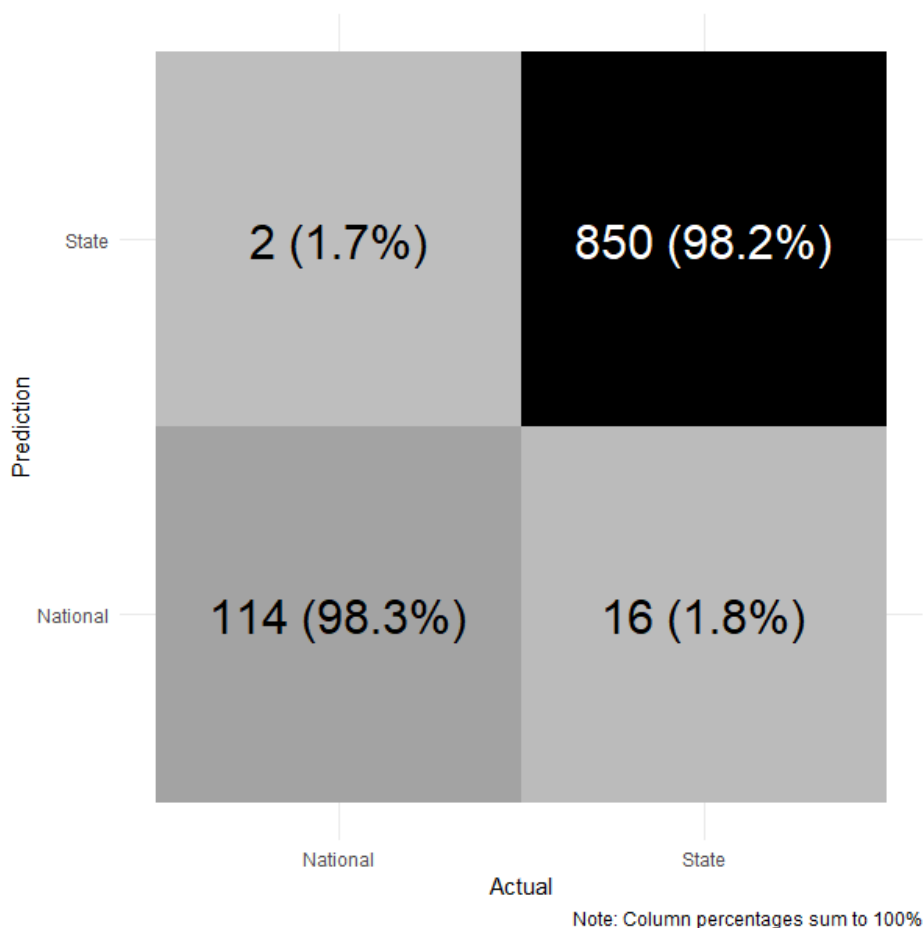


Figure 2: Lasso Confusion Matrix for Training Documents

The next step is to fit the trained STM model to the test corpus (the presidential and gubernatorial debates) to obtain the estimated topic proportions for each debate. To obtain a simple indicator of nationalization, I apply the trained lasso model to the estimated topic proportions in each debate to predict the state or national origin of debate. The results for

the pooled data are shown in the first panel confusion matrix in Figure 3. Overall, the model still does well in determining the origin of the debates; about 98% of presidential debates are classified as national, and 93% of gubernatorial debates are classified as state. This indicates gubernatorial debates still contain state-level content more so than national-level content. The right panel of Figure 3 shows the predictions for just gubernatorial debates over time. There is no obvious trend in debates being classified as more or less national, again indicating the predominance of state content in these debates.
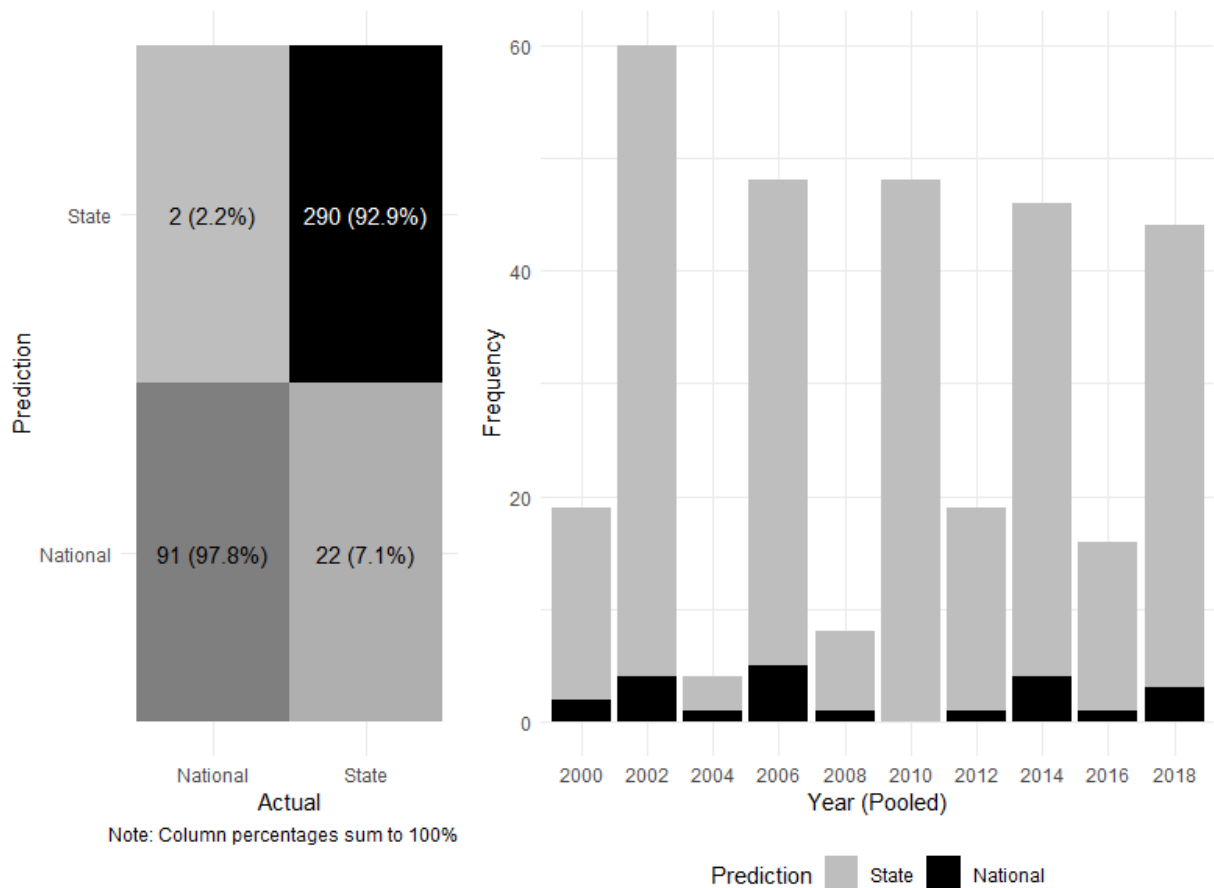


Figure 3: Lasso Confusion Matrix for Test Documents and State Classification Over Time

Of course, it may be the case that while gubernatorial candidates still spend a majority of their time discussing state-level issues, the relative weight given to national issues has increased over time. If this were true, the lasso model may still accurately predict the state

origin of gubernatorial debates despite rising nationalization. To determine if this may be occurring, I divide each debate in the test corpus into equal fifths and rerun the trained lasso model to the estimated topic proportions in each debate fifth. The results are shown in Figure 4. For the pooled results, there is greater evidence of nationalization, with almost 20% of the debate fifths classified as national, although the great majority are still correctly classified as state. This would be roughly equivalent to one fifth of every gubernatorial debate being spent discussing national content. Examining results over time, there is again little evidence of an increase in nationalized rhetoric. To the extent that rhetoric *is* nationalized, it has remained somewhat constant over the last two decades. If anything, the most nationalized rhetoric took place from 2004 to 2007 and has subsequently declined.
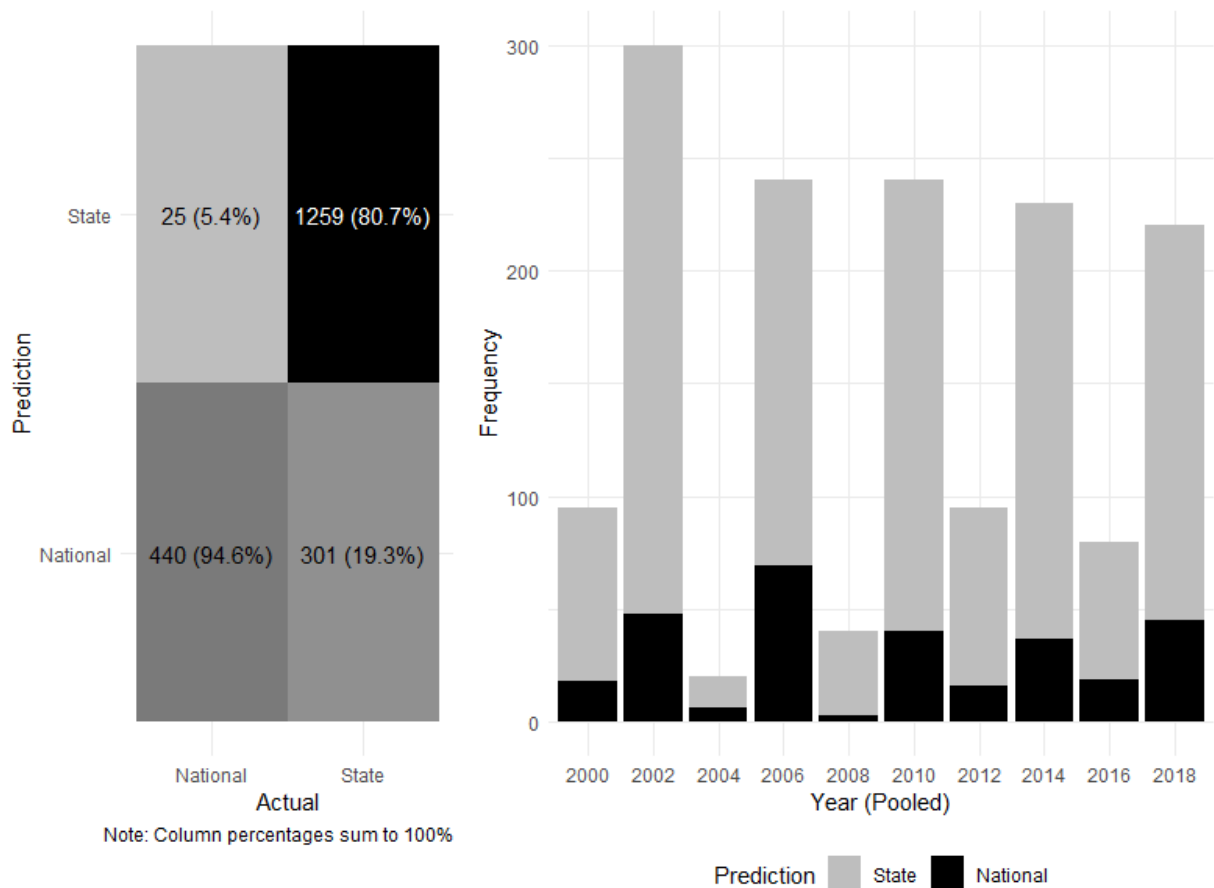


Figure 4: Lasso Confusion Matrix for Test Document Fifths and State Classification Over Time

The classification approach with lasso gives a simple binary solution to determining whether rhetoric has nationalized. However, we may want a more nuanced measure of nationalization. Debates (or chunks of debates) may be have more national or state content, but how much more? To approximate a more continuous measure of nationalization, I regress each debate topic proportion on a binary indicator of origin (presidential or gubernatorial debate) to get a simple difference of topic proportion means for each topic. Because certain topics are not consistently associated with state or national origin over repeated samples, I remove 26 topics where the coefficient p-value is greater than 0.05. I then compute as summary nationalization score as the following:

$$NatScore_i = \sum_{t=1}^{n}(TopicProp_{it} - YearProp_t) * PropDiff_t,$$

where $TopicProp_{it}$ is the topic proportion of topic $t$ in debate $i$, $YearProp_t$ is the expected topic $t$ proportion for the year of the debate, and $PropDiff_t$ is the difference in means coefficient from the aforementioned regression (a positive coefficient indicating a topic being more national, negative more state). Functionally, this takes each estimated debate topic proportion, controls for yearly topic variation, weights each topic by the difference in means coefficient, and sums over all topics. The average nationalization score for gubernatorial and presidential debates from 2000 to 2018 is shown in Figure 5, with confidence intervals from a one-sample t-test.[5] Positive values indicate higher national-level content and negative values indicate higher state-level content. Again, there does not seem to be a clear over-time trend in the level of nationalization, and state-level national content is well below national-level national content. There is an uptick in national content for gubernatorial debates in 2004 and 2006, although the magnitude of the nationalization is small.

---

[5]The large confidence intervals at 2002/3 and 2004/5 are largely due to small sample size.
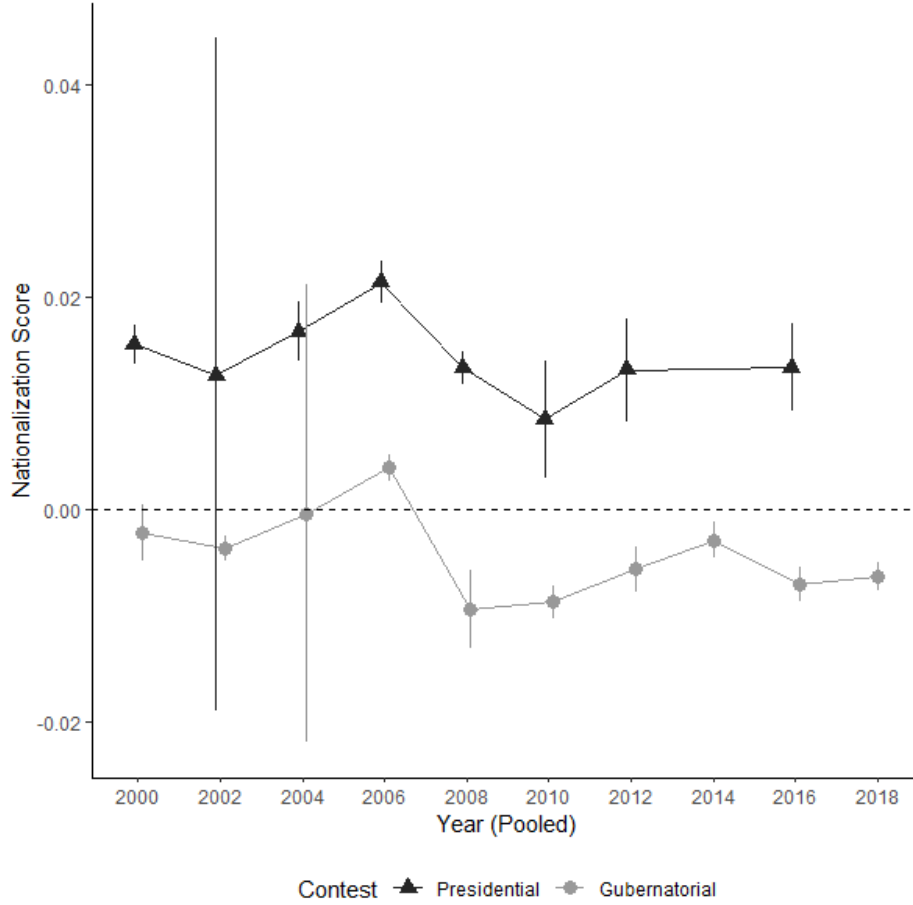
Figure 5: Nationalization in Gubernatorial/Presidential Debates, 2000-2018

Finally, given there is evidence of state-level campaign nationalization, is it related to the nationalization of election results? To answer this question, I create state-year election result nationalization scores using the method from Hopkins (2018) by regressing county-level two-party Democratic gubernatorial vote share on the most recent county-level two-party Democratic vote share and weighting by total votes cast in the presidential election. I regress this election nationalization score on the state-year campaign nationalization score from Figure 5. The resulting bivariate relationship is shown in Figure 6. The resulting bivariate relationship ($\beta = 3.04$, p = 0.08) is positive but small and uncertain.[6] A one standard deviation increase in campaign nationalization is associated with only a 0.1 standard

---

[6]This result is robust to the addition of state and year fixed effects.

deviation increase in result nationalization. Thus, campaign nationalization only accounts for a very small portion (if any) of the nationalization of election results.
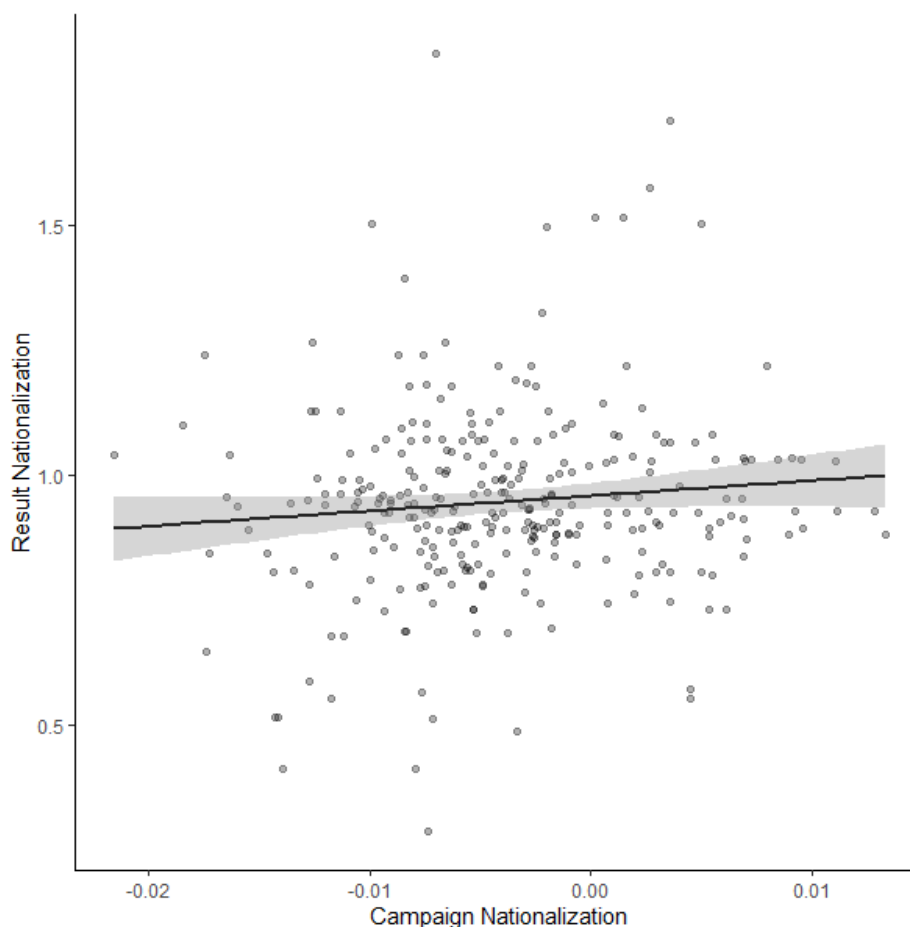


Figure 6: Campaign and Election Result Nationalization

# 5    Discussion

The preceding results cast doubt on the popular notion that gubernatorial campaigns are nationalized affairs. On the whole, gubernatorial candidates still talk about issues relevant to the state and have done so at about the same rate for the past two decades. While certain national topics do seem to find their way into gubernatorial campaigns, the share of time spent debating them is relatively small, and the degree to which such nationalized topics do

come up is only weakly associated with subsequent nationalization of election results.

What does this mean for the study of nationalization? Seemingly, while voters do still behave in a nationalized manner at the ballot box, their behavior is not influenced by short-term elite messaging. Of course, one could also describe the same phenomena in a more negative light; *despite* campaigns largely being about state-level issues, voters still behave in a nationalized manner. More work must be done to analyze the information voters have access to and how they prioritize such information (either nationalized or not) in their decision-making.

Additionally, while nationalization may account for a small share of campaign rhetoric, it may still serve a strategic purpose beyond convincing voters at the ballot box. It is possible that mentions of national issues during the campaign receive a disproportionate amount of attention from the national media, so while the candidates themselves are not nationalizing the race, the media is. This sort of coverage may draw the attention of national donor networks, whose involvement may sway the election (Gimpel, Lee, and Kaminski 2006; Hopkins 2018; Reckhow et al. 2017).

# References

Abramowitz, Alan I., and Steven Webster. 2016. "The rise of negative partisanship and the nationalization of U.S. elections in the 21st century." *Electoral Studies* 41(March): 12–22.

Arceneaux, Kevin. 2006. "The Federal Face of Voting: Are Elected Officials Held Accountable for the Functions Relevant to Their Office?" *Political Psychology* 27(October): 731–754.

Benedictis-Kessner, Justin De, and Christopher Warshaw. N.d. "Accountability for the Local Economy at All Levels of Government in United States Elections." *American Political Science Review.* Forthcoming.

Benjamin, Daniel J, and Jesse M Shapiro. 2009. "Thin-Slice Forecasts of Gubernatorial Elections." *The Review of Economics and Statistics* 91(July): 523–536.

Benoit, William L., Glenn J. Hansen, and Rebecca M. Verser. 2003. "A meta-analysis of the effects of viewing U.S. presidential debates." *Communication Monographs* 70(December): 335–350.

Benoit, William L., LeAnn M. Brazeal, and David Airne. 2007. "A Functional Analysis of Televised U.S. Senate and Gubernatorial Campaign Debates." *Argumentation and Advocacy* 44(September): 75–89.

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

Broockman, David E., and Daniel M. Butler. 2017. "The Causal Effects of Elite Position-Taking on Voter Attitudes: Field Experiments with Elite Communication." *American Journal of Political Science* 61(1): 208–221.

Brown, Adam R. 2010. "Are Governors Responsible for the State Economy? Partisanship, Blame, and Divided Federalism." *The Journal of Politics* 72(July): 605–615.

Brunner, Jim. 2016. "Gov. Jay Inslee defeats Republican challenger Bill Bryant." *The Seattle Times* (November).

Carsey, Thomas M. 2001. *Campaign dynamics: the race for governor.* Ann Arbor: University of Michigan Press.

Catanese, David. 2011. "Tomblin wins W.Va. gov race." *Politico* (October).

Chawla, N. V., K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer. 2002. "SMOTE: Synthetic Minority Over-sampling Technique." *Journal of Artificial Intelligence Research* 16(June): 321–357.

Cryer, J. 2019. Navigating Identity in Campaign Messaging: The Influence of Race & Gender on Strategy in U.S. Congressional Elections. SSRN Scholarly Paper ID 2863215 Social Science Research Network Rochester, NY: .

Das, Sanmay, Betsy Sinclar, Steven Webster, and Hao Yan. 2019. "All (Mayoral) Politics is Local?".

Gimpel, James G., Frances E. Lee, and Joshua Kaminski. 2006. "The Political Geography of Campaign Contributions in American Politics." *The Journal of Politics* 68(August): 626–639.

Hayes, Danny, and Jennifer L. Lawless. 2018. "The Decline of Local News and Its Effects: New Evidence from Longitudinal Data." *The Journal of Politics* 80(January): 332–336.

Hillygus, D. Sunshine, and Simon Jackman. 2003. "Voter Decision Making in Election 2000: Campaign Effects, Partisan Activation, and the Clinton Legacy." *American Journal of Political Science* 47(4): 583–596.

Hopkins, Daniel J. 2018. *The increasingly United States: how and why American political behavior nationalized.* Chicago studies in American politics Chicago: The University of Chicago Press.

Huddy, Leonie, and Alexa Bankert. 2017. "Political Partisanship as a Social Identity." *Oxford Research Encyclopedia of Politics* (May).

Jensen, Amalie, William Marble, Kenneth Scheve, and Matthew Slaughter. 2019. "City Limits to Partisan Polarization in the American Public." *Working Paper* (July).

Kousser, Thad. 2014. "How America's "devolution revolution" reshaped its federalism." *Revue française de science politique* 64(2): 265.

Manchester, Julia. 2019. "GOP seeks to nationalize gubernatorial elections." *The Hill* (November).

Martin, Gregory J., and Joshua McCrain. 2019. "Local News and National Politics." *American Political Science Review* 113(May): 372–384.

Martin, Jonathan. 2017. "Primary for Virginia Governor Tests Power of an Anti-Trump Campaign." *The New York Times* (February).

Martin, Jonathan. 2019. "Kentucky Governor's Race Tests Impact of Impeachment in States." *The New York Times* (October).

Mason, Lilliana. 2015. ""I Disrespectfully Agree": The Differential Effects of Partisan Sorting on Social and Issue Polarization." *American Journal of Political Science* 59(1): 128–145.

McKinley Jr., James C. 2010. "Perry Re-elected in Texas Governor Race." *The New York Times* (November).

Mullinix, Kevin J. 2016. "Partisanship and Preference Formation: Competing Motivations, Elite Polarization, and Issue Importance." *Political Behavior* 38(June): 383–411.

Parthasarathy, Ramya, Vijayendra Rao, and Nethra Palaniswamy. 2019. "Deliberative Democracy in an Unequal World: A Text-As-Data Study of South India's Village Assemblies." *American Political Science Review* 113(August): 623–640.

Reckhow, Sarah, Jeffrey R. Henig, Rebecca Jacobsen, and Jamie Alter Litt. 2017. ""Outsiders with Deep Pockets": The Nationalization of Local School Board Elections." *Urban Affairs Review* 53(September): 783–811.

Roberts, Margaret E., Brandon M. Stewart, and Edoardo M. Airoldi. 2016. "A Model of Text for Experimentation in the Social Sciences." *Journal of the American Statistical Association* 111(July): 988–1003.

Rojas, Rick, and Jeremy Alford. 2019. "In Louisiana, a Narrow Win for John Bel Edwards and a Hard Loss for Trump." *The New York Times* (November).

Sievert, Joel, and Seth C. McKee. 2019. "Nationalization in U.S. Senate and Gubernatorial Elections." *American Politics Research* 47(September): 1055–1080.

Tesler, Michael. 2015. "Priming Predispositions and Changing Policy Positions: An Account of When Mass Opinion Is Primed or Changed." *American Journal of Political Science* 59(4): 806–824.

Vavreck, Lynn. 2009. *The message matters: the economy and presidential campaigns*. Princeton, N.J: Princeton University Press.

Zhang, Kaiping. 2019. "Encountering Dissimilar Views in Deliberation: Political Knowledge, Attitude Strength, and Opinion Change." *Political Psychology* 40(2): 315–333.

# 6    Appendix

## 6.1    A1: SMOTE

SMOTE uses a k nearest neighbors approach amongst minority class instances to create synthetic minority instances; in this particular case, the data structure supplied to SMOTE has one observation per document, with variables included for year (pooled into two-year intervals), document origin (state or national), and a count variable for every word in the corpus. For example, imagine I'm creating a model to distinguish between pop songs and nursery rhymes, and one of my observations in the training set is the following nursery rhyme as an example: "Mary had a little lamb, little lamb, little lamb." The data version of this would have one factor variable column for the class (nursery rhyme) and five others, one for each distinct word. The "Mary," "had," and "a" columns would each contain the value 1 since each word occurs only once, while the "little" and "lamb" columns would each contain the value 3. You can then imagine each instance as an observation in n-dimensional space (each dimension being the count for a distinct word). SMOTE creates synthetic instances by looping over each instance, choosing one of the K closest neighbors, connecting them by a line through the n-dimensional space, and synthesizing a new instance. Because topical content is likely to vary between years in my speech data, I split the data by year, apply SMOTE separately to each data-year to solve class imbalance within that year, and recompile. Importantly, SMOTE can only create synthetic instances with words that existed in the true minority instance documents, so the process does not introduce words from the majority class.

## 6.2 A3: Regression models

Table 1:

| | Dependent variable: | |
|---|---|---|
| | estimate | |
| | (1) | (2) |
| nat.score | 3.043* | 4.241* |
| | (1.755) | (2.445) |
| | | |
| Constant | 0.958*** | 0.908*** |
| | (0.013) | (0.110) |
| | | |
| State/Year Fixed Effects? | No | Yes |
| Observations | 296 | 296 |
| $R^2$ | 0.010 | 0.432 |
| Adjusted $R^2$ | 0.007 | 0.283 |
| Residual Std. Error | 0.186 (df = 294) | 0.158 (df = 234) |

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