

Author Clustering and Topic Estimation for Short Texts

Graham Tierney*, Christopher Bail†, Alexander Volfovsky*

June 20, 2022

Abstract

Analysis of short text, such as social media posts, is extremely difficult because of their inherent brevity. In addition to classifying topics of such posts, a common downstream task is grouping the authors of these documents for subsequent analyses. We propose a novel model that expands on the Latent Dirichlet Allocation by modeling strong dependence among the words in the same document, with user-level topic distributions. We also simultaneously cluster users, removing the need for post-hoc cluster estimation and improving topic estimation by shrinking noisy user-level topic distributions towards typical values. Our method performs as well as — or better — than traditional approaches, and we demonstrate its usefulness on a dataset of tweets from United States Senators, recovering both meaningful topics and clusters that reflect partisan ideology. We also develop a novel measure of echo chambers among these politicians by characterizing insularity of topics discussed by groups of Senators and provide uncertainty quantification.

1 INTRODUCTION

The proliferation of micro-blogs and other short text on social media platforms, product reviews, and online advertisements has increased interest in models that can extract useful information from short documents. Standard text analysis tools require large numbers of long documents to estimate useful quantities. One of the most widely used models for text data is the Latent Dirichlet Allocation (LDA) topic model from [Blei et al. \(2003\)](#), which estimates topics by within-document word co-occurrences. If document word counts are extremely sparse — as is the case with short text — LDA and other topic models will suffer in topic estimation. Moreover, when texts themselves are short, researchers often want to jointly model both text and other features such as similarities among the authors of the texts.

We motivate the development in this paper by considering a dataset of all tweets by United States Senators from the first five months of 2020 and use it to measure partisanship across topics expressed in the tweets. This range covers the second impeachment of President Trump through the early stages of the coronavirus pandemic. Tweeting has become an important communication tool for US politicians who issue positions, interact with constituents, and attempt to set the political agenda on a daily basis ([Jaidka et al., 2019](#)). Just as the introduction of radio addresses in 1923 and televised debates in 1960 marked transition points in American political discourse, Twitter use has now become a focal point of US politics. Legislators respond to topics raised by their constituents on Twitter ([Barberá et al., 2019](#)), and President Trump used Twitter to distract from negative news stories ([Lewandowsky et al., 2020](#)). Clustering users in this political context can reveal the presence or absence of partisan topics, which can measure “echo chambers” and ideological segregation ([Bakshy et al., 2015](#); [Bail et al., 2018](#)). Many analyses of echo chambers on social media focus on account-level homophily, e.g. identifying that Twitter users typically follow accounts with political ideology similar to their own ([Mosleh et al., 2021](#); [Barberá, 2015](#)). This tie formation would affect exposure to different kinds of political messages only if politicians discuss distinct topics. As such, our model focuses on the content of politician messages. However, characterizing behavior and topics discussed at the tweet-level is still extremely challenging due to the brevity of the messages. For example, in our application, the post-term matrix, where the ij entry is the count of appearances of word j in post i , is 99.8% zeros.

*Statistical Science and †Sociology, Duke University

In this paper we develop a model that directly targets short text while leveraging similarities among users to improve topic coherence. Specifically, we cluster U.S. Senators by Twitter activity, and discover meaningful topics and groupings of accounts. We collected a dataset of all tweets by sitting US Senators from January 1st, 2020 to May 31st, 2020, which, after pre-processing described in Section 4, contains 61,241 tweets and 2,644 unique words. Our model captures partisan differences in coverage of the coronavirus pandemic and economic stimulus. Moreover, our clustering of Senators based only on their Twitter activity maps onto distinct ideologies that can transcend party identification as well: Senators who tweet similarly have similar political beliefs. Those Senators who are grouped with out-partisans are frequently outliers within their own party, such as the moderate Republican Senators Susan Collins and Lisa Murkowski and the Democratic Socialist Senator Bernie Sanders. We also quantify insularity of topics discussed by the identified clusters by measuring how many tweets it would take each cluster to mention every topic. We find that the more liberal Democrats and nationally-focused Republicans have the broadest topic coverage on average, but individual Democrats are much less varied than individual Republicans in these categories, which means that individual Democrats tend to cover every topic faster. In contrast to two-stage procedures that first learn topics then cluster users, our estimated topics are more coherent and our clusters better match established theoretical measures of Senator ideology (Poole and Rosenthal, 2017).

The most common method researchers have used to overcome the brevity of messages on social media platforms is to merge all short documents by the same user (or user type) into a single, long document and estimate a topic model on those combined documents (Hong and Davison, 2010; Pennacchiotti and Popescu, 2011; Steinskog et al., 2017; Barberá et al., 2019). More recent theoretical work has estimated which tweets to merge from the data by borrowing methods from the information retrieval literature (Hajjem and Latiri, 2017). Central to these methods is the assumption that no information is contained in the fact that certain words are *chosen* to belong to the same short document. A Twitter user facing significant character-limits choosing to write two words in conjunction should imply that their topics are more similar than the same two words used in distinct tweets.

Given the downsides of combining short texts into a single long text per user, Zhao et al. (2011) first introduced a single topic per-short-text LDA (called Twitter-LDA) to compare the content of New York Times articles to Twitter. Each user tweeted about a mixture of topics, and each tweet had only a single latent topic.¹ They found that using a single topic per tweet — rather than a mixture of topics — produced more semantically coherent topic distributions. This is very intuitive: it is highly likely that only a single message can be conveyed in under 280 characters. The Twitter-LDA model proposed by Zhao et al. (2011) is nearly identical to the Dirichlet-Multinomial Mixture (DMM) model in Yin and Wang (2014). In the latter model, each user produces a single document, which is mapped to a single topic distribution over words. This is equivalent to single-topic LDA where each document comes from a unique user. DMM methods traditionally ignore user information, and model a single corpus-level mixture over topics, rather than a user-level mixture. Both DMM and Twitter-LDA were derived as collapsed Gibbs samplers, making them impractical for even moderate sample sizes. To overcome the issue that each text is produced by a unique user in DMM, the recently proposed LapDMM, which also derives a collapsed variational algorithm, introduces a regularization term that makes semantically similar texts, identified by word-embedding methods discussed later, more likely to have the same topic (Li et al., 2019). This procedure requires knowing or learning this semantic distance *a priori* which again suffers from the same constraints that all short text modeling suffers from.

Other data augmentation approaches have recently combined learned topic information from longer documents with shorter documents in terms of keywords (Sahami and Heilman, 2006) or actual text (Hu et al., 2009; Jin et al., 2011). A significant limitation of these methods is that they require a corpus of long text that is *a priori* compatible with the short-text documents being studied. Yan et al. (2013) proposed a biterm model (BTM) for short text where topic-level word co-occurrence relationships are learned for the whole corpus. Each document is decomposed into the set of all two-word pairs (biterns) that can be formed from the document’s contents. Each bitern is modeled as being drawn from a single topic. The lack of a meaningful generative model makes the results less interpretable and the method cannot leverage information that the same user generated multiple documents. Moving beyond topic models altogether, neural networks

¹Incorporating authorship information into topic models is not new; the author-topic model of Rosen-Zvi et al. (2004) associated each author in a corpus with a distribution over topics that was identified by multi-author documents. Multi-author documents, however, are typically not found in social media.

have been introduced to estimate word embeddings in a latent vector space, then classify short texts based on a distance metric between texts (Mikolov et al., 2013; Rangarajan Sridhar, 2015). The word mover’s distance (Kusner et al., 2015) is one popular method to compute a distance matrix for a corpus. More recent work has leveraged word embedding methods to augment topic models, either with embeddings learned from a large external corpus (Li et al., 2016) or the corpus of short documents being modeled (Shi et al., 2018). These methods generally ignore user information (the distance metric does not account for authorship), and, more importantly, rely on pre-trained word embeddings from larger corpora, which may not reflect the specific vocabulary and semantic usage of words in the corpus in question. This reliance especially limits the ability of these methods to capture new or emerging topics.

Thematically similar to LapDMM, clustering techniques have been introduced into the LDA universe to improve the performance of topic modeling. Barnard et al. (2003) proposed a mixture of LDA models to group documents with an associated image. Within a cluster, each document’s topic distribution is a draw from a cluster-specific Dirichlet distribution. Wallach (2008) proposed a similar clustering model that placed a hierarchical prior on the cluster-specific Dirichlet parameters. Xie and Xing (2013) expanded that level of clustering to model global topics shared across all documents and local topics used only within a cluster. Just as with LapDMM, the clustering here is defined at the document level and not at the user level.

1.1 Our contribution

To improve upon these methods for short text we combine a generalization of the short text LDA topic model with unsupervised clustering of authors of short documents. *This hierarchical model is able to share information at multiple levels leading to higher quality estimates of per-author topic distributions, per-cluster topic distribution centers, and author cluster assignments.* Our method fuses the clustering of **both users and documents**. We leverage the dependencies of words in a single post and borrow strength across users to inform user-specific topic distributions. Our method outperforms standard tools on simulated data even when our model is misspecified, and provides unique insights in a real-world application. To facilitate computation we derive both a collapsed Gibbs sampler that directly targets the posterior as well as a variational approximation to that posterior. We note that prior work on short text has heavily relied on slow and unscalable Gibbs samplers (Zhao et al., 2011; Yan et al., 2013) and that the LapDMM variational approximation heavily relies on a one document per-user assumption. Our method can handle multiple short texts per user and generalizes the variational approximation derived for traditional LDA such that it can be naively parallelized. We show that this approximation improves, in terms of inference being comparable to that of the Gibbs sampler, as the scale of the data increases. This computational gain extends the model to otherwise intractable datasets.

The remainder of the paper is organized as follows: Section 2 outlines the details of the proposed short text LDA with clustering (stLDA-C) model and shows its’ relationship to previously proposed models. Section 3 derives the collapsed Gibbs sampler and the variational approximation to the posterior. We apply the model in Section 4 to the dataset of US Senator tweets, identifying distinct clusters of politicians and heterogeneity in topics by partisanship. We leverage the form of our model to define a novel metric of the echo-chamber effect on social media and discuss the potential heterogeneity of this effect. Section 5 validates the proposed model and estimation procedures on simulated datasets, both when the model is correctly specified and when it is misspecified.

2 MODEL

The general framework for latent Dirichlet allocation models specifies a probabilistic relationship between a collection of topics, a corpus of words, and the documents that are comprised of those words. Traditional LDA (Blei et al. (2003), referenced as tLDA below) models each document as a mixture over topics: Each word has its own latent topic assignment, and documents have unique topic distributions. The tLDA model conceptualizes topics by considering words co-occurring in the same document, but when the documents are short those co-occurrences are extremely rare. This baseline approach further ignores the information about whether documents were produced by the same author, which is particularly common for short text from social media platforms. Our proposed model captures this user-level dependence and improves topic estimation by leveraging the dependence of the latent topic variables of words in the same tweet.

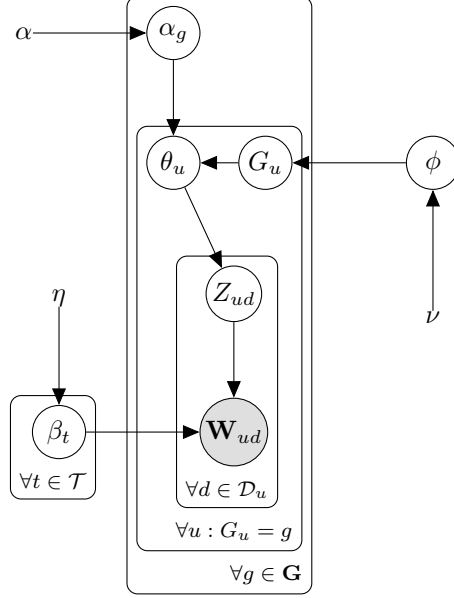


Figure 1: Plate diagram for the stLDA with clustering model. The diagram further illustrates the connection between the model proposed in this manuscript with previous models in the literature. Observed quantities are illustrated in gray and objects without a border around them are hyperparameters.

Concretely, our method adds *three* elements to the traditional LDA model. To model short texts: (1) we restrict each document to a single topic, rather than a mixture over topics. (2) We allow mixing of topics at the user-level, where observed author information allows us to learn more accurate topics by borrowing strength across documents authored by the same user. Lastly, (3) we incorporate unsupervised clustering to group users by similarity in their topic choices. This allows for information sharing between users with similar topic distributions and shrinkage of user-topic distributions with high uncertainty towards a common mean. We refer to our model as stLDA-C.

Figure 1 displays the plate diagram of our hierarchical model. Prior distributions over cluster-level topics, user cluster assignments, and topic distributions over words are governed by fixed parameters α , ν , and η respectively. α_g is the vector parameter of a Dirichlet distribution over topics choices for users in cluster g , identified by $G_u = g$. Each user-specific topic distribution θ_u is a draw from $\text{Dir}(\alpha_g)$. For each tweet d by user u , its topic, Z_{ud} , is a single draw from θ_u , and all words in tweet ud are sampled from the topic distribution over words, β_t , where $Z_{ud} = t$. ϕ encodes the proportion of users in each group and forms a prior distribution for G_u . The generative process is described below:

1. For each cluster $g = 1, \dots, G$, sample $\alpha_g \stackrel{\text{iid}}{\sim} f_\alpha$, where f_α is the prior over α_g
2. For each topic $t = 1, \dots, T$, sample $\beta_t \stackrel{\text{iid}}{\sim} \text{Dir}(\eta)$
3. Sample $\phi \sim \text{Dir}(\nu)$
4. For each user u :
 - (a) Sample $G_u \sim \text{Cat}(\phi)$: $P(G_u = g) = \phi_g$
 - (b) Sample $\theta_u \sim \text{Dir}(\alpha_{G_u})$
 - (c) For each document $d = 1, \dots, n_u$
 - i. Sample the topic $Z_{ud} \sim \text{Cat}(\theta_u)$
 - ii. Sample words $\mathbf{W}_{ud} \sim \text{Multi}(n_{ud}, \beta_{Z_{ud}})$

$\text{Multi}(n, p)$ is the multinomial distribution of size n and probability vector p . $\text{Cat}(p)$ is the categorical distribution, a multinomial of size one. $\text{Dir}(\zeta)$ is the Dirichlet distribution with parameter ζ . Below we discuss how our model generalizes aspects of other LDA models and specializes them to short text.

Introducing a user clustering layer. On Twitter, researchers frequently want to not only learn topics

of tweets, but also group users into categories. Previous work has estimated user-level topic frequencies, and then used those estimates in a separate clustering procedure. Our method links the topic modeling and clustering task, allowing full propagation of uncertainty and combining the two inference steps. This clustering layer is reminiscent of the document clustering in tLDA due to Wallach (2008), but because there are many documents per user, the clustering is more naturally described on the user layer. Each user u 's cluster is denoted by G_u . For users in the same cluster g , $\theta_u \sim \text{Dir}(\alpha_g)$. This hierarchical model allows for shrinkage of noisily estimated θ_u parameters from users who post infrequently to a common mean and borrowing strength for that estimation across users with similar topics. This is especially important when there is large heterogeneity in the number of posts from each user, which we demonstrate in the application.

Using multiple clusters, rather than a single hierarchical prior, is important in the context of social media. The scale and structure of social media data often means that many users discuss completely distinct sets of topics. Literature on echo chambers on social media platforms document this phenomenon (Bakshy et al., 2015). Having a single hierarchical prior, where all θ_u are drawn from the same distribution, will shrink estimates of θ_u to an average taken across the entire population of users. With multiple clusters, noisily estimated θ_u parameters are shrunk towards the average of users they are most similar to.

Our model treats a corpora as a mixture of a finite number of clusters and finite number of topics, both of which are assumptions that have been relaxed in the Bayesian non-parametric literature using Dirichlet processes to allow for countably infinite numbers of mixture components. The multi-corpora Hierarchical Dirichlet Process model of Teh et al. (2006), for example, can be thought of as an infinite dimensional extension of our model when cluster labels (corpora in their terminology) are fixed and known.² As described below, the scenarios and data we designed our model to analyze require a smaller number of interpretable topics, so allowing for large numbers of topics is not necessarily an improvement. Moreover, when the subjects of individual tweets are of particular interest and domain experts can assess meaningfulness of topics, as in our application, human validation of topic quality and utility for downstream tasks is often preferable to the data-driven selection in the Hierarchical Dirichlet Process model.

Connections to single topic LDA. Our model naturally generalizes the two different single topic per document models, Dirichlet-Multinomial Mixture (DMM) and single topic LDA (stLDA). In the DMM model, each short text is drawn from a single latent topic and the corpus has a single set of mixture proportions (Nigam et al., 2000). In this model, each document belongs to a single latent topic and words in the document are a multinomial draw from a cluster-specific probability distribution over words. Several estimation methods have been developed, including EM and collapsed Gibbs sampling (Nigam et al., 2000; Yin and Wang, 2014). DMM is a simpler version of the model proposed here (and by Zhao et al. (2011)). Indeed, the DMM model can be recovered by removing the user- and cluster-specific plates in Figure 1.

stLDA from Zhao et al. (2011) overcomes the sparsity of word co-occurrences by ensuring that all words in a document contribute to the estimation of the word distribution of a single topic, t . By placing the topic mixing at the user rather than document level, stLDA captures dependence among documents by the same user and improves topic estimation by borrowing information across these documents. Topics of documents by the same user are conditionally independent given the user's topic distribution θ_u , but become dependent after integrating θ_u out, whereas in tLDA the topics of individual documents remain independent when each document's topic mixture is integrated out. To learn the topic of a given tweet by user u , stLDA places significantly more weight on words and documents produced by u because they inform beliefs about θ_u . This is in contrast to tLDA in which all words are given equal importance, regardless of the user who produced them. stLDA can be recovered from our model by integrating over cluster parameters α_g , G_u , and ϕ in Figure 1.

Connections to combined LDA. The way many researchers apply topic models to Twitter data is to combine all documents by the same user into one large document and estimate tLDA models on the smaller set of longer documents (Hong and Davison, 2010; Pennacchiotti and Popescu, 2011; Steinskog et al., 2017; Barberá et al., 2019). We call this method cLDA. The topics of words in the same tweet should be related, even given a user's average topic frequency, because tweets are often focused on a narrow subject

²A large literature on discrete, supervised LDA has similar structural form to our model here when class or cluster labels are observed, such as the recent work Wang et al. (2021). These models, while similar in form, are not directly comparable to our model because we must infer the clustering from the data. Dirichlet process models that use latent clustering, as in the Nested Dirichlet Process (Rodriguez et al., 2008), do not allow sharing of atoms across clusters. In our case, that would mean all topics are cluster-specific, which is not desirable because we want to compare topic choice across clusters.

due to character limits. However, because cLDA uses a simple bag-of-words model, the information that certain words were used in different tweets is lost. In cLDA, the topics of words in the same tweet by one user are conditionally independent given the user’s topic frequency. Suppose words w and w' appear in the same tweet by user u and they have potentially distinct topics Z_w and $Z_{w'}$. Conditional on θ_u , the user’s topic distribution, knowing the topic of w' provides no additional information about the topic of w . $P(Z_w = t | Z_{w'} = t, \theta_u) = P(Z_w = t | \theta_u) = \theta_{ut}$, the marginal probability that user u picks topic t . This is true for all words authored by a single user, regardless of whether they appeared in the same tweet. This result does not capture how most users create posts. Suppose a politician talks about national politics half of the time and local issues the other half. Knowing that one of the words in a tweet is about national politics should increase the belief that other words in the tweet are also about national politics, even knowing the user spends on average half the time on each issue. cLDA resolves the issue of short documents by assuming a level of conditional independence that fails to capture the document generating process.

In stLDA, the topics of words in the same tweet are identical, which more accurately captures how tweets are created. This means that, in stLDA, knowing the topic of one word in a tweet identifies the topics of all other words. In the same situation as above, where words w and w' appear in the same tweet by user u , $Z_w = Z_{w'} = Z_{ud}$. Thus, conditional on $Z_{w'}$, Z_w is independent of θ_u . $P(Z_w = t | Z_{w'} = t, \theta_u) = 1$. *One can thus think of cLDA and stLDA as two extremes for modeling the dependence of topics of words in the same tweet, conditional on the user’s average topic selection.* stLDA still allows for borrowing strength across the topics of the user’s other tweets. θ_u is learned from those topics, which strongly influences the topic probabilities of new tweets.

DMM models have also been augmented with word embeddings to improve topic estimation by identifying semantically similar words (Li et al., 2016). Again, these models rely on either a pre-trained set of embeddings or an auxiliary dataset known to overlap, at least partially, in subject matter. We feel that these are quite strong assumptions, particularly given the target application of our model is social media data where semantic content is different from longer text. A common application of short text models is to news headlines, which have a more plausible link to longer text documents, but our model is designed for a different area of application.

3 PARAMETER ESTIMATION

We develop two estimation procedures: first, an exact but (potentially) slow collapsed Gibbs sampler and, second, a faster but approximate variational inference method.

3.1 Collapsed Gibbs

We estimate the latent topic parameters Z_{ud} with collapsed Gibbs sampling. Griffiths and Steyvers (2004) developed the collapsed Gibbs sampler for tLDA by analytically integrating out θ_d and β_t , so only samples of the latent topic variables were required. Below we derive the collapsed sampler for stLDA-C. We are interested in $P(Z_{ud} = t | \mathbf{Z}_{-ud}, \mathbf{W}, G_u = g, \alpha_g)$, where \mathbf{Z}_{-ud} denotes the topics of all tweets except for tweet d by user u , \mathbf{W} denotes the word counts for all tweets, and G_u and α_g contain the cluster assignment and parameters for u . Conditional on θ_u and β_t , there is no dependence on topics or words in other tweets. Hence we can write $P(Z_{ud} = t | \theta_u, \beta_t, \mathbf{Z}_{-ud}, \mathbf{W}) \propto P(\mathbf{W}_{ud} | Z_{ud} = t, \beta_t) P(Z_{ud} = t | \theta_u) \propto \prod_i \beta_{ti}^{W_{udi}} \theta_{ut}$, where the product is taken over all $i = 1, \dots, V$ unique words in the corpus. This density can be integrated over the posteriors for β_t and θ_u given the words and topics of other tweets and the α_g hyper-parameters. With independent Dirichlet priors on all β_t and hierarchical Dirichlet priors on θ_u , the posteriors are also Dirichlet by conjugacy of the Dirichlet and Multinomial distributions. As such we have $\theta_u | \mathbf{Z}_{-ud} \sim \text{Dir}(\alpha_g + \mathbf{Z}_{-d}^{(u)})$, where $\mathbf{Z}_{-d}^{(u)}$ is the vector of topic counts in only user u ’s tweets excluding tweet d and $\beta_t | \mathbf{Z}_{-ud}, \mathbf{W} \sim \text{Dir}(\eta + \mathbf{W}_{-ud}^{(t)})$ where η is the parameter of the prior and $\mathbf{W}_{-ud}^{(t)}$ is the word counts of tweets with topic t excluding tweet ud . Putting these together we get:

$$p(Z_{ud} = t | \mathbf{Z}_{-ud}, \mathbf{W}) = \int \theta_{ut} p(\theta_u | \mathbf{Z}_{-ud}) d\theta_u \int \prod_i \beta_{ti}^{W_{udi}} p(\beta_t | \mathbf{Z}_{-ud}, \mathbf{W}_{-ud}) d\beta_t$$

These integrals are equivalent to the posterior predictive probabilities of observing $Z_{ud} = t$ and $\mathbf{W}_{\mathbf{ud}}$ given the topics and words of other tweets. The integral over θ_u is the same as in [Griffiths and Steyvers \(2004\)](#) and it simplifies to:

$$\frac{\alpha_{gt} + (Z_{-d}^{(u)})_t}{\sum_j \alpha_{gj} + (Z_{-d}^{(u)})_j}$$

The second integral simplifies to:

$$\frac{\Gamma\left(N\eta + \sum_i (W_{-ud}^{(t)})_i\right)}{\prod_i \Gamma\left(\eta + (W_{-ud}^{(t)})_i\right)} \frac{\prod_i \Gamma\left(\eta + (W_{-ud}^{(t)})_i + W_{udi}\right)}{\Gamma\left(N\eta + \sum_i (W_{-ud}^{(t)})_i + \sum_i W_{udi}\right)}$$

For more detailed proofs of this derivation see Appendix A. These two quantities are calculated for each topic, then normalized to sum to 1 across all topics.

Corpus-level cluster proportions ϕ are estimated from conditionally conjugate prior distributions. Given the counts of cluster membership \mathbf{G} and a $\text{Dir}(\nu)$ prior on ϕ , $\phi|\mathbf{G} \sim \text{Dir}(\nu + \mathbf{G})$. ϕ provides a prior distribution for each G_u , the cluster indicator for user u .

We note that $p(G_u = g|\phi, \mathbf{Z}^{(\mathbf{u})}, \alpha_g) \propto P(\mathbf{Z}^{(\mathbf{u})}|\alpha_g, G_u = g)p(G_u = g)$. The first term is simply the Dirichlet-Multinomial probability of the given topic counts, and the second term is ϕ_g . It is worth noting that the ϕ_g term effectively builds into the model a preference for fewer clusters. The posterior for ϕ_g will be roughly proportional to the number of other users in cluster g , which means that when sampling a new cluster for user u , the model will prefer to place u into a larger cluster. This results in a property where some clusters will be empty if the data are better modeled by fewer clusters than the number the researcher has specified. We will demonstrate this property in our simulations below.

We sample α_g with a Metropolis update within the collapsed Gibbs sampler. Any prior distribution can be chosen for α_g , but it is easiest to think of α_g as $\mathbf{m}_g c_g$ where \mathbf{m}_g is a point on the simplex in \mathbb{R}^T and $c_g > 0$ is a concentration parameter. A Dirichlet prior on \mathbf{m}_g can capture the prior expected value of θ_u and the prior on c_g determines how concentrated the prior is around that expected value. An important consideration in estimation is if the concentration for group g is small, very different user topic distributions have similar likelihoods in that group. Informative priors on this parameter can be useful to ensure that within-group topic variability is low. Alternatively, if a diffuse prior is used, the model may estimate that certain clusters have zero users in them while other clusters have highly variable topic distributions. This feature can be useful to determine if the data are more effectively described as coming from fewer clusters.

The most computationally demanding part of the MCMC is the collapsed sampler for Z_{ud} . The topic of each tweet depends on the topic assignment of all other tweets, even for tweets by other users, because those topics inform beliefs about β_t . Thus, this step does not parallelize across users or groups. Consequently, it is useful to compute many updates to the other parameters, α_g , ϕ , and G_u , after each iteration through Z_{ud} to improve mixing.³

3.2 Variational Inference

Variational inference is commonly used to estimate LDA topic models ([Blei et al., 2003](#)). The key difficulty in estimating the full posterior for stLDA-C, as with traditional LDA, is the link between β_t and θ_u . The dependence between these two parameters, clearly observed in the above Gibbs sampler, can lead to poor mixing and a computational bottleneck. The link is why the collapsed Gibbs sampler needs to update each post’s topic sequentially rather than in parallel. In variational inference, the true posterior is approximated by a simpler distribution, typically with less dependence among the parameters. We follow the same approach and approximate the true posterior by minimizing the KL divergence between the variational distribution and the true posterior.

We first describe the variational approximation for a single user, the analog of a single document in tLDA, then describe how that result generalizes. We wish to identify three latent, user-specific variables given higher-order parameters and the data: $p(g_u, \theta_u, \mathbf{Z}_u|\phi, \beta, \alpha, \mathbf{W}_u)$. g_u is user u ’s cluster membership, θ_u

³Updating certain parameters more frequently does not change the stationary distribution of the Markov chain. It simply helps find the local optima of the re-sampled parameters given the topics.

is the user’s topic choice probabilities, and \mathbf{Z}_u contains the topics of the user’s tweets. ϕ is the corpus-level cluster proportions, β is the $T \times V$ matrix of T topic distributions over the V words in the corpus, α is the $G \times T$ matrix of G cluster-specific, T -dimensional Dirichlet parameters. \mathbf{W}_u is the $n_u \times V$ document-term matrix for user u . For clarity the u subscript will be omitted for the remainder of this section.

We approximate the distribution p above with the variational distribution q :

$$q(g, \theta, \mathbf{Z}|\lambda, \gamma, \xi) = q(g|\lambda)q(\theta|\gamma) \prod_{d=1}^n q(Z_d|\xi_d)$$

where $q(g|\lambda)$ is categorical over the different clusters with probabilities λ , $q(\theta|\gamma)$ is Dirichlet with parameter γ , and $q(Z_d|\xi_d)$ is categorical over the different topics with probabilities ξ_d . Note that these are also user-specific quantities. In this section, the parameters of the user-specific variational distributions are referred to as the variational parameters, and the parameters α , β , and ϕ that are shared across users are referred to as model parameters. The variational inference approach (summarized by Figure 2) iterates between learning user specific parameters (in parallel over users) and updating the overall model parameters.

1. Initialize model parameters ϕ , α , and β .
2. Initialize variational parameters λ_u , γ_u , and ξ_u .
3. **repeat**
4. **for** each user u :
5. Update variational parameters with Equations (1) - (3)
6. Update model parameters with equations (4) - (5) and Appendix B results.
7. **until** convergence

Figure 2: Full Variational Algorithm for stLDA-C

3.2.1 Variational Parameter Estimation

We want to set the variational parameters λ , γ , and ξ to minimize the KL divergence from the variational distribution to the posterior. Using $D(q||p)$ to denote this quantity, the optimization problem can be expressed as:

$$(\lambda^*, \gamma^*, \xi^*) = \arg \min_{\lambda, \gamma, \xi} D(q(g, \theta, \mathbf{Z}|\lambda, \gamma, \xi) || p(g, \theta, \mathbf{Z}|\phi, \beta, \alpha, \mathbf{W}))$$

Detailed in Appendix B, the minimization can be computed with an iterative fix point method. The update rules are shown below.

$$\lambda_g \propto \phi_g \exp(E_q[\log p(\theta|g, \alpha_g)]) \tag{1}$$

$$\gamma_t = \sum_g \lambda_g \alpha_{gt} + \sum_d \xi_{dt} \tag{2}$$

$$\xi_d \propto \exp(E_q[\log \theta_t]) \prod_{j=1}^V \beta_{tj}^{W_{dj}} \tag{3}$$

Next, we explain each equation and the interpretation of the terms where possible. The expectations of functions of θ are tractable and computed analytically. Equation (1) is the variational probability of cluster membership in cluster g . The update rule is interpretable as the prior ϕ_g times the likelihood of the user’s topic distribution under cluster g , with the likelihood approximated by the exponential of the expected log-likelihood under the variational distribution. Equation (2) is the variational approximation to the user’s topic distribution. The Dirichlet parameter for topic t is the weighted average of each cluster’s Dirichlet

parameter (α_{gt}), with weights given by the variational probability of cluster membership (λ_g), plus the total variational probability of posts by the user having topic t . This result is analogous to the variational update in Blei et al. (2003) Equation (7) with the single model parameter α_i replaced with the cluster-weighted average. Equation (3) is the variational probability that post d by the user is from topic t . The update rule is interpretable as proportional to the prior θ_t times the likelihood of the post coming from topic t , with θ_t approximated by the variational distribution. This is again analogous to the variational update in Blei et al. (2003) Equation (6), except because all words share the same topic, the post-likelihood term is multiplicative.

A key result here is that because of the independence in the variational distribution, each user’s update only depends on model and user-specific variational parameters. Thus, the operations can be computed in parallel if such computing resources are available. This is a significant computational improvement over the collapsed Gibbs sampler where the update of \mathbf{Z} , each post’s topic, has to be computed in sequence. This, along with only needing a single value of the variational parameters, rather than a sample from the posterior, are what make this approximation method faster.

3.2.2 Model Parameter Estimation

We can estimate the remaining model parameters with an empirical Bayes approach. The marginal likelihood $p(\mathbf{W}|\phi, \alpha, \beta)$ is intractable, but the variational posterior provides a lower-bound on the likelihood that we can use instead. Essentially, the general procedure is the same variational EM method as in variational estimation of traditional LDA. First, use the above results to find $\{\lambda^*, \gamma^*, \xi^*\}$ for each user. Second, compute the values of $\{\phi, \beta, \alpha\}$ to maximize the lower bound on the likelihood computed in the first step. Iterate between these steps until convergence.

The derivations of these results are shown in Appendix B. Because these parameters are shared across users, we re-introduce the u subscripts indicating user and \mathcal{D}_u indicating the set of documents produced by user u .

$$\phi_g \propto \sum_u \lambda_{ug} \tag{4}$$

$$\beta_{tj} \propto \sum_u \sum_{d \in \mathcal{D}_u} \xi_{udt} W_{udj} \tag{5}$$

The update for α is more complicated and does not have a clear intuition, so we omit the statement here. It is available in Appendix B and is computed via an efficient Newton-Raphson method very similar to traditional LDA.

Equation (4) is clearly interpretable as setting ϕ_g proportional to the total variational probability of users belonging to cluster g . Equation (5) for β_{tj} is similarly the total variational probability that word j is assigned to topic t across all users and documents.

A particular concern with one-topic-per-word models is when attempting to model a new document, new words cannot be assigned a topic because $\beta_{tj'} = 0$ where j' is not in the training data. To address this issue, Blei et al. (2003) extends the variational approximation to cover the topic distributions over words. They perform adaptive smoothing where the update rule is $\beta_{tj} \propto \eta_t + f(t, j)$, essentially adding a topic-specific constant η_t . This additional approximation is not needed for the one-topic-per-document model proposed here. There is no need to assign a unique topic to word j' . The topic of the new document can be inferred from the words that appear both in the new document and in the training data.

Note that when initializing model parameters, α and β , uniform starts are not advisable. If rows of α are identical, then each user is equally likely to be placed in each cluster (the clusters are all identical), and the resulting λ_u values will all be uniform. Similar results for ξ_u will occur if β is initialized as uniform.

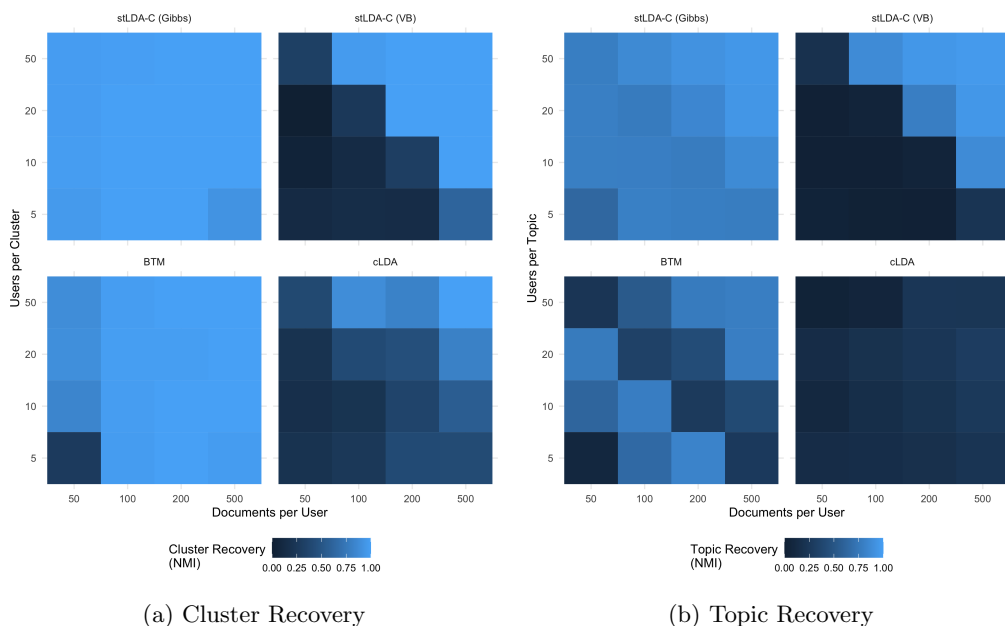
3.3 Comparison

In this section, we demonstrate the quality of the variational approximation to the true posterior, with additional comparisons to a very flexible topic model, the Biterm Topic Model (BTM, Yan et al. (2013)), and the model commonly used in practice, what we refer to as combined LDA (cLDA). We find that the

variational approximation improves as the size of the data increases, which is also when the computational advantage of the variational approximation is most useful.

Figure 3 plots cluster and topic recovery, showing that the variational approximation is quite close to the Gibbs estimation given sufficient data. We use stLDA-C as a generative model and simulate users from 3 clusters and 10 topics. Two of the clusters discuss disjoint sets of 5 topics, and the third cluster discusses all topics. This simulation is quite favorable to stLDA-C, with well-separated clusters and the single-topic per document assumption holds. We know the Gibbs sampler will quickly converge to the true posterior, which enables us to evaluate the quality of the variational approximation. The top-right panel in each sub-figure of Figure 3 shows that once the data are of sufficient size, the variationally-approximated and Gibbs-estimated posteriors match and recover the ground truth. Our recommendation is to only use the variational method when the amount of documents is prohibitively large, which is the exact scenario when the variational approximation is of high quality. BTM and cLDA are able to separate users by cluster fairly well, but still struggle with topic estimation.⁴ More extensive model validation, with more comparison methods and simulations where stLDA-C is misspecified, are provided in Section 5.

Figure 3: Gibbs and Variational Bayes (VB) Comparison



Notes: Data was simulated from the stLDA-C model with three clusters, two that post about disjoint sets of five topics and one that covers all ten topics. The number of posts per user and number of users per cluster were varied for each simulation. Documents and users are classified into the cluster and topic with the highest posterior probability. Normalized Mutual Information is used to evaluate each classification. Gibbs estimation of stLDA-C fully recovers the truth. Variational estimation will recover the truth, matching the correct posterior, with at least 5,000 total documents. Results are averaged over 7 replications for all but the largest setting, which was replicated only twice due to computation time for the Gibbs samplers.

3.4 Analyst set variables

Two key parameters are chosen by the analyst, the number of topics, T , and the number of clusters G . In traditional topic models, T is often chosen to be large, with the understanding that some topics will pick up on background noise and likely be uninterpretable. When each document covers multiple topics, this approach does not hinder inference because the high-probability topics are the ones that are most important and the less-interpretable ones tend to be lower probability across documents. In a topic-per-document model, this behavior is less desirable because the model may map an entire document to a non-interpretable

⁴Clustering for these models is performed after estimation and is described in detail in Section 5

or incoherent topic. We generally recommend choosing fewer topics than in traditional topic models, and potentially lowering the number of topics if many documents are mapped to non-interpretable topics.

G is also an important variable, and closely tied to ν , the prior on the corpus-level proportions of each cluster. A large ν encodes prior belief that the number of users in each cluster will be similar, and small ν encodes prior belief that there will be one large cluster and several smaller clusters. Depending on the analyst’s intention for downstream analysis using the clusters, this flexibility can be important. A key feature here as well is that the model can identify that fewer than the specified number of clusters can fit the data well. In simulations with very little heterogeneity across users, our model correctly identifies that users come from only one cluster and unique set of α_g parameters. As such, we recommend that analysts using this model think about setting ν specific to their prior beliefs and desired uses of cluster information, as well as potentially setting G larger than desired to identify if clusters are emptied out. Various hold-out methods can be used to evaluate different choices of T and G , where model parameters are learned on a training set of documents, then a measure of model fit is computed on a set of held-out documents. We demonstrate such methods in Sections 4 and 5.

Identifiability in topic models, as in all latent factor models, is a concern. A rotation of the topics produces equivalent likelihoods, which could hinder inference and Gibbs sampling in particular. Literature on the identifiability in topic models focuses on methods that use separable non-negative matrix factorization of the document-term matrix or some transformation of it (Anandkumar et al., 2013; Huang et al., 2016). One solution is to specify anchor words that appear in one and only one topic (Donoho and Stodden, 2004). When this assumption is not appropriate, other assumptions on higher-order moments of the corpus or topic distributions can suffice (Anandkumar et al., 2012, 2013; Huang et al., 2016).

Specifying anchor words requires strong domain knowledge, and they may be difficult to learn when the number of topics is large. These issues are less concerning in this single topic per (short) document model. One could specify anchor-documents with less domain knowledge because more context is provided; one only needs to identify a sets of documents that the analyst would like to separate into different topics. Per the above discussion, generally our method is better with fewer topics than in topic-per-word models, so this task might not be too onerous for the analyst. We do not rely on anchoring topics in our simulations or application, but the R code provided has functionality to pre-specify topics of a subset of documents. We use the Gibbs sampler only in Sections 5 and 3.3, and we do not observe evidence of label switching in the posterior draws. Topic or cluster rotations would not affect the variational approximation used in Section 4.

4 APPLICATION

We collected a dataset of tweets from sitting U.S. Senators from January 1st, 2020 to May 31st, 2020. To form an analyzable dataset, some pre-processing steps are taken: all words are made lower case, then stop words, hyperlinks, and words used in fewer than 0.1% of tweets are removed, and finally all documents with one or fewer words remaining are dropped. The final dataset includes 61,241 tweets covering 2,644 words. We set T , the number of topics, to 30 and G , the number of clusters, to 4 and proceed with estimation using the variational approximation derived in Section 3.2. Four clusters were chosen to allow for separation of the two parties and some within-party heterogeneity. We also compared topic coherence on a set of held-out tweets for models with 2, 4, 6, 8, and 10 clusters.⁵ The model with 4 clusters had the best score. Topic coherence has been shown to correlate with human-judged semantic coherence and is a common evaluation metric for topic models (Mimno et al., 2011). Moreover, pure held-out likelihood measures are known to have *negative* correlation with human-judged coherence (Chang et al., 2009). The four cluster model also had the best clustering NMI relative to party labels, reflecting the best separation of senators by partisanship.

stLDA-C estimates meaningful topics: tweets from the same topic relate to similar legislative priorities and political events, with distinct partisan- and non-partisan topics. The clusters are also meaningful, mapping onto senator ideology measured by DW-NOMINATE scores. First, we will describe the topics estimated, then second discuss the senator clustering.

⁵The hold-out set was created by randomly sampling 10% of tweets from each Senator with at least 10 tweets, excluding only Senator Ben Sasse.

4.1 Estimated Topics

Figure 4 displays the top seven words from each topic with shading indicating how many of the tweets in that topic are sent by Democrats. Topic 1 is Democratic Senators discussing Trump’s health care plan, and Topic 3 focuses on climate change. The emerging global pandemic is a focus of several topics, however, each pandemic-related topic captures different aspects of the discussion. Topic 4 focuses on Democrats’ efforts to pass a relief bill and Topic 5 is their efforts to improve testing. Topic 17 captures bipartisan efforts to encourage Americans to stay at home and stop the spread. The most Republican coronavirus topic, Topic 29, captures references to the virus as “chinese” and “communist.” Our topic model is able to group tweets into coherent units, separating both by subject matter and partisan valance.

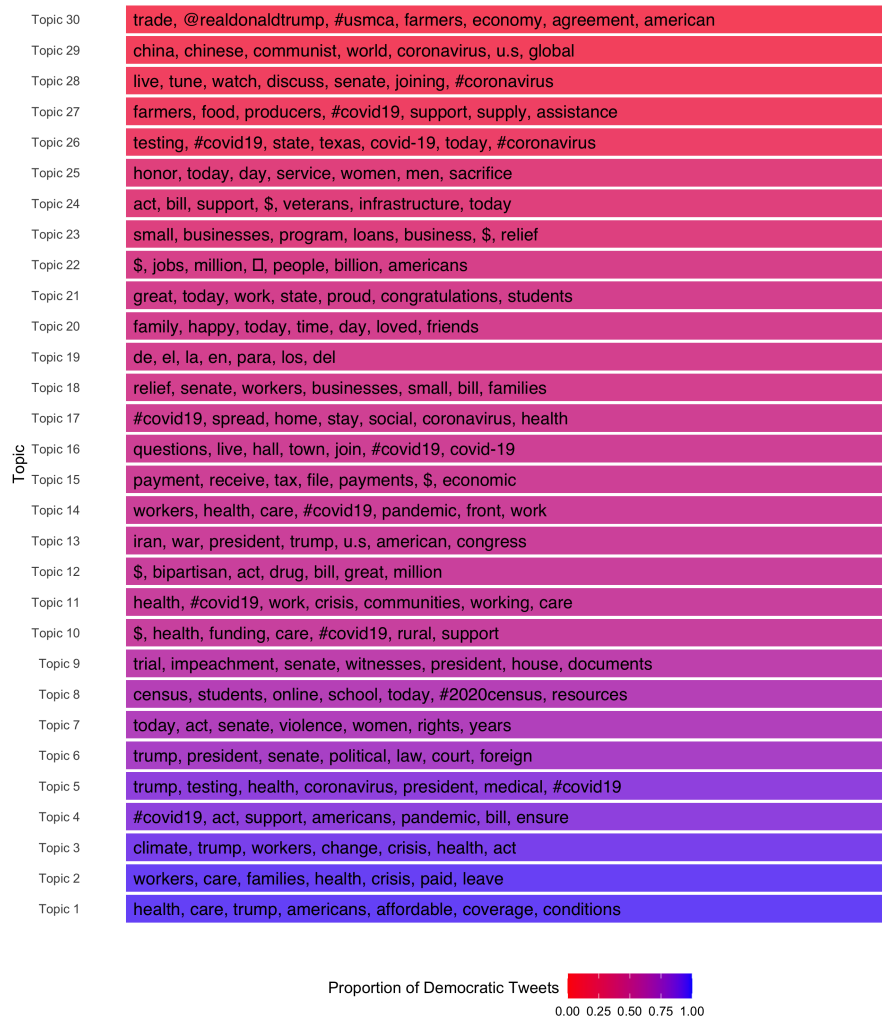


Figure 4: Top Topic Words and Topic Partisanship

Notes: Each topic is represented by its 7 most probable words. Topics are ordered and shaded by the percentage of tweets in that topic authored by Democratic Senators. We observe meaningful differences between topics in content, as well as different takes on the same content by partisanship.

We also examine the changes in the occurrence of topics across time in Figure 5. This gives additional evidence that the top words from each topic are identifying the typical tweets in said topic. Topic 9, the one with the most tweets, has top words indicating that it is about Democrats discussing President Trump’s first impeachment trial, which ran from mid-January to early February. Indeed, we see a large spike in that topic in the same time-frame, and almost no discussion afterwards. The coronavirus-related topics, particularly

Topic 17, spike mid-March and continue at an elevated level afterwards. Topic 4, Democrats lobbying for passage of legislative items to address the coronavirus, has more tweets from before the pandemic hit because many of the key words overlap with support for passage of other legislative items. Topic 5, focusing on testing, has less general words, so exhibits a more notable increasing trend. The most-partisan topics, 1 and 30, have quite different temporal structures. Topic 1 is tweets from Democrats talking about Trumps lack of a health care policy and defending the Affordable Care Act (ACA), which transition into tweets highlighting the importance of ACA reforms in light of the coronavirus pandemic.

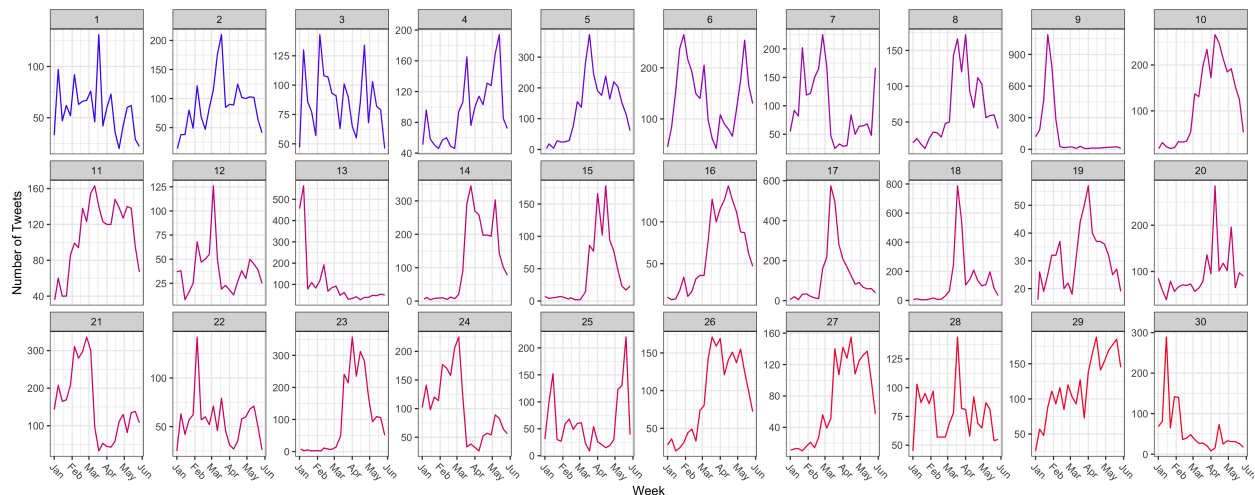


Figure 5: Weekly Topic Counts

Notes: Each panel shows the weekly number of tweets in each topic over the study period. Topics are ordered and lines are shaded by the percentage of tweets in each topic that are sent by Democratic Senators. Spikes and other time trends indicate that the topic is identifying real variation in how Senators are talking on Twitter.

4.2 Clustering

We ask our model to find four clusters among the 100 Senators active in this period of the 116th Congress.

Table 1: Senators by Cluster and Party (stLDA-C)

Cluster	Party		DW-NOMINATE	
	Democrat	Republican	1st Dim.	2nd Dim.
1	14	2	-0.19	-0.05
2	31	1	-0.35	-0.19
3	0	38	0.48	0.09
4	2	12	0.45	-0.06

Notes: Summary values for each cluster of Senators. Columns 2 and 3 report the number of Senators in each cluster by party. Columns 3 and 4 show the average DW-NOMINATE scores for Senators in each cluster. The differences between Democratic Senators on Twitter appears follow ideological divisions in the party, while the differences in Republican Twitter behavior are not.

To assess whether the clustering maps onto ideological differences, we plot the clustering and DW-NOMINATE scores for each Senator in Figure 6 and compute party counts and average DW-NOMINATE scores in Table 1. DW-NOMINATE scores are a widely-used latent measure of a politician’s ideology learned from their roll call votes. The first dimension, plotted on the x-axis, captures a liberal-conservative split and the second dimension captures other differences, commonly related to social issues (Lewis et al., 2021).⁶

⁶These scores are commonly used in the political science literature, and have even been used as a reference point to evaluate

Senators whose party is the minority in their cluster are highlighted.

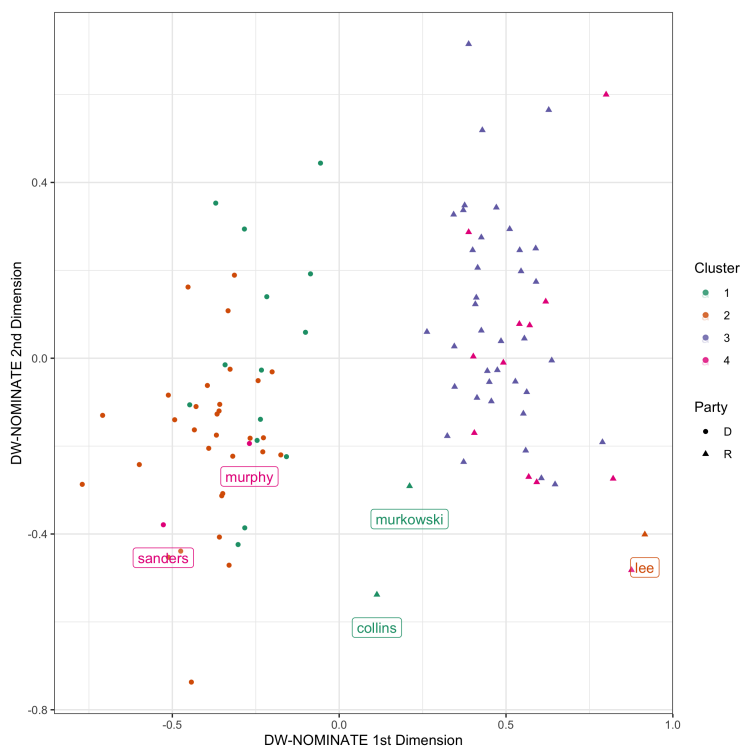


Figure 6: stLDA-C Clusters and DW-NOMINATE Scores

Notes: Each point represents a Senator in the 116th Congress and their DW-NOMINATE scores, which measure political ideology from voting behavior. Points are colored by which cluster they are classified into by stLDA-C using their Tweets in early 2020. Senators who are classified into a cluster where their party is the minority party are labeled.

Cluster 1 contains more Democrats who score higher (more conservative) on the 1st dimension and slightly higher on the 2nd dimension, and cluster 2 contains most of the remaining Democrats. Clusters 3 and 4 split the Republican Party, but do not do so as clearly along the DW-NOMINATE ideology measures. Republican Senators Murkowski and Collins of Alaska and Maine are grouped with the more moderate Democratic Senators in Cluster 1. Republican Senator Lee of Utah is in cluster 2, likely because nearly 20% of his tweets are estimated to come from topic 6, which contains many tweets from Democrats about President Trump’s first impeachment trial, focusing on threats to democracy from foreign interference. Many of the tweets after the impeachment vote are about the re-authorization Foreign Intelligence Surveillance Act (FISA), which was related to impeachment via the use of a FISA court to authorize surveillance of the Trump Campaign. The Democratic Senators clustered with Republicans are Senator Sanders of Vermont and Senator Chris Murphy of Connecticut. Senator Sanders’s most common topics are topic 2 and 22, covering 50% of his tweets. Topic 22 is 80% Republican tweets that focus on economic issues, which Senator Sanders also highlights.

To visualize the results, Figure 7 shows all individual users’ posterior expected topic distribution and the cluster-level expected topic distributions across all four clusters. Recall that Clusters 1 and 2 are majority Democrats, while 3 and 4 are majority Republicans. Topics are ordered such that Topic 1 has the largest proportion of Democratic tweets and Topic 30 the smallest. The smaller Democratic cluster, Cluster 1, has more frequent users of the mixed topics numbered 10 to 25, which is also reflected in their larger (more conservative) 1st dimension of the DW-NOMINATE score. Cluster 2 Democrats correspondingly concentrate on the lower numbered topics, Topics 1-9 especially, which are more ideologically pure. The Republican clusters are less separated on partisan topics, but rather more on variance. Cluster 3 has much

unsupervised modeling of US Senator’s Twitter networks (Barberá, 2015).

more consistency in topic usage across individuals, while Cluster 4 contains very highly variable distributions. The Senators in Cluster 4 tend to be Republicans with a larger, national profile, such as Senators Mitch McConnell, Ted Cruz, and Rand Paul. Cluster 4 Senators have more than 10 times as many Twitter followers and twice as many tweets in the time period than Cluster 3 Senators, on average.

This distinction is not surprising given the differences in party organization. Political Scientists have documented that the Republican Party has more ideological unity, while the Democratic Party is a coalition of social groups (Grossmann and Hopkins, 2016). As such, Democrats on Twitter separate into a more left-wing and a more moderate group, identified by consistent but distinct topics in their tweets. Most Republicans, on the other hand, discuss the same set of topics, but a few with national profiles use Twitter to post about a variety of highly variable topics.

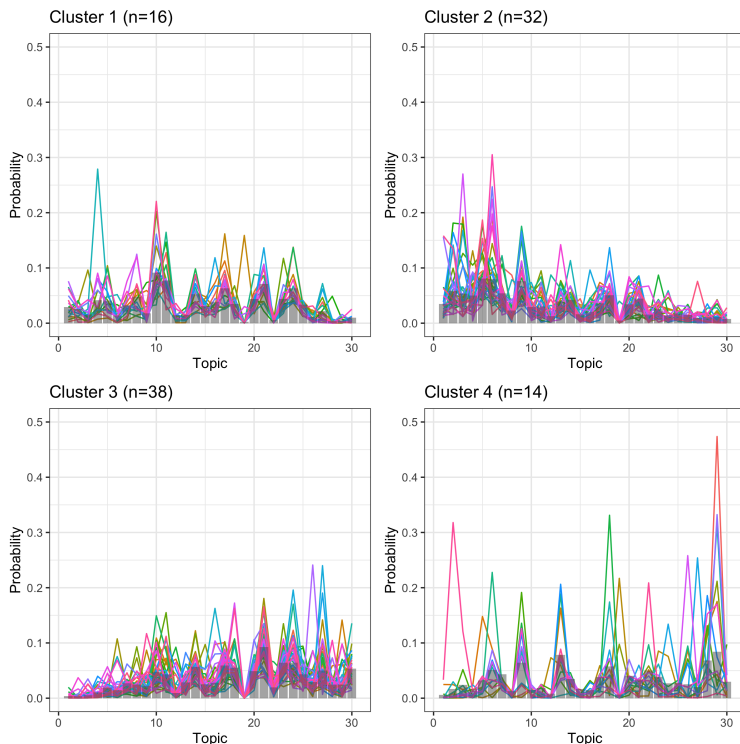


Figure 7: stLDA-C Clusters-Level and User-Level Topic Distributions

Notes: Each panel shows for a given cluster, all individual-level posterior expected topic distributions as colored lines and the cluster-level expected topic distribution as grey bars. Clusters 1 and 2 are majority Democrats, while 3 and 4 are majority Republicans. Topics are ordered such that Topic 1 has the largest proportion of Democratic tweets and Topic 30 the smallest.

4.3 Measuring the echo chamber

Given the focus on echo chambers, we calculate the number of tweets a user would need to read until they have seen at least one tweet on each of the 30 topics to measure insularity across clusters. Specifically, suppose a Twitter user follows only a single account that represents an “average” member of a given cluster g . The tweets they see are multinomial distributed with probabilities $\alpha_{gt} / \sum_t \alpha_{gt}$. We measure diversity on this Twitter user’s timeline by how many tweets do they need to view until they have seen one tweet on every topic. Work on echo chambers usually compares media consumption across predetermined groups, typically partisan affiliation, but our method allows us to measure heterogeneity of echo chambers within party for data-driven groupings of users.

If $\alpha_{gt} = 0$, clearly this number is infinite. Given our estimation method does not compute an exact zero for any cluster-topic pair, we can simulate this finite random variable for each cluster.⁷ This quantity

⁷This setup is a variant of the coupon collector problem with unequal probabilities (Ferrante and Saltalamacchia, 2014).

is summarized in columns 2-4 of Table 2. Clusters 2 and 4 would cover all 30 topics in a little bit less than 300 tweets, Cluster 1 would take almost 500, and Cluster 3 would take almost 1,500. The more liberal Democrats and national-focused Republicans cover the widest variety of topics. The local-focused Republicans in Cluster 3 have the narrowest topic choices, so provide the least diverse coverage in their tweets, reflecting the ideological unity more typical in the Republican Party (Grossmann and Hopkins, 2016).

This cluster-level “average” behavior ignores different user-level variability by cluster, which is quite relevant for the two Republican clusters. As such, we also model a Twitter user who follows n_g accounts all from the same cluster where n_g is the number of Senators in that cluster reported in Table 1. The topics of tweets they see are simulated by first drawing n_g user-topic distributions from the cluster-specific Dirichlet distribution, then sampling each topic by randomly choosing one of the n_g user-topic distributions to draw a topic from. The topics thus follow a uniform mixture of multinomial distributions. We compute the same statistics and report them in in columns 5-7 of Table 2. Accounting for this cluster-level variability notably increases the right tail of this distribution, as evidenced in the upper bounds of the 95% CIs. Accounting for this variability increases the median by about 30 tweets for the more liberal Democrat cluster and almost 150 for the national Republican cluster. Cluster 4 has the highest variability. The sum of the Dirichlet parameters is about 22, while for all other clusters it is around 55. Even though the average Senator in the cluster covers all topics at a similar rate to the liberal Democrats in Cluster 2, the within-cluster heterogeneity often produces user-level distributions that concentrate on small numbers of topics.

Rows 5 to 10 in Table 2 show the same statistics for combinations of two clusters. The best combination, the one that covers all topics the fastest, is combing either Democratic cluster with Cluster 4 Republicans. These combinations improve topic coverage over following any cluster individually. The median time is reduced by approximately 20% when following an average Senator and 30% when following individual Senators. The only inter-party combination that is worse than both intra-party combinations (combinations 1 & 2 and 3 & 4) is the moderate Democrats (Cluster 1) and local-focused Republicans (Cluster 3). Even so, the combination is a notable improvement over a user who follows Cluster 3 Republicans alone.

Table 2: Expected Number of Tweets until Full Topic Coverage.

Cluster	Average Senator			Individual Senators		
	Median	95% CI	Pr(First)	Median	95% CI	Pr(First)
1	488	(187, 1650)	9%	646	(192, 5014)	13%
2	281	(122, 794)	41%	308	(130, 997)	54%
3	1467	(246, 7530)	2%	2078	(267, 83901)	3%
4	266	(121, 721)	47%	408	(145, 2877)	30%
1 & 2	314	(128, 993)	10%	310	(114, 1204)	16%
1 & 3	493	(137, 2454)	7%	784	(135, 6493)	6%
1 & 4	218	(95, 824)	29%	245	(96, 1208)	29%
2 & 3	298	(97, 1429)	17%	324	(101, 1807)	16%
2 & 4	213	(98, 674)	28%	236	(99, 869)	26%
3 & 4	332	(117, 1022)	9%	536	(143, 3976)	6%

Notes: Columns report the number of tweets needed to cover all 30 topics by cluster(s). Columns 2-4 report statistics for a user who follows a single politician who tweets at the cluster-level expected value, i.e. the topics tweets they see are multinomial distributed with probabilities $p_t = \alpha_{gt} / \sum_t \alpha_{gt}$, for rows 1-4, or the mean of two cluster-level expected values for rows 5-10. Columns 5-7 report statistics from accounting for cluster-level variability, a twitter user who follows all politicians from a single cluster, i.e. the topics of tweets they see are a n_g -component uniform mixture of multinomial distributions with multinomial probabilities drawn from the cluster-specific Dirichlet parameters α where n_g is the number of politicians in the cluster g listed in Table 1. The median and 95% credible intervals are reported. Pr(First) computes the probability that a user who follows politicians from only the given cluster(s) sees all topics first among a cohort of users who each follow only politicians from a unique cluster (rows 1-4) or only politicians from two unique clusters (rows 5-10).

Computing the expectation exactly for 30 topics requires a summation over 2^{30} terms and is computationally challenging. As such, we use 1000 Monte Carlo samples that simulate the process to compute estimates and uncertainty measures.

4.4 Comparison to Other Methods

Table 3: Summary of Alternative Methods

(a) Senators by Cluster and Party (BTM)

Cluster	Party		DW-NOMINATE	
	Democrat	Republican	1st Dim.	2nd Dim.
1	7	32	0.33	0.10
2	19	13	0.01	-0.15
3	21	7	-0.14	-0.12
4	0	1	0.79	-0.19

(b) Senators by Cluster and Party (Sea-NMF)

Cluster	Party		DW-NOMINATE	
	Democrat	Republican	1st Dim.	2nd Dim.
1	6	12	0.21	0.10
2	7	23	0.29	0.06
3	34	17	-0.06	-0.15
4	0	1	0.79	-0.19

Notes: Summary values for each cluster of Senators with topic modeling using BTM or Sea-NMF and PAM clustering method. Columns 2 and 3 report the number of Senators in each cluster by party. Columns 3 and 4 show the average DW-NOMINATE scores for Senators in each cluster. Clusters are more variable in size and do not map as well onto ideological divisions as clusters identified by stLDA-C.

Next, we construct Table 1 and Figure 6 using BTM and Sea-NMF. We note that clustering in both of these procedures happens as a second step to topic estimation. Two-step procedures typically do not separate political parties.

Table 3a and Figure 8a describe the analysis using BTM. We estimate the model with the same number of topics and identify the same number of clusters as in the application of stLDA-C. Cluster 1 is 32 Republicans and 7 moderate Democrats who score high on the social policy dimension as well. After that, the clusters are much less clearly based on ideology. Cluster 2 is slightly more Democratic.

Table 3b and Figure 8b describe the analysis using the Sea-NMF method. Again, the clustering does not follow ideological lines as well as the stLDA-C clusters. The majority party in each cluster composes 67% to 77% of the members, and 30 Senators are in clusters where their party is in the minority. Using stLDA-C, the smallest majority is 86% and only 5 Senators are in the minority in their cluster.

Both two-step clustering procedures place Nebraska Senator Ben Sasse into his own cluster. He sent only 4 tweets from his official Senate account during the time frame, so his topic distribution is very noisily estimated when Tweets are assumed to be independent, as in both the BTM and Sea-NMF methods. Without any shrinkage, given the topics of a user’s n tweets, their topic distribution can only be integer multiples of $1/n$ across topics. When n is large, this is not a significant restriction. But when n is small, the potential topic distributions look quite different from users with large n and struggle to capture a continuous parameter like θ_u . In our method, given the topics \mathbf{Z}_u and cluster g_u , $\theta_u | \mathbf{Z}_u, g_u \sim \text{Dir}(\mathbf{Z}_u + \boldsymbol{\alpha}_{g_u})$. As the number of tweets increases, the posterior mean shifts from the normalized $\boldsymbol{\alpha}_{g_u}$ values to \mathbf{Z}_u . Thus, stLDA-C clusters Senator Sasse with other Republicans rather than clustering him alone because it recognizes that his θ_u parameter is very noisily estimated and that the few tweets he does send are similar to other Republican tweets. Social media data frequently has this kind of heterogeneity among users, some people post very frequently and others very rarely. Our model handles both types of users well by borrowing strength across multiple levels of data and shrinking noisy estimates towards typical values.

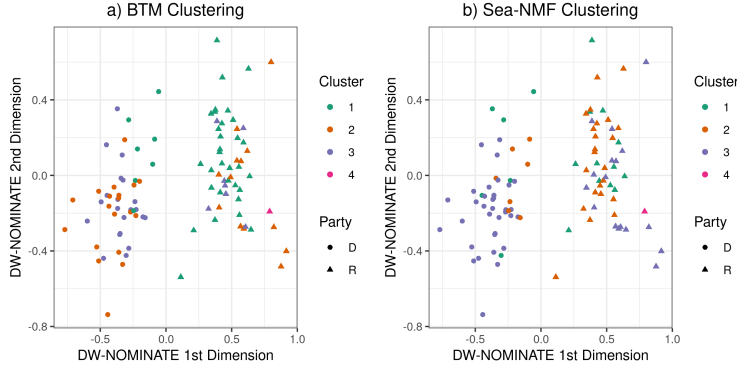


Figure 8: Alternative Method Clusters and DW-NOMINATE Scores

Notes: Each point represents a Senator in the 116th Congress and their DW-NOMINATE scores, which measure political ideology from voting behavior. Points are colored by which cluster they are classified into by BTM (panel a) or Sea-NMF (panel b) using their Tweets in early 2020.

5 MODEL VALIDATION

Tested under known conditions, the model recovers the ground truth for both topic distributions and user clusters better than other methods (where clustering might be performed as a post-estimation step) and produces more coherent topics. Section 5.1 uses data simulated from our model, and Section 5.2 uses data simulated from alternative models where stLDA-C is misspecified.

We fit stLDA-C using the collapsed Gibbs sampler derived in Section 3.1. The closest comparison model is stLDA, single topic LDA without clustering (fit using a collapsed Gibbs sampler). We also compare our model to two more recent methods, BTM, which directly models word-co-occurrence (Yan et al., 2013), and the Sea-NMF model, which uses word embeddings learned on the corpus to augment a topic model (Shi et al., 2018).⁸ We also consider tLDA on the full corpus and cLDA on documents created by combining each user’s tweets. Standard variational methods were used to estimate these models. We also compare our method to the word mover’s distance (WMD) by computing a distance matrix between documents, then running standard clustering algorithms on the results to recover topics.

We compare each model’s ability to accurately recover the true document topics and user clusters, and the coherence of the learned topics. Topic coherence is a standard comparison tool for topic models that measures how frequently high-probability words from the same topic co-occur (Mimno et al., 2011). For a given number of highest probability words from each topic, V (15 in our work), the coherence score of the topic is: $\sum_{i=2}^V \sum_{j=1}^{i-1} \log \frac{p(w_i, w_j) + 1/D}{p(w_j)}$, where $p(w_i, w_j)$ is the probability of words w_i and w_j co-occurring in the same document and $p(w_k)$ is the marginal probability of word w_k occurring. Smaller values of this score indicate greater coherence. The usual estimates of the probabilities are the co-occurrence or occurrence counts divided by the number of documents D . In short texts, this metric is known to be noisy because of the sparsity of word co-occurrences even for high-probability words (Shi et al., 2018; Quan et al., 2015). The traditional solution is to estimate the probabilities on an external corpus of longer-documents. To ensure topical similarities, we simulate 100 long documents (100 words) from each user’s observed true topic distribution using traditional LDA and calculate the coherence score on that corpus.

5.1 Simulation Study

We simulate data according to the generative model for single topic LDA with clustering (stLDA-C). To simulate realistic topic distributions, we estimated 10 word-topic distributions over 4,534 unique words from the Guardian news corpus included in the quanteda R package and used the posterior mean as the “true” topic distributions for our simulations (Watanabe and Muller, 2019). We then simulated data from 4 clusters

⁸The BTM model is implemented in the R package BTM using a Gibbs sampler (Wijffels, 2019), and the Sea-NMF model is implemented using the replication code provided by the authors on Github.

Table 4: Topic Coherence

Estimation	Simulation Model			
	Method	stLDA-C	stLDA	cLDA
cLDA		-364.02	-365.57	-362.74
tLDA		-358.49	-368.49	-363.96
BTM		-270.47	-260.89	-269.18
Sea-NMF		-268.13	-267.80	-318.65
stLDA		-256.20	-268.61	-309.87
stLDA-C		-254.59	-259.47	-242.74
True		-249.56	-253.61	-253.83

Notes: Topic coherence (Mimno et al., 2011) is calculated as $\sum_{i=2}^V \sum_{j=1}^{i-1} \log \frac{p(w_i, w_j) + 1/D}{p(w_j)}$, where $p(w_i, w_j)$ is the number of documents where words w_i and w_j both occur divided by the number of documents D . This metric is calculated on an external corpus of long documents generated from the same user-topic distributions. “True” uses the true topic distributions. We find that our method produces the most coherent topics across all estimation methods and simulation conditions.

with 10 users per cluster. We choose parameters for the simulation such that there are both similar and distinct clusters: $\alpha_{1t} = 10$ for all topics, $\alpha_{2t} = 20$ for $t \leq 5$ and 0 otherwise, $\alpha_{3t} = 20$ for $t > 5$ and 0 otherwise, and $\alpha_{4t} = 25$ for $t \in \{2, 3, 4, 5\}$ and 0 otherwise. θ_u in clusters 2 and 3 will place positive probability on only disjoint sets of five topics. Cluster 4 is extremely similar to cluster 2; it only differs by one topic. We draw 100 tweets of 13 words each per user, for a total of 4,000 documents.⁹ Note that while the data generating process is the stLDA-C model, the parameters lie outside of the prior with 0 values in certain cluster-level parameters. We use diffuse, non-informative priors of $\eta = \nu = 1$ and $\alpha_g = \mathbf{m}_g c_g$ with $\mathbf{m}_g \sim Dir(1)$ and $c_g \sim N(100, 50)$.¹⁰

Coherence. Our model produces the most coherent topic distributions. Table 4 shows the topic coherence scores averaged across topics for each applicable model. The true distributions are obviously the most coherent, but our method and stLDA without clustering are not far behind. Sea-NMF and BTM significantly outperform cLDA and tLDA. The alternative simulation models are discussed in the next section.

Topic recovery. The first column of Table 5 compares document topic recovery for our method and comparison methods under the stLDA-C generative model. Each iteration of the Gibbs sampler provides a sample of Z_{ud} , so the proportion of times $Z_{ud} = t$ is a posterior estimate of $P(Z_{ud} = t | \mathbf{W})$. Documents are classified into the highest posterior-probability topic. Clustering is evaluated with normalized mutual information (NMI), which measures the mutual information between the estimated and ground truth clusters normalized to be between 0 (no information shared) and 1 (perfect correlation). The more novel BTM and Sea-NMF methods are able to recover many topics accurately, but they are notably worse than the single-topic models. stLDA without clustering is able to recover document topics with comparable accuracy to stLDA-C.

Cluster recovery. Our approach has the best cluster recovery with each user being mapped to the correct cluster. Table 6 shows the results. To highlight the places other methods fail in the clustering, we show cluster-specific true- and false-positive rates rather than a global metric. The other methods must rely on a two-stage cluster estimation procedure: user-topic distributions are calculated as the proportion of each users’ tweets that are classified into each topic, then those distributions are passed to the PAM clustering algorithm, a more stable variant of k-means clustering (Reynolds et al., 2006). The principal confusion is between clusters 2 and 4, which only differ by one topic. All alternative models identify that clusters 2 and 4 differ on topic 1, but some users from cluster 2 who talk about topic 1 less than the average user are grouped into cluster 4. The independent clustering necessarily does not leverage the correct distance metric between topic distributions. Only our method stLDA-C, which unifies the estimation of clusters and topics, is able to recover the clusters with 100% accuracy because it recognizes the differences across topic 1 are the most significant differences between clusters 2 and 4.

⁹The average tweet in the data application has 13 words.

¹⁰The same priors were used in Section 3.3.

Table 5: Document Topic Recovery (NMI)

Method	stLDA-C DGP	stLDA DGP
tLDA	0.052	0.023
WMD	0.247	0.233
BTM	0.890	0.878
Sea-NMF	0.925	0.853
stLDA	0.952	0.904
stLDA-C	0.955	0.953

Notes: Topic recovery is evaluated with normalized mutual information, a measure of cross-entropy from the estimated and true labels. An NMI of 1 indicates perfect recovery of the ground truth. For topic models, documents are classified by the highest posterior probability topic. For WMD, the partitioning around medoids (PAM) clustering method is used, a more robust version of k-means (Reynolds et al., 2006). stLDA-C recovers nearly every document’s topic, while WMD and tLDA are not much better than random guesswork.

Table 6: Cluster Recovery from stLDA-C Simulated Data

True Cluster	WMD		BTM		Sea-NMF		stLDA		stLDA-C	
	TPR	FPR	TPR	FPR	TPR	FPR	TPR	FPR	TPR	FPR
1	100%	3%	100%	0%	100%	0%	100%	0%	100%	0%
2	60%	23%	70%	0%	70%	0%	70%	0%	100%	0%
3	100%	0%	100%	0%	100%	0%	100%	0%	100%	0%
4	70%	20%	100%	10%	100%	10%	100%	10%	100%	0%

Notes: TPR for cluster c is the largest proportion of users with cluster c mapped to a single estimated cluster. FPR is the proportion of users with that estimated cluster among users in cluster $c' \neq c$. For stLDA-C, users are classified into the highest posterior probability cluster. For other methods, user-level topic distributions are passed to the the PAM clustering method to identify 4 clusters (Reynolds et al., 2006). stLDA-C recovers all clusters accurately, while two-step procedures struggle to separate clusters 2 and 4. Cluster 4 contains users who talk about all but one of the topics used by users in cluster 2.

A useful visualization of the results compares the expected topic distribution of each cluster and the variability of users within a cluster. Figure 9 shows each estimated user topic distribution on top of the cluster mean. Our model recovers the ground truth and captures the fact that there are four groups of users with similar but not identical topic distributions.

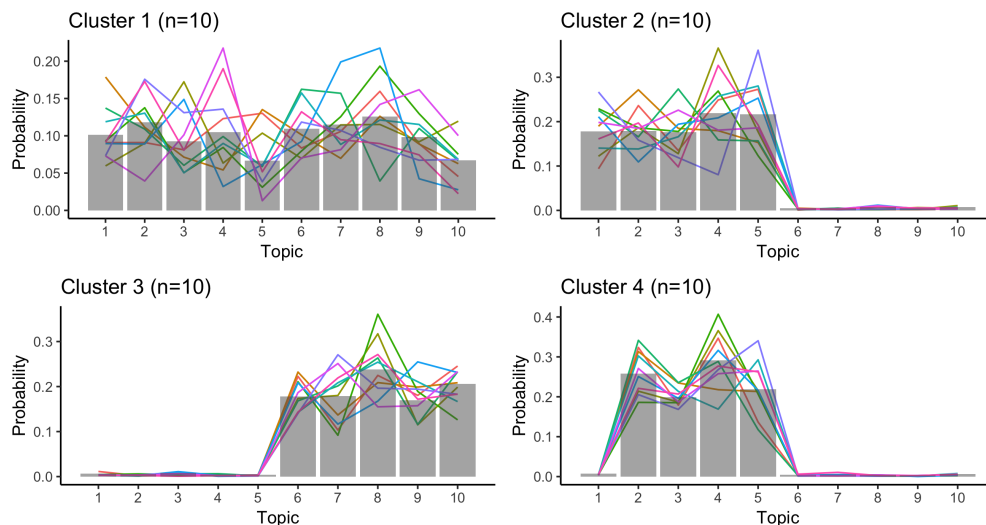


Figure 9: User Topic Distributions by Cluster (Simulated Data)

Notes: Cluster-level expected values are shown in grey. Each user’s estimated topic distribution is shown with a colored line. stLDA-C correctly identifies which topics each cluster discusses, capturing both the cluster-level averages and user-level variability.

5.2 Model Miss-Specification

In the above simulations, we generated data according to the stLDA-C generative model. In this section, we simulate data according to two alternative models and test whether our method still performs well. We preserve the size of the simulated data from simulations for stLDA-C to allow comparison of results across simulations.

The first alternative simulation model is simple stLDA, a submodel of stLDA-C. Effectively, data are simulated from only one cluster. Each θ_u is sampled from a uniform distribution over the 10 dimensional simplex, then 100 topics are sampled from θ_u , and one document for each sampled topic is generated. We compare the same models in the last section in terms of document-topic recovery, user-cluster recovery, and topic coherence.

We use the same parameter settings for our comparison models. Notably, we allow stLDA-C to use the same maximum number of clusters (four), and discover that the model correctly identifies that only one cluster is needed. Other comparison models rely entirely on a researcher-specified number of clusters, rather than specifying a maximum number and letting the data identify the true number of clusters needed.

Table 4 shows the topic coherence scores for each model in this simulation in the second column. stLDA-C and BTM are again quite similar, with our model estimating slightly more coherent topics. Column 2 in Table 5 shows the NMI for learned document topics for each method under stLDA simulations. All four of the most advanced methods have high NMI values. Our model has the highest of the four, even higher than the correctly-specified stLDA method. This simulation shows that our setup is able to adapt to sub-models.

Next, we simulate data from the cLDA generative process by generating one long document for each user, then breaking it up into short, tweet-length documents. Thus, our model is miss-specified since it assumes a single topic for each document, whereas the truth is that multiple topics are present. We generate user-specific topic distributions with the same clustering as in the stLDA-C simulations: one cluster discusses all topics, one cluster only half of the topics, another cluster discuss the other half, and the final cluster

Table 7: Cluster Recovery from cLDA Simulated Data

True Cluster	WMD		BTM		Sea-NMF		stLDA		stLDA-C	
	TPR	FPR	TPR	FPR	TPR	FPR	TPR	FPR	TPR	FPR
1	70%	7%	90%	0%	100%	0%	90%	0%	100%	0%
2	100%	20%	100%	37%	50%	10%	100%	37%	100%	33%
3	100%	3%	50%	0%	100%	0%	60%	0%	100%	0%
4	40%	40%	100%	37%	70%	17%	100%	37%	100%	33%

TPR for cluster c is the largest proportion of users with cluster c mapped to a single estimated cluster. FPR is the proportion of users with that estimated cluster among users in cluster $c' \neq c$. For stLDA-C, users are classified into the highest posterior probability cluster. For other methods, user-level topic distributions are passed to the the PAM clustering method to identify 4 clusters (Reynolds et al., 2006). stLDA-C has the best cluster recovery. It separates users into three clusters, successfully separating all members of clusters 1 and 3, and it combining clusters 2 and 4 which only differ on one topic. All other procedures fail to fully separate 1 and 3 and do not correctly identify the similarities among clusters 2 and 4.

discusses all but one of topics another cluster uses. We compare all models on cluster recovery and topic coherence.

Our model correctly identifies that clusters 1 and 3 are distinct, but merges the clusters that differ on only one topic into a single cluster. Only three of the maximum of four clusters have users in them. Other methods are able to somewhat replicate this recovery, but none are as good as stLDA-C. They confuse outliers in clusters two and four with members of the other cluster and include some members of cluster 1, who discuss all topics, in the combined clusters of 2 and 4. The results are shown in Table 7.

Our model also performs best in terms of topic coherence. We again simulate many long documents from the same user-topic distributions to compare the quality of estimated topics. Column 3 in Table 4 shows that our method even marginally out-performs the true topic distributions in terms of coherence.

6 DISCUSSION

Social media data and short text more broadly are becoming an important area of study for researchers. Standard models for long text already have to address issues regarding high dimensionality when treating words as features. Models for short text have to address those same issues with the added complication of increased sparsity in observed word co-occurrences and frequent desire to use model results in a downstream task. stLDA-C solves all of those issues by linking the topics of words used in the same short text, borrowing strength across documents generated by the same user, and clustering users with a hierarchical prior that accounts for uncertainty in user-level parameters.

Our model is able to recover topics and clusters of users in simulated data, whereas other state-of-the-art methods struggle with one or both tasks. We developed the unified estimation procedure for topics and clusters because of the frequent desire to categorize at both the post and user level simultaneously. Our variational approximation allows our model to scale to large corpora typical of social media data. As demonstrated in the application, those user-level groupings are often related to important user features such as political party. A further extension of our algorithm could use the information on these features in the cluster estimation itself, such as in Roberts et al. (2013). We use only topic frequencies to group users, but one could build upon our clustering step models that leverage demographic information as well.

Even without including external information about users, our application demonstrates that such information can be gleaned from the clustering. We capture distinct partisan topics and group US Senators into clusters that reflect partisan ideology. Our model separates tweets about several different aspects of the coronavirus pandemic, and identifies distinct subgroups within each party that tweet differently. The novel measure of the echo-chamber characterizes the degree of this separation and has important ramifications for understanding the scope and character of exposure to cross-partisans on social media. To learn what the other side really thinks about current events, citizens will need to actively seek out discussion from leaders in the opposing party.

References

- Anandkumar, A., Foster, D. P., Hsu, D. J., Kakade, S. M., and Liu, Y.-k. (2012). “A Spectral Algorithm for Latent Dirichlet Allocation.” In Pereira, F., Burges, C. J. C., Bottou, L., and Weinberger, K. Q. (eds.), *Advances in Neural Information Processing Systems*, volume 25. Curran Associates, Inc. [11](#)
- Anandkumar, A., Hsu, D. J., Janzamin, M., and Kakade, S. M. (2013). “When are Overcomplete Topic Models Identifiable? Uniqueness of Tensor Tucker Decompositions with Structured Sparsity.” In *Advances in Neural Information Processing Systems*, 1986–1994. [11](#)
- Bail, C. A., Argyle, L. P., Brown, T. W., Bumpus, J. P., Chen, H., Hunzaker, M. B. F., Lee, J., Mann, M., Merhout, F., and Volfovsky, A. (2018). “Exposure to opposing views on social media can increase political polarization.” *Proceedings of the National Academy of Sciences*, 115(37): 9216–9221. [1](#)
- Bakshy, E., Messing, S., and Adamic, L. A. (2015). “Exposure to ideologically diverse news and opinion on Facebook.” *Science*, 348(6239): 1130–1132. [1](#), [5](#)
- Barberá, P. (2015). “Birds of the same feather tweet together: Bayesian ideal point estimation using twitter data.” *Political Analysis*, 23(1): 76–91. [1](#), [14](#)
- Barberá, P., Casas, A., Nagler, J., EGAN, P. J., Bonneau, R., Jost, J. T., and Tucker, J. A. (2019). “Who leads? who follows? measuring issue attention and agenda setting by legislators and the mass public using social media data.” *American Political Science Review*. [1](#), [2](#), [5](#)
- Barnard, K., Duygulu, P., Forsyth, D., de Freitas, N., Blei, D. M., and Jordan, M. I. (2003). “Matching Words and Pictures.” *J. Mach. Learn. Res.*, 3: 1107–1135. [3](#)
- Blei, D. M., Ng, A. Y., and Jordan, M. I. (2003). “Latent Dirichlet Allocation.” *J. Mach. Learn. Res.*, 3: 993–1022. [1](#), [3](#), [7](#), [9](#)
- Chang, J., Boyd-Graber, J., Gerrish, S., Wang, C., and Blei, D. M. (2009). “Reading tea leaves: How humans interpret topic models.” *Advances in Neural Information Processing Systems 22 - Proceedings of the 2009 Conference*, 288–296. [11](#)
- Donoho, D. and Stodden, V. (2004). “When does non-negative matrix factorization give a correct decomposition into parts?” *Advances in Neural Information Processing Systems*, 16: 1141–1148. [11](#)
- Ferrante, M. and Saltalamacchia, M. (2014). “The coupon collector’s problem.” *Materials matemàtics*, 0001–35. [15](#)
- Griffiths, T. L. and Steyvers, M. (2004). “Finding scientific topics.” *Proceedings of the National Academy of Sciences*, 101(suppl 1): 5228–5235. [6](#), [7](#)
- Grossmann, M. and Hopkins, D. A. (2016). *Asymmetric politics: Ideological Republicans and group interest Democrats*. Oxford University Press. [15](#), [16](#)
- Hajjem, M. and Latiri, C. (2017). “Combining IR and LDA Topic Modeling for Filtering Microblogs.” *Procedia Computer Science*, 112: 761 – 770. Knowledge-Based and Intelligent Information & Engineering Systems: Proceedings of the 21st International Conference, KES-20176-8 September 2017, Marseille, France. [2](#)
- Hong, L. and Davison, B. D. (2010). “Empirical study of topic modeling in twitter.” In *Proceedings of the first workshop on social media analytics*, 80–88. acm. [2](#), [5](#)
- Hu, X., Zhang, X., Lu, C., Park, E. K., and Zhou, X. (2009). “Exploiting Wikipedia As External Knowledge for Document Clustering.” In *Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD ’09, 389–396. New York, NY, USA: ACM. [2](#)
- Huang, K., Fu, X., and Sidiropoulos, N. D. (2016). “Anchor-free correlated topic modeling: Identifiability and algorithm.” *Advances in Neural Information Processing Systems*, 1794–1802. [11](#)

- Jaidka, K., Zhou, A., and Lelkes, Y. (2019). “Brevity is the soul of Twitter: The constraint affordance and political discussion.” *Journal of Communication*, 69(4): 345–372. [1](#)
- Jin, O., Liu, N. N., Zhao, K., Yu, Y., and Yang, Q. (2011). “Transferring topical knowledge from auxiliary long texts for short text clustering.” In *Proceedings of the 20th ACM international conference on Information and knowledge management*, 775–784. ACM. [2](#)
- Kusner, M., Sun, Y., Kolkin, N., and Weinberger, K. (2015). “From word embeddings to document distances.” In *International conference on machine learning*, 957–966. [3](#)
- Lewandowsky, S., Jetter, M., and Ecker, U. K. (2020). “Using the president’s tweets to understand political diversion in the age of social media.” *Nature communications*, 11(1): 1–12. [1](#)
- Lewis, J. B., Poole, K., Rosenthal, H., Boche, A., Rudkin, A., and Sonnet, L. (2021). *Voteview: Congressional Roll-Call Votes Database*. [13](#)
- Li, C., Wang, H., Zhang, Z., Sun, A., and Ma, Z. (2016). “Topic modeling for short texts with auxiliary word embeddings.” In *Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval*, 165–174. [3](#), [6](#)
- Li, X., Zhang, J., and Ouyang, J. (2019). “Dirichlet multinomial mixture with variational manifold regularization: Topic modeling over short texts.” In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, 7884–7891. [2](#)
- Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013). “Efficient estimation of word representations in vector space.” *arXiv preprint arXiv:1301.3781*. [3](#)
- Mimno, D., Wallach, H. M., Talley, E., Leenders, M., and McCallum, A. (2011). “Optimizing semantic coherence in topic models.” In *Proceedings of the conference on empirical methods in natural language processing*, 262–272. Association for Computational Linguistics. [11](#), [18](#), [19](#)
- Mosleh, M., Martel, C., Eckles, D., and Rand, D. G. (2021). “Shared partisanship dramatically increases social tie formation in a Twitter field experiment.” *Proceedings of the National Academy of Sciences*, 118(7). [1](#)
- Nigam, K., McCallum, A. K., Thrun, S., and Mitchell, T. (2000). “Text classification from labeled and unlabeled documents using EM.” *Machine learning*, 39(2-3): 103–134. [5](#)
- Pennacchiotti, M. and Popescu, A.-M. (2011). “A machine learning approach to twitter user classification.” In *Fifth International AAAI Conference on Weblogs and Social Media*. [2](#), [5](#)
- Poole, K. T. and Rosenthal, H. (2017). *Ideology & congress: A political economic history of roll call voting*. Routledge. [2](#)
- Quan, X., Kit, C., Ge, Y., and Pan, S. J. (2015). “Short and sparse text topic modeling via self-aggregation.” In *Twenty-Fourth International Joint Conference on Artificial Intelligence*. [18](#)
- Rangarajan Sridhar, V. K. (2015). “Unsupervised Topic Modeling for Short Texts Using Distributed Representations of Words.” In *Proceedings of the 1st Workshop on Vector Space Modeling for Natural Language Processing*, 192–200. Denver, Colorado: Association for Computational Linguistics. [3](#)
- Reynolds, A. P., Richards, G., de la Iglesia, B., and Rayward-Smith, V. J. (2006). “Clustering Rules: A Comparison of Partitioning and Hierarchical Clustering Algorithms.” *Journal of Mathematical Modelling and Algorithms*, 5(4): 475–504. [19](#), [20](#), [22](#)
- Roberts, M. E., Stewart, B. M., Tingley, D., Airolidi, E. M., et al. (2013). “The structural topic model and applied social science.” In *Advances in neural information processing systems workshop on topic models: computation, application, and evaluation*, volume 4, 1–20. Harrahs and Harveys, Lake Tahoe. [22](#)

- Rodriguez, A., Dunson, D. B., and Gelfand, A. E. (2008). “The nested Dirichlet process.” *Journal of the American statistical Association*, 103(483): 1131–1154. 5
- Rosen-Zvi, M., Griffiths, T., Steyvers, M., and Smyth, P. (2004). “The author-topic model for authors and documents.” In *Proceedings of the 20th conference on Uncertainty in artificial intelligence*, 487–494. AUAI Press. 2
- Sahami, M. and Heilman, T. D. (2006). “A Web-based Kernel Function for Measuring the Similarity of Short Text Snippets.” In *Proceedings of the 15th International Conference on World Wide Web*, WWW ’06, 377–386. New York, NY, USA: ACM. 2
- Shi, T., Kang, K., Choo, J., and Reddy, C. K. (2018). “Short-Text Topic Modeling via Non-negative Matrix Factorization Enriched with Local Word-Context Correlations.” In *Proceedings of the 2018 World Wide Web Conference*, WWW ’18, 1105–1114. Republic and Canton of Geneva, Switzerland: International World Wide Web Conferences Steering Committee. 3, 18
- Steinskog, A., Therkelsen, J., and Gambäck, B. (2017). “Twitter Topic Modeling by Tweet Aggregation.” In *Proceedings of the 21st Nordic Conference on Computational Linguistics*, 77–86. Gothenburg, Sweden: Association for Computational Linguistics. 2, 5
- Teh, Y. W., Jordan, M. I., Beal, M. J., and Blei, D. M. (2006). “Hierarchical dirichlet processes.” *Journal of the american statistical association*, 101(476): 1566–1581. 5
- Wallach, H. M. (2008). “Structured topic models for language.” Ph.D. thesis, University of Cambridge Cambridge, UK. 3, 5
- Wang, F., Zhang, J. L., Li, Y., Deng, K., and Liu, J. S. (2021). “Bayesian Text Classification and Summarization via A Class-Specified Topic Model.” *Journal of Machine Learning Research*, 22(89): 1–48. 5
- Watanabe, K. and Muller, S. (2019). “Quanteda Tutorials.” <https://tutorials.quanteda.io>. 18
- Wijffels, J. (2019). *BTM: Biterm Topic Models for Short Text*. R package version 0.2.1. 18
- Xie, P. and Xing, E. P. (2013). “Integrating document clustering and topic modeling.” *arXiv preprint arXiv:1309.6874*. 3
- Yan, X., Guo, J., Lan, Y., and Cheng, X. (2013). “A Biterm Topic Model for Short Texts.” In *Proceedings of the 22Nd International Conference on World Wide Web*, WWW ’13, 1445–1456. New York, NY, USA: ACM. 2, 3, 9, 18
- Yin, J. and Wang, J. (2014). “A Dirichlet Multinomial Mixture Model-based Approach for Short Text Clustering.” In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD ’14, 233–242. New York, NY, USA: ACM. 2, 5
- Zhao, W. X., Jiang, J., Weng, J., He, J., Lim, E.-P., Yan, H., and Li, X. (2011). “Comparing twitter and traditional media using topic models.” In *European conference on information retrieval*, 338–349. Springer. 2, 3, 5

A Collapsed Gibbs Sampler

Proof for the collapsed Gibbs sampler. Z_{ud} is the topic of tweet d by user u . θ_u is the user's topic distribution, β_t is the topic distribution t over words. We desire $P(Z_{ud} = t | \mathbf{Z}_{-\mathbf{ud}}, \mathbf{W})$, the distribution of Z_{ud} given the topics of all tweets except tweet d by user u and the word counts of all tweets \mathbf{W} . Given the topic of a tweet $Z_{ud} = t$, the word counts $\mathbf{W}_{\mathbf{ud}}$ are a multinomial draw from topic distribution β_t . θ_u identifies the unconditional, prior distribution for Z_{ud} . Via Bayes' rule, the full conditional for $P(Z_{ud} = t | \theta_u, \beta_t, \mathbf{Z}_{-\mathbf{ud}}, \mathbf{W})$ is proportional to the product of this multinomial likelihood and θ_u prior. Thus, given β_t and θ_u , the posterior for Z_{ud} does not depend on any words or topics besides those in tweet ud .

$$P(Z_{ud} = t | \theta_u, \beta_t, \mathbf{Z}_{-\mathbf{ud}}, \mathbf{W}) \propto P(\mathbf{W}_{\mathbf{ud}} | Z_{ud} = t, \beta_t) P(Z_{ud} = t | \theta_t) \propto \prod_i \beta_{ti}^{W_{ud}^{ti}} \theta_{ut}$$

From the paper:

$$\begin{aligned} p(Z_{ud} = t | \mathbf{Z}_{-\mathbf{ud}}, \mathbf{W}) &= p(Z_{ud} = t | \mathbf{Z}_{-\mathbf{ud}}, \mathbf{W}_{\mathbf{ud}}, \mathbf{W}_{-\mathbf{ud}}) \\ &\propto p(\mathbf{W}_{\mathbf{ud}} | Z_{ud} = t, \mathbf{Z}_{-\mathbf{ud}}, \mathbf{W}_{-\mathbf{ud}}) p(Z_{ud} = t | \mathbf{Z}_{-\mathbf{ud}}, \mathbf{W}_{-\mathbf{ud}}) \\ &= \int p(\mathbf{W}_{\mathbf{ud}} | Z_{ud} = t, \beta_t) p(\beta_t | \mathbf{Z}_{-\mathbf{ud}}, \mathbf{W}_{-\mathbf{ud}}) d\beta_t \times \\ &\quad \int p(Z_{ud} = t | \theta_u) p(\theta_u | \mathbf{Z}_{-\mathbf{ud}}) d\theta_u \\ &= \int \prod_i \beta_{ti}^{W_{ud}^{ti}} p(\beta_t | \mathbf{Z}_{-\mathbf{ud}}, \mathbf{W}_{-\mathbf{ud}}) d\beta_t \times \\ &\quad \int \theta_{ut} p(\theta_u | \mathbf{Z}_{-\mathbf{ud}}) d\theta_u \end{aligned}$$

Conditional on $G_u = g$ and α_g , the unconditional distribution of θ_u is: $p(\theta_u) \sim \text{Dir}(\alpha_g)$. The Dirichlet is conjugate with multinomial data. Let $\mathbf{Z}_{-\mathbf{d}}^{(\mathbf{u})}$ denote the topic counts of user u excluding the topic of tweet d .

$$\begin{aligned} p(\theta_u | G_u = g, \alpha_g, \mathbf{Z}_{-\mathbf{ud}}) &\propto p(\mathbf{Z}_{-\mathbf{d}}^{(\mathbf{u})} | \theta_u) p(\theta_u | G_u = g, \alpha_g) \\ &\propto \prod_{t=1}^T \theta_{ut}^{Z_{-d,t}^{(u)}} \prod_{t=1}^T \theta_{ut}^{\alpha_{gt}-1} \\ &= \prod_{t=1}^T \theta_{ut}^{Z_{-d,t}^{(u)} + \alpha_{gt}-1} \\ &\propto \text{Dir}(\mathbf{Z}_{-\mathbf{d}}^{(\mathbf{u})} + \alpha_g) \end{aligned}$$

Using this result in the first integral, we have:

$$\begin{aligned} \int \theta_{ut} p(\theta_u | \mathbf{Z}_{-\mathbf{ud}}) d\theta_u &= \frac{\Gamma(\sum_j Z_{-d,j}^{(u)} + \alpha_{gj})}{\prod_j \Gamma(Z_{-d,j}^{(u)} + \alpha_{gj})} \int \theta_{ut} \prod_{j=1}^T \theta_{uj}^{Z_{-d,j}^{(u)} + \alpha_{gj}-1} d\theta_u \\ &= \frac{\Gamma(\sum_j Z_{-d,j}^{(u)} + \alpha_{gj})}{\prod_j \Gamma(Z_{-d,j}^{(u)} + \alpha_{gj})} \frac{\prod_j \Gamma(Z_{-d,j}^{(u)} + \alpha_{gj} + I(j=t))}{\Gamma(\sum_t Z_{-d,j}^{(u)} + \alpha_{gj} + 1)} \\ &= \frac{Z_{-d,t}^{(u)} + \alpha_{gt}}{\sum_j Z_{-d,j}^{(u)} + \alpha_{gj}} \end{aligned}$$

This is the result in the paper. $\Gamma(\cdot)$ denotes the gamma function and $I(\cdot)$ the indicator function. The integral is computed by recognizing the kernel of a Dirichlet distribution with the same parameters as the posterior with 1 added the t 'th parameter.

The second integral is computed similarly. For documents with $Z_{ud} = t$, the word counts are a multinomial draw of size n from the probability distribution β_t . Thus, the posterior can be expressed as: $p(\beta_t | \mathbf{Z}_{-ud}, \mathbf{W}_{-ud}) \propto p(\beta_t | \mathbf{W}_{-ud}^{(t)})$, conditioning only on the word counts from documents with $Z_{ud} = t$. With a $\text{Dir}(\eta)$ prior on β_t , the posterior also Dirichlet with parameters $\mathbf{W}_{-ud}^{(t)} + \eta$ by the same conjugate update shown above. Thus, the second integral simplifies to:

$$\begin{aligned} \int \prod_i \beta_{ti}^{W_{ud,i}} p(\beta_t | \mathbf{Z}_{-ud}, \mathbf{W}_{-ud}) d\beta_t &= \frac{\Gamma(N\eta + \sum_i W_{-ud,i}^{(t)})}{\prod_i \Gamma(\eta + W_{-ud,i}^{(t)})} \int \prod_i \beta_{ti}^{W_{ud,i}} \prod_i \beta_{ti}^{W_{-ud,i} + \eta} d\beta_t \\ &= \frac{\Gamma(N\eta + \sum_i W_{-ud,i}^{(t)})}{\prod_i \Gamma(\eta + W_{-ud,i}^{(t)})} \frac{\prod_i \Gamma(\eta + W_{-ud,i}^{(t)} + W_{ud,i})}{\Gamma(N\eta + \sum_i W_{-ud,i}^{(t)} + W_{ud,i})} \end{aligned}$$

This is the result in the paper. The integral is computed by recognizing the kernel of a Dirichlet distribution with the same parameters as the posterior for β_t with the word counts for tweet ud added to the respective parameters. Some simplification of the gamma functions are possible. They do not add to the intuitive understanding, so we do not show them here.

B Variational stLDA-C

This section derives the variational approximation to the posterior from Section 3.2 in the paper. This follows the same general procedure as the variational approximation derived in Blei et al. (2003), and we note the similarities where appropriate.

B.1 ELBO Bound

Here we derive a evidence lower bound (ELBO) on the KL divergence between the variational distribution and the posterior. First, we will build the ELBO for an individual user after re-introducing the notation from the main paper.

\mathbf{W}_u is the $(n_u \times V)$ document-term matrix for user u . The latent variables are: g_u their cluster membership, θ_u their topic distribution, and \mathbf{Z}_u the topics of their posts. For simplicity, I omit the $_u$ notation for the remainder of this section. The model parameters are: α the $G \times T$ matrix of Dirichlet parameters for each cluster, β the $T \times V$ matrix of topic distributions, and ξ the corpus-level cluster proportions. The variational distribution q is expressed as:

$$q(g, \theta, \mathbf{Z} | \lambda, \gamma, \phi) = q(g | \lambda) q(\theta | \gamma) \prod_{d=1}^{n_u} q(z_d | \phi_d)$$

Here $q(g | \lambda)$ is $\text{Cat}(\lambda)$, $q(\theta | \gamma)$ is $\text{Dir}(\gamma)$, and $q(z_d | \phi_d)$ is $\text{Cat}(\phi_d)$. Both λ and ϕ_d are points on the probability simplex; γ is a T dimensional vector of positive numbers. Although not made explicit in the notation, $\{\lambda, \gamma, \phi\}$ are all *user-specific* quantities. For consistency with Blei et al. (2003) notation, I let ξ denote the cluster mixing proportions (ϕ in the main paper) and let ϕ_d be the post-specific variational parameters. Now, we express the probability of the user's posts.

$$\begin{aligned}
\log p(\mathbf{W}|\alpha, \beta, \xi) &= \log \sum_g \int_{\theta} \sum_z p(\theta, \mathbf{Z}, g|\alpha, \beta, \xi) \\
&= \log \sum_g \int_{\theta} \sum_z p(\theta, \mathbf{Z}, g|\alpha, \beta, \xi) \frac{q(g, \theta, \mathbf{Z})}{q(g, \theta, \mathbf{Z})} \\
&= \log E_q \left[\frac{p(\theta, \mathbf{Z}, g|\alpha, \beta, \xi)}{q(g, \theta, \mathbf{Z})} \right] \\
&\geq E_q[\log p(g, \theta, \mathbf{Z}, \mathbf{W}|\alpha, \beta, \xi)] - E_q[\log q(g, \theta, \mathbf{Z})] := L
\end{aligned}$$

The last step is from applying Jensen's inequality to push the log through the expectation. E_q notes the expectation is taken with respect to the variational distribution q .

As in normal LDA, the difference between the left and right hand terms ($\log p(\mathbf{W}|\alpha, \beta, \xi)$ and the lower bound, L) is the KL divergence between the variational posterior and the true posterior. Let $D(A||B)$ denote the KL divergence between A and B.

$$\begin{aligned}
E_q[\log p(g, \theta, \mathbf{Z}, \mathbf{W}|\alpha, \beta, \xi)] - E_q[\log q(g, \theta, \mathbf{Z})] &= \log p(\mathbf{W}|\alpha, \beta, \xi) + E_q[\log p(g, \theta, \mathbf{Z}|\mathbf{W}, \alpha, \beta, \xi)] - E_q[\log q(g, \theta, \mathbf{Z})] \\
&= \log p(\mathbf{W}|\alpha, \beta, \xi) - D(q(g, \theta, \mathbf{Z})||p(g, \theta, \mathbf{Z}|\mathbf{W}, \alpha, \beta, \xi))
\end{aligned}$$

The likelihood of the data and latent parameters can be factored into simpler conditional distributions.

$$p(g, \theta, \mathbf{Z}, \mathbf{W}|\alpha, \beta, \xi) = p(g|\xi)p(\theta|g, \alpha)p(\mathbf{Z}|\theta)p(\mathbf{W}|\mathbf{Z}, \beta)$$

Thus, the ELBO (L) can be expressed as (reintroducing variational parameters):

$$\begin{aligned}
L(\lambda, \gamma, \phi; \alpha, \beta, \xi) &= E_q[\log p(g, \theta, \mathbf{Z}, \mathbf{W}|\alpha, \beta, \xi)] - E_q[\log q(g, \theta, \mathbf{Z})] \\
&= E_q[\log p(g|\xi)] + E_q[\log p(\theta|g, \alpha)] + E_q[\log p(\mathbf{Z}|\theta)] + E_q[\log p(\mathbf{W}|\mathbf{Z}, \beta)] \\
&\quad - E_q[\log q(g)] - E_q[\log q(\theta)] - E_q[\log q(\mathbf{Z})]
\end{aligned}$$

Each expectation can be simplified. Each line below corresponds to expansion of the seven terms above, in order.

$$L(\lambda, \gamma, \phi; \alpha, \beta, \xi) = \sum_g \lambda_g \log(\xi_g) \tag{6}$$

$$+ \sum_g \lambda_g \left(\log \Gamma\left(\sum_{t=1}^T \alpha_{gt}\right) - \sum_{t=1}^T \log \Gamma(\alpha_{gt}) + \sum_{t=1}^T (\alpha_{gt} - 1) E_q[\log \theta_t] \right) \tag{7}$$

$$+ \sum_{d=1}^n \sum_{t=1}^T \phi_{dt} E_q[\log \theta_t] \tag{8}$$

$$+ \sum_{d=1}^n \sum_{t=1}^T \phi_{dt} \left[\log \left(\frac{(\sum_{j=1}^V W_{dj})!}{\prod_{j=1}^V W_{dj}!} \right) + \sum_{j=1}^V W_{dj} \log \beta_{tj} \right] \tag{9}$$

$$- \sum_g \lambda_g \log \lambda_g \tag{10}$$

$$- \log \Gamma\left(\sum_{t=1}^T \gamma_t\right) + \sum_{t=1}^T \log \Gamma(\gamma_t) - \sum_{t=1}^T (\gamma_t - 1) E_q[\log \theta_t] \tag{11}$$

$$- \sum_{d=1}^n \sum_{t=1}^T \phi_{dt} \log(\phi_{dt}) \tag{12}$$

These are nearly the same terms derived by Blei et al. (2003) Appendix A.3, Equation (15). Lines (3), (6), and (7) are the same. Lines (1) and (5) are new for cluster membership. (2) is more complicated because of the summation over the clusters; the \sum_g component is new. (4) is only different because the normalizing constant for the multinomial distribution is added; the post is a draw of size n_d from the relevant topic distribution, rather than each word being a draw of size one. Line (4) is the only change required if stLDA (no clustering) is used rather than traditional LDA.

B.2 Variational Parameter Estimation (λ, γ, ϕ)

Here I derive the ELBO-maximizing values of each variational parameter. Note that, as with vanilla LDA, the update rules are dependent, so the estimation needs to iterate between all of them until convergence. This section relies on computing $E_q[\log \theta_t]$. The value is expressed algebraically in Blei et al. (2003) Appendix A.1.

B.2.1 λ

Recall that λ_g is the variational probability that the user is in cluster g . Maximizing L with respect to λ only involves lines (1), (2), and (5) with the constraint that $\sum_g \lambda_g = 1$. The derivative of L with respect to λ_g is:

$$\log(\xi_g) + f(\alpha_g, \gamma) - \log \lambda_g - 1 + c$$

Where $f(\alpha_g, \gamma)$ is from line (2) above and c is a constant from the Lagrangian. Setting equation equal to zero gives

$$\lambda_g \propto \xi_g \exp(f(\alpha_g, \gamma)) = \xi_g \exp(E_q[\log p(\theta|g, \alpha_g)])$$

This is interpretable as proportional to the ‘‘prior’’ ξ_g times the exponential of the variational expected log likelihood of θ if $\theta \sim Dir(\alpha_g)$. This is interpretable as computing $p(g|\theta)$ as $\propto p(\theta|g)p(g)$ with the likelihood approximated via the variational posterior.

B.2.2 ϕ

Recall that ϕ_{dt} is the variational probability that document d has topic t . Maximizing L with respect to ϕ only involves lines (3), (4), and (7) with the constraint that $\forall d, \sum_t \phi_{dt} = 1$. The terms in L that are functions of ϕ_{dt} are isolated below, with c_d to represent the multinomial normalizing constant.

$$\phi_{dt} E_q[\log \theta_t] + \phi_{dt} \log(c_d) + \sum_{j=1}^V W_{dj} \log \beta_{tj} - \phi_{dt} \log(\phi_{dt})$$

The derivative of L with respect to ϕ_{dt} including Lagrangian term ζ is:

$$E_q[\log \theta_t] + \log(c_d) + \sum_{j=1}^V W_{dj} \log \beta_{tj} - \log(\phi_{dt}) - 1 + \zeta$$

Setting equal to zero gives:

$$\begin{aligned} \log \phi_{dt} &= E_q[\log \theta_t] + \log(c_d) + \log \left(\prod_{j=1}^V \beta_{tj}^{W_{dj}} \right) - 1 + \zeta \implies \\ \phi_{dt} &\propto \exp(E_q[\log \theta_t]) \prod_{j=1}^V \beta_{tj}^{W_{dj}} \end{aligned}$$

This is again interpretable as the likelihood times the prior, but this time the prior is approximated via the variational distribution.

B.2.3 γ

Maximizing L with respect to γ involves lines (2), (3), and (6), noting that $E_q[\log \theta_t] = g_t(\gamma)$, a known formula with well established numerical approximation methods. The components of L involving γ_t are listed below:

$$\sum_g \lambda_g \sum_t (\alpha_{gt} - 1) g_t(\lambda) + \sum_d \sum_t \phi_{dt} g_t(\lambda) - \log \Gamma \left(\sum_{t=1}^T \gamma_t \right) + \log \Gamma(\gamma_t) - \sum_{t=1}^T (\gamma_t - 1) g_t(\lambda)$$

This simplifies to:¹¹

$$\sum_t g_t(\lambda) \left(\sum_g \lambda_g \alpha_{gt} + \sum_d \phi_{dt} - \gamma_t \right) - \log \Gamma \left(\sum_t \gamma_t \right) + \log \Gamma(\gamma_t)$$

Before taking the derivative, it is necessary to expand g_t . $g_t(x) = \Psi(x_t) - \Psi(\sum_t x_t)$ where Ψ is the digamma function (first derivative of the log gamma function). Note that this is the same expression as in Blei et al. (2003) A.3.2 with $\sum_g \lambda_g \alpha_{gt} := \alpha_i$, the sum of each cluster's relevant Dirichlet parameter weighted by (variational) probability of cluster membership λ_g . This will result in the same maximizing value with the substitution because the substitution just changes a constant term. For clarity, however, I will continue to fully derive the result here. For ease of notation, let $\tilde{\alpha}_t = \sum_g \lambda_g \alpha_{gt}$.

Substituting the expression for g_t and $\tilde{\alpha}_t$ gives:

$$\sum_t \left(\tilde{\alpha}_t + \sum_d \phi_{dt} - \gamma_t \right) \left(\Psi(\lambda_t) - \Psi \left(\sum_t \gamma_t \right) \right) - \log \Gamma \left(\sum_t \gamma_t \right) + \log \Gamma(\gamma_t)$$

Taking a derivative with respect to γ_t gives:

$$\begin{aligned} \frac{\partial L}{\partial \gamma_t} &= \left(\tilde{\alpha}_t + \sum_d \phi_{dt} - \gamma_t \right) \Psi'(\gamma_t) - \Psi(\gamma_t) \\ &\quad - \sum_i \left(\tilde{\alpha}_i + \sum_d \phi_{di} - \gamma_i \right) \Psi' \left(\sum_i \gamma_i \right) + \Psi' \left(\sum_i \gamma_i \right) \\ &\quad - \Psi \left(\sum_t \gamma_t \right) + \Psi(\gamma_t) \\ &= \left(\tilde{\alpha}_t + \sum_d \phi_{dt} - \gamma_t \right) \Psi'(\gamma_t) - \sum_i \left(\tilde{\alpha}_i + \sum_d \phi_{di} - \gamma_i \right) \Psi' \left(\sum_i \gamma_i \right) \end{aligned}$$

Setting equal to zero yields a maximum (for all t) of $\gamma_t = \sum_g \lambda_g \alpha_{gt} + \sum_d \phi_{dt}$.

B.3 Model Parameter Estimation (α, β, ξ)

The above section was just for a single user (the analog of a single document in traditional LDA). To estimate the parameters shared across users, cluster-specific Dirichlet parameters α_g , topic distributions β , and cluster proportions ξ , we have to sum the ELBO bound derived above across users. In another abuse of notation return the u do denote individual users.

Let the full likelihood be expressed as:

$$L(\alpha, \beta, \xi) = \sum_u L_u(\lambda_u, \gamma_u, \phi_u; \alpha, \beta, \xi)$$

L_u is the function defined by equations (6) – (12) above with the user-dependence made explicit. This is the function that needs to be maximized with respect to $\{\alpha, \beta, \xi\}$ holding $\{\lambda_u, \gamma_u, \phi_u\}$ constant.

¹¹The -1 parenthetical terms cancel because $\sum_g \lambda_g = 1$

B.3.1 ξ

This is the easiest parameter. It appears only in term (6). Taking a derivative and using a Lagrangian to enforce $\sum_g \xi_g = 1$ gives:

$$\xi_g \propto \sum_u \lambda_{ug}$$

B.3.2 β

This is similar to the variational multinomial in A.4.1 of Blei et al. (2003), with only the added complexity of having multiple words drawn from the same topic.

The full likelihood terms with β plus Lagrangian term is expressed below. $d \in \mathcal{D}_u$ indexes over all documents produced by user u .

$$\sum_u \sum_{d \in \mathcal{D}_u} \sum_t \sum_j \phi_{udt} W_{udj} \log \beta_{tj} - \sum_t \zeta_t \left(1 - \sum_j \beta_{tj} \right)$$

Taking derivatives with respect to β_{tj} results in:

$$\beta_{tj} \propto \sum_u \sum_{d \in \mathcal{D}_u} \phi_{udt} W_{udj}$$

The intuition is that β_t is proportional to the word occurrences weighted by the variational probability that the word was drawn from topic t .

B.3.3 α

These terms are more complicated but can be computed in parallel over clusters. Recall that α is a $G \times T$ matrix of positive values. The ELBO terms involving α are in line (7).

$$\sum_u \sum_g \lambda_{ug} \left(\log \Gamma \left(\sum_{t=1}^T \alpha_{gt} \right) - \sum_{t=1}^T \log \Gamma(\alpha_{gt}) + \sum_{t=1}^T (\alpha_{gt} - 1) g_t(\gamma_u) \right)$$

Let Ψ denote the digamma function, the first derivative of the log gamma function. Taking a derivative with respect to α_{gt} gives:

$$\frac{\partial L}{\partial \alpha_{gt}} = \sum_u \lambda_{ug} \left[\Psi \left(\sum_t \alpha_{gt} \right) - \Psi(\alpha_{gt}) + g_t(\gamma_u) \right] = M_g \left[\Psi \left(\sum_t \alpha_{gt} \right) - \Psi(\alpha_{gt}) \right] + \sum_u \lambda_{ug} g_t(\gamma_u)$$

Where $M_g = \sum_u \lambda_{ug}$. Note that this depends on $\alpha_{gt'}$ for $t' \neq t$ but not on any $\alpha_{g't'}$ where $g' \neq g$. Thus, the update step can be computed simultaneously for each cluster.

This has the same form as the derivative in traditional LDA so the same Newton-Raphson algorithm can be used based on the structure of the Hessian matrix.

$$\frac{\partial L}{\partial \alpha_{gt} \alpha_{gt'}} = -I(t = t') M_g \Psi'(\alpha_{gt}) + M_g \Psi' \left(\sum_t \alpha_{gt} \right)$$

Thus, the Hessian for cluster g can be expressed as: $\mathbf{H}_g = \text{diag}(-M_g \Psi'(\boldsymbol{\alpha}_g)) + \mathbf{1}\mathbf{1}^T M_g \Psi'(\sum_t \alpha_{gt})$, a diagonal matrix plus a constant matrix.