

# Coupling Latent Dirichlet Allocation to multivariate change point time series models to study macroecological time series patterns

Juniper L. Simonis<sup>1,†</sup>, Erica M. Christensen<sup>1,2</sup>, David J. Harris<sup>1,3</sup>, Renata Diaz<sup>1</sup>, Hao Ye<sup>1</sup>, Ethan P. White<sup>1</sup>, and S. K. Morgan Ernest<sup>1</sup>

<sup>1</sup>Weecology Lab, University of Florida

<sup>2</sup>Current institution: the Jornada Rangeland Research Program, New Mexico State University

<sup>3</sup>Current institution: Wayfair, Inc.

<sup>†</sup>Corresponding author: juniper.simonis@weecology.org

## INTRODUCTION

We endeavor to develop methods for analyzing time series of high-dimensional data, and are motivated context by the study of ecological communities comprised of species, where samples of organisms are collected over time (Christensen *et al.* 2018). Specifically, we are interested in determining if the composition of a community (relative composition of species) changes over the course of the study, and if it does, we seek to quantify those dynamics, which may occur abruptly (Williams *et al.* 2011) or smoothly (Tingley *et al.* 2009). However, ecological communities are typically composed of many species relative to the number of samples collected (*i.e.*, the data are high-dimensional; McCune and Grace 2002), which presents a challenge to time series modeling. To address this problem, we reduce the dimensionality of the community data prior to time series analysis (Christensen *et al.* 2018). We accomplish this through a two-stage analysis referred to here as LDATS. The first stage (LDA) uses Latent Dirichlet Allocation (Blei *et al.* 2003) to find the optimally simplified, latent representation of the data, which is then analyzed in the second stage (TS) using Bayesian change point Time Series models (Western and Kleykamp 2004) that we extend for multinomial data using softmax regression (Venables and Ripley 2002). This manuscript describes the two-stage LDATS model in a unified mathematic setting and accompanies an LDATS R package (Simonis *et al. In Development*)

Latent Dirichlet Allocation is a hierarchical Bayesian model that uses a generative classifier (*i.e.*, it uses the joint probability of the inputs and outputs with Bayes' rules to calculate the posterior, as opposed to a discriminative classifier that calculates the posterior directly; see Ng and Jordan 2002) to decompose high-dimension data into a reduced number of latent groups (Blei *et al.* 2003) also known as topics (*i.e.*, LDA is a topic model). The LDA model originally derived and developed by Blei *et al.* (2003) for analysis of textual corpora has since been successfully applied to ecological data (Valle *et al.* 2014, Christensen *et al.* 2018, Valle *et al.* 2018), the domain of interest here. In relation to the original linguistic LDA description and notation (Blei *et al.* 2003), organisms within a sample are like words in a document, species are like terms (word options) in a vocabulary, component communities are like linguistic topics, samples are like documents using the terms, and the whole study is like the corpus of documents (Valle *et al.* 2014). Importantly, LDA is a mixed-membership model, such that terms (species) can be associated with multiple topics (component communities). For the sake of maintaining the relationship between our two-stage LDATS model and the topic model derivation of LDA, we retain the original naming (*e.g.*, observations of words within documents, latent grouping of terms into topics) of Blei *et al.* (2003)

The TS models used here to analyze the decomposed (via LDA) sample data build upon the Bayesian change point model of Western and Kleykamp (2004), which allows for discrete (change point) and continuous temporal changes as well as covariate impacts. This approach allows estimation of the continuous dynamics and covariate impacts to be estimated unconditionally with respect to discrete changes (*i.e.*, the model includes change point uncertainty when estimating regression coefficients; Western and Kleykamp 2004), but had a few components that needed expansion for our application. The original model included a single change point (Western and Kleykamp 2004), but we recognize that ecological communities can undergo multiple discrete shifts during a study (including temporary changes) or may not undergo any abrupt shifts (*i.e.*, change gradually or not at all; Ratajczak, *et al.* 2018). Therefore, we expanded the model to include

potentially 0, as well as multiple (but currently restricted to a maximum of five for purely computational reasons), change points (Ruggieri 2013). Further, the original model assumed a univariate normal response variable (Western and Kleykamp 2004), whereas the output from the LDA in Stage 1 are the proportional parameters of a multinomial response, which are multivariate and non-normal. For use in LDATS, we therefore generalized (à la generalized linear modeling; McCullagh and Nelder 1989) the model of Western and Kleykamp (2004) using softmax regression (Venables and Ripley 2002) to predict multinomial probability variables. The TS models are fit with Bayesian techniques using parallel tempering Markov Chain Monte Carlo (ptMCMC) methods (Earl and Deem 2005) to locate the change points and neural networks (Ripley 1996) to estimate continuous time and covariate parameters.

By combining dimension reduction and time series analysis into a single mathematical framework and software pipeline, LDATS provides a robust and user-friendly methodology for evaluating complex dynamics of high-dimension timeseries (Fig. 1). Although we understand the pressing need for these models within ecology (Andersen *et al.* 2009, Karssenberg *et al.* 2017, Ratajczak, *et al.* 2018), we recognize their application may be broader, given the interest in regime shifts in financial, political, and engineering sectors, for example (Scheffer 2009, Gal and Anderson 2010). We therefore keep the model and its coding implementation as general as possible to facilitate application to other systems of interest. This manuscript details the mathematical model underlying the LDATS methodology and an example application from the motivating system (based on Christensen *et al.* 2018). The accompanying LDATS software package implements the model in R (R Core Team 2018). The code is currently available via the Weecology Lab GitHub repository (<https://github.com/weecology/LDATS>).

## MODEL DEVELOPMENT

### *General Terminology and Notation*

Because of the overlap in notation between LDA (Blei *et al.* 2013), the time series models used here (Western and Kleykamp 2004), and the ptMCMC method used to fit the time series models (Earl and Deem 2005) (*e.g.*, all use  $\beta$  but with different meanings), we create a notational set for use here with an effort to minimize name reuse (Table 1). Given that LDATS specifically uses LDA as the first stage and that our methods build upon topic models, in instances of notational overlap between the LDA and TS or ptMCMC components, we defer to the LDA usage. We do make one important deviation from the original LDA notation (Blei *et al.* 2003), however, to clarify the dimensionality of state variables and parameters. Specifically, we use the lowercase, regular type letter (*e.g.*,  $\beta$ ) to indicate a singular value; the lowercase, boldface letter (*e.g.*,  $\boldsymbol{\beta}$ ) to indicate a vector of values; and the capital, boldface letter (*e.g.*,  $\mathbf{B}$ ) to indicate a matrix of values.

A corpus (set of documents)  $D$  consists of  $M$  total documents comprising  $N$  total words from  $V$  total terms. Each document  $d$  (in  $1 \dots M$ ) consists of  $N_d$  words ( $n$  in  $1 \dots N_d$ ) assigned to one of  $v$  in  $1 \dots V$  terms. The total number of words in the corpus is the sum of the words within each document. The weight of document  $u_d$  is the number of words it has relative to the maximum number of words in any document:

$$u_d = \frac{N_d}{\arg \max_{1 \leq j \leq M} (N_j)} \quad [1]$$

allowing us to account for variable numbers of words among documents (with vector  $\mathbf{u}_{1 \dots M}$ , or just  $\mathbf{u}$ ).

LDA involves grouping the terms into  $i$  in  $1 \dots k$  total latent (unobserved) component topics, where “component topic” means a group of terms that tend to be found together in specific proportions. The allocation process (Blei *et al.* 2003) allows individual terms to be assigned to multiple component topics. The total number of latent topics is also unknown, and for the present approach is fixed *a priori* within a given

Stage 1 (LDA) model  $m_1$  in  $1 \dots \mathcal{M}_1$  (note the difference between  $\mathcal{M}_1$ , the number of Stage 1 models and  $M$ , the number of documents) as  $k_{m_1}$ .

Each word within a document has an observed term identity  $w_{d,n_d}$  and a latent topic membership that is fit in the Stage 1 model  $m_1$  as  $z_{m_1,d,n_d}$ . Because there are varying numbers of words in each document, we use a vector structure to hold word-level data across the corpus. The term identities of all words within document  $d$  are  $w_{d,1 \dots N_d}$  (or  $w_d$ ) and the term identities of all words across all documents  $M$  are  $w_{1 \dots M, 1 \dots N_d}$  (or  $w$ ), an  $N$ -length vector. Similarly, the topic identities of all words in document  $d$  are  $z_{m_1,d,1 \dots N_d}$  (or  $z_{m_1,d}$ ) and the topic identities of all words across all documents  $M$  are  $z_{m_1,1 \dots M, 1 \dots N_d}$  (or  $z_{m_1}$ ), an  $N$ -length vector. Thus  $z_{m_1}$  contains the topic identity (latent state) and  $w$  the term identity (observed state) for all words in the corpus. Note that the term identity does not contain a rank-index subscript associated with the model ( $m_1$ ), indicating that the term identity stays the same for all models, whereas the topic identity does have the  $m_1$  subscript, indicating that it varies among the  $\mathcal{M}_1$  Stage 1 models.

We are interested in temporal changes in topic composition, and so define the time of document  $d$  to be  $t_d$  and the vector of all document times to be  $t_{1 \dots M}$  or simply  $t$ .  $t$  defines the temporal relationship among documents, is used to locate change points, and presently must be a discrete (or discretizable) variable. For a given time series model  $m_2$  (in  $1 \dots \mathcal{M}_2$ ), we also collate  $C_{m_2}$  total covariates (including an overall intercept), indexed as  $c$  in  $1 \dots C_{m_2}$ , and measured for each document. The value of a particular covariate  $c$  for a specific document  $d$  is  $x_{m_2,d,c}$  and the set of  $C_{m_2}$  covariates for the document is a vector  $x_{m_2,d,1 \dots C_{m_2}}$  or simply  $x_{m_2,d}$ . All of the covariates (including the intercept) across all of the documents are held in  $X_{m_2}$ , an  $M \times (C_{m_2} + 1)$  matrix, that can vary among the  $\mathcal{M}_2$  time series models (hence the rank-index subscript).

#### Stage One: Dimension Reduction

##### Latent Dirichlet Allocation: Single Model

The first stage of the LDATS analysis reduces the raw, high-dimensional data (counts of terms in documents over time) to a lower dimensional representation of the information contained in the data using Latent Dirichlet Allocation (LDA; Blei *et al.* 2003). Specifically, we use the Variational Expectation Maximization (Jordan *et al.* 1999) version of the LDA model derived and developed first by Blei *et al.* (2003). For a Stage 1 model  $m_1$  with a total number of topics  $k_{m_1}$ , the distribution of topics within a document  $d$  is a  $k_{m_1}$ -dimension categorical random variable described by probabilities  $\theta_{m_1,d,1} \dots \theta_{m_1,d,k_{m_1}}$  held in the vector  $\theta_{m_1,d}$  ( $\sum \theta_{m_1,d} = 1$ ) and collated across documents into the  $M \times k_{m_1}$  matrix  $\Theta_{m_1}$ . Thus, the realized topic identity ( $z$ ) of word  $n$  within document  $d$  under model  $m_1$  is

$$z_{m_1,d,n} \sim \text{Cat}_{k_{m_1}}(\theta_{m_1,d}) \quad [2]$$

The vector of topic probabilities within a document ( $\theta_{m_1,d}$ ) is defined by a  $k_{m_1}$ -dimensional Dirichlet distribution with concentration parameters  $\alpha_{m_1,d}$ , which (following Blei *et al.* 2003) we assume do not change among documents (*i.e.*,  $\alpha_{m_1,1} = \dots = \alpha_{m_1,M} = \alpha_{m_1}$ ) and are symmetric (*i.e.*,  $\alpha_{m_1,d,1} = \dots = \alpha_{m_1,d,k_{m_1}} = \alpha_{m_1}$ ), reducing the set to a single unknown parameter  $\alpha_{m_1}$  for the model. Thus, the topic probabilities are

$$\theta_{m_1,d} \sim \text{Dir}_k(\alpha_{m_1}) \quad [3]$$

which is drawn separately for each document (Blei *et al.* 2003). The word-level term distribution within a document is a  $V$ -dimension categorical random variable contingent upon the topic identity of the word and

defined by probabilities  $\beta_{m_1,d,1,1} \dots \beta_{m_1,d,k_{m_1},V}$ , where  $\sum_v \beta_{m_1,d,i} = 1$ . The probabilities across all topics within a document are held in a  $k_{m_1} \times V$  matrix ( $\mathbf{B}_{m_1,d}$ ), which we assume is constant across documents, ( $\mathbf{B}_{m_1,d} = \dots = \mathbf{B}_{m_1,M} = \mathbf{B}_{m_1}$ ; Blei *et al.* 2003). The word-level term identity ( $\mathbf{w}$ ) is then generally defined:

$$w_{d,n} \sim \text{Cat}_V(z_{m_1,d,n}, \mathbf{B}_{m_1}) \quad [4]$$

$\mathbf{w}$  is therefore a function of unknown parameters  $\alpha_{m_1}$  (a scalar) and  $\mathbf{B}_{m_1}$  (a  $k_{m_1} \times V$  matrix).

Despite the fact that the observations ( $\mathbf{w}$ ; word-level term identities) are a statistical function of the unknown parameters  $\alpha_{m_1}$  and  $\mathbf{B}_{m_1}$ , we are actually interested in the latent components of the model, not the base parameters. Specifically, we would like to estimate the posterior probability distribution for the latent topic probabilities  $\boldsymbol{\theta}_{m_1}$  (and thus also states  $\mathbf{z}_{m_1}$ ) given the observations  $\mathbf{w}$ , but that distribution depends on the parameters fit by the model, and so our inferential problem becomes determining the posterior probability distribution of the latent components, given the observations and the estimated parameters:

$$\mathcal{P}(\boldsymbol{\theta}_{m_1}, \mathbf{z}_{m_1} | \mathbf{w}, \alpha_{m_1}, \mathbf{B}_{m_1}) = \frac{\mathcal{P}(\boldsymbol{\theta}_{m_1}, \mathbf{z}_{m_1}, \mathbf{w} | \alpha_{m_1}, \mathbf{B}_{m_1})}{\mathcal{P}(\mathbf{w} | \alpha_{m_1}, \mathbf{B}_{m_1})} \quad [5]$$

where  $\mathcal{P}$  is used generally as a probability distribution function. Clearly, estimating the probability distribution of the latent states necessitates calculation of  $\alpha_{m_1}$  and  $\mathbf{B}_{m_1}$ , which facilitates expansion of the probability of observations  $\mathbf{w}$  given the parameters  $\alpha_{m_1}$  and  $\mathbf{B}_{m_1}$  ( $\mathcal{P}(\mathbf{w} | \alpha_{m_1}, \mathbf{B}_{m_1})$ ; Appendix 1):

$$\mathcal{P}(\mathbf{w} | \alpha_{m_1}, \mathbf{B}_{m_1}) = \prod_{d=1}^M \left[ \int \mathcal{P}(\boldsymbol{\theta}_{m_1,d} | \alpha_{m_1}) \left( \prod_{n=1}^{N_d} \sum_{z_{m_1,d,n}} \mathcal{P}(w_{d,n} | z_{m_1,d,n}, \mathbf{B}_{m_1}) \mathcal{P}(z_{m_1,d,n} | \boldsymbol{\theta}_{m_1,d}) \right) d\boldsymbol{\theta}_{m_1,d} \right] \quad [6]$$

This equation highlights the problematic coupling of  $\boldsymbol{\theta}_{m_1,d}$  (and thus  $\alpha_{m_1}$ ) and  $\mathbf{B}_{m_1}$  in the summation over latent topics (Blei *et al.* 2003), which prevents direct, tractable estimation of parameters (and latent states).

To circumvent this issue, we use a variational approximation (Jordan *et al.* 1999) to the equations that decouples the parameters, and which we fit using the expectation-maximization routine (aka VEM for Variational Expectation Maximization; Blei *et al.* 2003; Appendix 2). To accomplish this, we endow the model with free latent variational parameters  $\boldsymbol{\Gamma}_{m_1}$  and  $\boldsymbol{\Phi}_{m_1}$  (Appendix 2) that decouple the terms and characterize a family of distributions ( $\mathcal{Q}$  to distinguish from  $\mathcal{P}$ ) providing a lower bound on the probabilities (Jordan *et al.* 1999, Blei *et al.* 2003). Once the VEM algorithm has converged, we achieve approximate maximum likelihood estimates for the model parameters ( $\alpha_{m_1}^*$  and  $\mathbf{B}_{m_1}^*$ ) given the full set of observations ( $\mathbf{w}$ ) for model  $m_1$ . This estimation procedure is executed using the LDA function in the topicmodels package (v0.2-7; Grün and Hornik 2011) in R (v 3.5.1; R Core Team 2018), which leverages C code written by Blei *et al.* (2003).

#### Latent Dirichlet Allocation: Multi-Model Inference

Given the fit of a specific Stage 1 model ( $m_1$ ), we can then consider multiple Stage 1 models to determine the model with the most parsimonious number of topics  $k_{m_1}^*$ . Specifically, we use AIC as our Stage 1 model selection criterion (Christensen *et al.* 2018), defined for a specific LDA model  $m_1$ :

$$\text{AIC}_{m_1} = -2\ell(\alpha_{m_1}^*, \mathbf{B}_{m_1}^* | \mathbf{w}) + 2\ell_{m_1} \quad [7]$$

where  $\ell$  is the log likelihood ( $\ell(\alpha_{m_1}^*, \mathbf{B}_{m_1}^* | \mathbf{w}) = \log \mathcal{P}(\mathbf{w} | \alpha_{m_1}^*, \mathbf{B}_{m_1}^*)$ ) and  $\ell_{m_1} = (1 + k_{m_1}V)$  is the

number of parameters in the model: 1 for  $\alpha_{m_1}^*$  and  $k_{m_1}V$  corresponding to each entry in  $\mathbf{B}_{m_1}^*$ , a  $k_{m_1} \times V$  matrix (Grün and Hornik 2011). If small sample size is a concern with respect to the degrees of freedom being consumed by the model, one can use the AICc correction based on the number of observations, here corpus size ( $M$ ; Grün and Hornik 2011):

$$\text{AICc}_{m_1} = -\text{AIC}_{m_1} + \frac{2k_{m_1}^2 + 2k_{m_1}}{M - k_{m_1} - 1} \quad [8]$$

Because of the use of multiple iterative optimization routines (which require starting values to be drawn at random) to estimate otherwise intractable probability distributions, it is critical to account for the potential influence of starting values on analytical results. Here, we accomplish this by running multiple models with the same number of topics ( $k_{m_1}$ ) using different starting values, assigned through the random number generator seed ( $\mathfrak{J}$ ). Specifically, we use  $\mathcal{N}$  replicates ( $\mathfrak{n}$  in  $1 \dots \mathcal{N}$ ) at each number of topics from 2 to  $k_{\max}$ , the total number of topics to be explored. The minimal number of topics is set to 2 by the current coding implementation of the LDA algorithm (Blei *et al.* 2003, Grün and Hornik 2011), although the underlying mathematics can include the limiting case of a single topic (*i.e.*, no dimension reduction). Thus, the total number of models in Stage 1 ( $\mathcal{M}_1$ ) is

$$\mathcal{M}_1 = \mathcal{N}(k_{\max} - 1) \quad [9]$$

The optimal (according to AIC) LDA model ( $m_1^*$ ) is determined by

$$m_1^* = \arg \min_{m_1 \in 1 \dots \mathcal{M}_1} \text{AIC}(m_1) \quad [10]$$

and has the corresponding set of parameters  $\{\alpha_{m_1^*}^*, \mathbf{B}_{m_1^*}^*, k_{m_1^*}, \mathfrak{J}_{m_1^*}\}$ .

#### Reduced-Dimension Data

Having found the optimal Stage 1 model, we obtain the posterior estimates for the document-level topic probabilities (held in a  $k_{m_1^*}$ -length vector), which we will use as the response in our Stage 2 models. Because we are using a variational approach, we actually obtain posterior point estimates for the variational parameters  $\mathbf{\Gamma}_{m_1^*}$  (rather than the base model parameters  $\mathbf{\Theta}_{m_1^*}$ ; Grün and Hornik 2011), taken from the final step of the VEM algorithm in model  $m_1^*$ . Recognizing that  $\mathbf{\Gamma}_{m_1^*}$  contains Dirichlet concentration parameters within documents ( $\mathbf{\gamma}_{m_1^*,d}^*$ ), we must normalize the values so they are proper proportions (sum to one) and can be modeled using a multinomial distribution. We notate the normalized parameters with the overbar accent as  $\bar{\mathbf{\gamma}}_{m_1^*,d}^*$ :

$$\bar{\mathbf{\gamma}}_{m_1^*,d}^* = \frac{\mathbf{\gamma}_{m_1^*,d}^*}{\sum \mathbf{\gamma}_{m_1^*,d}^*} \quad [11]$$

The normalized posterior point estimates of the topic proportions across all of the  $M$  documents in the corpus are held in the  $M \times k_{m_1^*}$  matrix  $\bar{\mathbf{\Gamma}}_{m_1^*}$ , which corresponds to the optimal (according to AIC based on VEM inference) decomposition of the word-level data to topic-level data. This matrix forms the multivariate response variable analyzed in the time series model, as outlined in the next section.

#### *Stage Two: Multinomial Time Series*

The second stage of the LDATS model analyzes the time series of topic proportions estimated by the LDA ( $\bar{\Gamma}_{m_1}^*$ ) to quantify temporal dynamics and in particular to identify abrupt change point. The times of the documents ( $\mathbf{t} = t_1 \dots t_M$ ) have the potential to influence topic proportions in multiple manners: the time may control the application of the predictor variables in the model (in the case of discrete change points), may directly influence quantitative values of predictors (if continuous time impacts are included in the regression model), or may not impact  $\bar{\Gamma}_{m_1}^*$  at all (in the case of a model with no change points and no continuous time impacts). Presently, temporal autocorrelation is not included in the time series models, but is planned for future work (see **FUTURE DEVELOPMENTS**). We base our model on that indicator regression approach of Western and Kleykamp (2004), but make notable alterations, namely allowing for multiple change points, fitting segment-level regressions individually (following Ruggeri 2013), and generalizing the segment-level regression to model multinomial response data.

A Stage 2 model  $m_2$  has a non-negative integer number of discrete change points ( $P_{m_2}$ ) that divide the corpus into distinct temporal segments ( $S$  in  $1 \dots S_{m_2}$ , aka “chunks”) such that the number of chunks is always one more than the number of change points ( $S_{m_2} = P_{m_2} + 1$ ; see [Change Points: Segmenting the Time Series](#)). If there are change points (*i.e.*,  $P_{m_2} > 0$ ), then their locations (for the  $p$  in  $1 \dots P_{m_2}$  change points) are unknown parameters to be estimated (Western and Kleykamp 2004, Christensen *et al.* 2018). A specific change point  $p$ ’s location is represented by  $\rho_{m_2,p}$  and the set of change point locations is the  $P_{m_2}$  – length vector  $\rho_{m_2,1 \dots P}$  (or  $\rho_{m_2}$ ). In addition to the number of change points,  $m_2$  has a within-chunk regression model defined by a set of covariates (including the intercept indicator, temporal, and non-temporal covariates)  $\mathbf{X}_{m_2}$  that impact the topic proportions through a set of parameters  $\hat{\mathbf{H}}_{m_2}$  à la generalized linear models (McCullagh and Nelder 1989; see [Segment-Level Models: Multinomial Logistic Regression](#)). In combination then, we seek to estimate the posterior probability distribution of the change point locations and the regression coefficients given the fitted topic probabilities from  $m_1^*$  ( $\mathcal{P}(\rho_{m_2}, \hat{\mathbf{H}}_{m_2} | \bar{\Gamma}_{m_1}^*)$ ).

We estimate  $\mathcal{P}(\rho_{m_2}, \hat{\mathbf{H}}_{m_2} | \bar{\Gamma}_{m_1}^*)$  by leveraging the dual nature of the change points as both parameters of the model and parameters that structure the model ( $\hat{\mathbf{H}}_{m_2}$ ) and calculating the marginal posterior probabilities of  $\rho_{m_2}$  and  $\hat{\mathbf{H}}_{m_2}$  in sequence (Western and Kleykamp 2004). First, we estimate the marginal posterior probability distribution of the change point locations, given the topic proportions

$$\mathcal{P}(\rho_{m_2} | \bar{\Gamma}_{m_1}^*) = \int_{\hat{\mathbf{H}}_{m_2}} \mathcal{P}(\rho_{m_2}, \hat{\mathbf{H}}_{m_2} | \bar{\Gamma}_{m_1}^*) d\hat{\mathbf{H}}_{m_2} \quad [12]$$

using Markov Chain Monte Carlo (MCMC; see [Combining Segment-Level Models: Inference About Change Point Locations](#)). Then, we define model realizations ( $r$  in  $1 \dots R_{m_2}$ , as in  $m_2^r$  in  $m_2^1 \dots m_2^{R_{m_2}}$ ) to specify the actual change point locations ( $\rho_{m_2}^r$ ; see [Change Points: Segmenting the Time Series](#)). In a corpus of  $M$  documents with a range of time stamps from  $t_1$  to  $t_M$ , the  $P_{m_2}$  change points have  $\binom{(t_M-1)-t_1}{P_{m_2}}$  (“ $(t_M - 1) - t_1$  choose  $P_{m_2}$ ”) possible values for  $\rho_{m_2}$ , constituting  $R_{m_2}$  unique realizations of  $m_2$ . Each realization has a posterior probability value itself ( $\mathcal{P}(\rho_{m_2} = \rho_{m_2}^r | \bar{\Gamma}_{m_1}^*)$ ) and produces a conditional posterior probability distribution of the regressors given the change point locations it specifies ( $\mathcal{P}(\hat{\mathbf{H}}_{m_2} | \rho_{m_2}^r, \bar{\Gamma}_{m_1}^*)$ ). We then combine these conditional distributions via Bayesian Model Averaging (BMA; Bartels 1997, Western and Kleykamp 2004, Hobbs and Hooten 2015) considering the realizations as submodels and using their posterior probabilities as weights to estimate the marginal posterior probability of the regression coefficients, given just the topic proportions:

$$\mathcal{P}(\mathbf{H}_{m_2} | \bar{\mathbf{F}}_{m_1}^*) = \sum_{r=1}^{R_{m_2}} \mathcal{P}(\mathbf{H}_{m_2} | \boldsymbol{\rho}_{m_2}^r, \bar{\mathbf{F}}_{m_1}^*) \mathcal{P}(\boldsymbol{\rho}_{m_2} = \boldsymbol{\rho}_{m_2}^r | \bar{\mathbf{F}}_{m_1}^*) \quad [13]$$

(see Combining Segment-Level Models: Inference About Within-Segment Parameters).

### Change Points: Segmenting the Time Series

To deconstruct the time series into chunks, we augment the vector of change point locations  $\boldsymbol{\rho}_{m_2}$  with the time step before the minimum ( $\min \mathbf{t} - 1$ ) and the maximum time step ( $\max \mathbf{t}$ ), generating the  $(P_{m_2} + 2)$ -length vector  $\tilde{\boldsymbol{\rho}}_{m_2}$  where the overbrace indicates the addition of the fixed range. In the instance that there are no change points (*i.e.*,  $P_{m_2} = 0$ ),  $\tilde{\boldsymbol{\rho}}_{m_2}$  is still defined, but now is simply a length-2 vector including the minimum and maximum times, and therefore includes no unknown change point locations to be estimated. We assign the documents into segments via a mapping function ( $f_{d \rightarrow s}$ ), which returns an indication ( $\xi_{m_2, d, s}$ ) of whether or not document  $d$  belongs to chunk  $s$  (0 for no or 1 for yes) based on its timestamp  $t_d$  and the start ( $\tilde{\rho}_{m_2, s, 1}$ , the first time step after the previous change point) and end ( $\tilde{\rho}_{m_2, s, 2}$ , the timestep of the change point) times of the chunk:

$$\xi_{m_2, d, s} = f_{d \rightarrow s}(\tilde{\boldsymbol{\rho}}_{m_2}, t_d, s) = \begin{cases} 0, & t_d < \tilde{\rho}_{m_2, s, 1} \\ 1, & \tilde{\rho}_{m_2, s, 1} \leq t_d \leq \tilde{\rho}_{m_2, s, 2} \\ 0, & t_d > \tilde{\rho}_{m_2, s, 2} \end{cases} \quad [14]$$

For each chunk of time,  $f_{d \rightarrow s}$  produces a length- $M$  vector of 0s and 1s ( $\xi_{m_2, 1 \dots M, s}$  or  $\xi_{m_2, s}$ ), which are collated across chunks into an  $M \times S_{m_2}$  matrix ( $\Xi_{m_2}$ ) that identifies to which chunk each document belongs (a document only belongs to one segment, such that the columns of  $\Xi_{m_2}$  each sum to 1).  $\Xi_{m_2}$  deconstructs the Stage 1 output ( $\bar{\mathbf{F}}_{m_1}^*$ ) into  $S_{m_2}$  submatrices ( $\bar{\mathbf{F}}_{m_1, s}^*$  in  $\bar{\mathbf{F}}_{m_1, 1}^* \dots \bar{\mathbf{F}}_{m_1, S_{m_2}}^*$ ), corresponding to the chunks.

### Segment-Level Models: Multinomial Logistic Regression

In LDATS, the indicator function is used to segment the data prior to analyses, such that under a given model  $m_2$ , each of the chunks of documents is fit with a separate version of the same regression (*i.e.*, the regression models all include the same predictor variables). Our approach corresponds to the assumption of no covariance among parameters across segments, which is reasonable given our conceptualization of the change points as discrete and abrupt. In this, LDATS follows Ruggieri (2013) but deviates from the approach of Western and Kleykamp (2004), who fit a single regression across the (two) chunks that applied the indicator internally, thereby allowing non-zero covariance among parameters between the segments. However, Western and Kleykamp (2004) had a singular change point, modeled a normal response variable, and used a relatively simple linear regression model with few covariates, all assumptions that we relax in LDATS (we allow for multiple change points, have a multinomial response variable, and permit complex predictive models) leading to a substantial increase in the number of parameters fit by the model, which would be computationally prohibitive to fit under a single model allowing for full covariance due to the size of the variance-covariance matrix (Genz 1992). (Although see **FUTURE DEVELOPMENTS**.)

The within-chunk component of a Stage 2 model  $m_2$  predicts the matrix of topic proportions for the  $M_{m_2, s}$  documents belonging to segment  $s$  ( $\bar{\mathbf{F}}_{m_1, s}^*$ ) in terms of  $C_{m_2}$  predictors ( $\mathbf{X}_{m_2, s}$ ) and  $C_{m_2} k_{m_1}^*$  coefficients ( $\mathbf{H}_{m_2, s}$ ), recognizing that the segmentation (*i.e.*, what  $s$  indexes over) depends on the identification matrix  $\Xi_{m_2}$ , which varies among realizations of model  $m_2$ . Although the original change point

model (Western and Kleykamp 2004) assumed a univariate normal response variable, our response data are multivariate and non-normal. Specifically, our response variable is a set of  $k_{m_1}^*$  multinomial probabilities, each of which must be non-negative and which must sum to 1 within a document. We address these constraints on the response variable by taking a generalized linear model approach (McCullagh and Nelder 1989) and modeling the data using a log-linear multinomial (aka multinomial logit or softmax) regression (Ripley 1996, Venables and Ripley 2002) based on a set of augmented parameters  $\hat{\mathbf{H}}_{m_2,s}^r$  (the acute accent indicates augmentation) that define the first topic as a reference value (Appendix 3; Venables and Ripley 2002):

$$\mathbb{E}[\bar{\mathbf{T}}_{m_1,s}^*]_{m_2^r} = \text{softmax}(\mathbf{X}_{m_2,s} \hat{\mathbf{H}}_{m_2,s}^r) \quad [15]$$

where  $\mathbb{E}$  indicates the expected (predicted) value(s) of the proportion(s) and **softmax** is the normalized exponential function that generalizes the logistic function to multiple dimensions (Bishop 2006). This representation is aligned with the generalized linear model equation (McCullagh and Nelder 1989), wherein our link and inverse link functions are the multinomial logit and softmax (akin to the binomial logit and logistic functions for a logistic regression). Recognizing the uncertainty in the relationship between the “observations” ( $\bar{\mathbf{T}}_{m_1,s}^*$ ) and predictions ( $\mathbb{E}[\bar{\mathbf{T}}_{m_1,s}^*]_{m_2^r}$ ) allows us to estimate the probability distribution of the regression parameters for the segment given the change point locations and the segment’s topic proportions ( $\mathcal{P}(\hat{\mathbf{H}}_{m_2,s}^r | \boldsymbol{\rho}_{m_2,s}^r, \bar{\mathbf{T}}_{m_1,s}^*)$ ) and measure the segment-level contribution to the probability of the realization of change point locations (used to calculate  $\mathcal{P}(\boldsymbol{\rho}_{m_2} = \boldsymbol{\rho}_{m_2}^r | \bar{\mathbf{T}}_{m_1}^*)$ ; see [Combining Segment-Level Models: Inference About Within-Segment Parameters](#)). This flexible formulation of the regression model is a substantial generalization of the original LDATS application, which imposed a seasonal (within-year) dynamic modeled as a Fourier series that made biological sense for the system of focus (Christensen *et al.* 2018). However, many systems experience dynamics that are not seasonal (*e.g.*, decadal cycles) or are not stationary (*i.e.*, directional changes), and so generalizing the regression model is a key component to generalizing LDATS more broadly.

Following Bayes’ rule, the posterior probability distribution of the regression parameters for a segment is proportional to the probability distribution of the topic proportions given the parameters (*i.e.*, the model likelihood) times the prior probability distribution for the parameters:

$$\mathcal{P}(\hat{\mathbf{H}}_{m_2,s}^r | \boldsymbol{\rho}_{m_2,s}^r, \bar{\mathbf{T}}_{m_1,s}^*) \propto \mathcal{P}(\bar{\mathbf{T}}_{m_1,s}^* | \boldsymbol{\rho}_{m_2,s}^r, \hat{\mathbf{H}}_{m_2,s}^r) \mathcal{P}(\hat{\mathbf{H}}_{m_2,s}^r | \boldsymbol{\rho}_{m_2,s}^r) \quad [16]$$

where all of the probabilities are still conditional on the realized change point locations (*i.e.*, the submodel). The probability of a single document  $d_s$  given the parameters is the product of each topic’s predicted probability raised to power of the corresponding observed probability, which is equivalent to the negative cross entropy between the observed and predicted distributions (Berger *et al.* 1996, Malouf 2002). Of particular importance here is that the “observed” values (which are actually estimated latent topic proportions) are probabilities, and so we retain the full version of the probability equation (compared to when the observed data are categorized individuals and the equation can be simplified due to the fact that all states are 0 except one that is 1 for each observation). The probability of the  $M_{m_2,s}^r$  documents in the segment is then the weighted (by  $\mathbf{u}$ ) product of the document probabilities:

$$\mathcal{P}(\bar{\mathbf{T}}_{m_1,s}^* | \boldsymbol{\rho}_{m_2,s}^r, \hat{\mathbf{H}}_{m_2,s}^r) = \prod_{d_s=1}^{M_{m_2,s}^r} e^{u_{d_s}} \prod_{i=1}^{k_{m_1}^*} \left( \mathbb{E}[\bar{\mathbf{y}}_{m_1,d_s,i}^*]_{m_2^r} \right)^{\bar{\mathbf{y}}_{m_1,d_s,i}^*} \quad [17]$$

The weights are exponentiated on the probability scale and thus linear on the additive log-probability (*i.e.*, log-likelihood, loss) scale.



Generally speaking, we use a multivariate Gaussian prior for the regression parameters, following standard Bayesian approaches (Gelman and Hill 2007). The current iteration of LDATS assumes a mean-0 distribution with common precision (inverse variance) across segments  $\lambda$  and no covariance (but see **FUTURE DEVELOPMENTS**). Recalling that the model parameters  $\hat{\mathbf{H}}_{m_2,s}^r$  are defined with the first topic as a reference (intercept), however, there are some entries in  $\hat{\mathbf{H}}_{m_2,s}^r$  that are fixed at 0 and not free for estimation. Specifically, for any covariate  $c$ , the parameter associated with the first topic  $i = 1$  is 0:  $\hat{\eta}_{m_2,s,i=1,c} = 0$ . Thus, we define the variance-covariance matrix using the augmented identity matrix  $\hat{\mathbf{I}}$  where the accent matches  $\hat{\mathbf{H}}$  and signifies that all of the entries associated with the first topic are set to 0. Including the set-to-0 regressors, there are  $C_{m_2}k_{m_1}^*$  entries in  $\hat{\mathbf{H}}_{m_2,s}^r$ , and thus the multivariate normal distribution is of dimension  $C_{m_2}k_{m_1}^*$ :

$$\hat{\mathbf{H}}_{m_2,s}^r \sim \text{MVN}_{C_{m_2}k_{m_1}^*}(\mathbf{0}, \lambda \hat{\mathbf{I}}) \quad [18]$$

where  $\text{MVN}$  is defined by the mean vector ( $\mathbf{0}$ ) and precision matrix ( $\lambda \hat{\mathbf{I}}$ ). The original formulation of LDATS (Christensen *et al.* 2018) assumed  $\lambda = \mathbf{0}$  (*i.e.*, a fully vague prior) but including a small increase in precision ( $\lambda \approx 10^{-4} - 10^{-2}$ ) can aid in finding the optimal solution as long as all coefficients have been scaled to about  $[0,1]$  (*e.g.*, normalized; Ripley 1993, Ripley 1996). This approach to estimating  $\mathcal{P}(\hat{\mathbf{H}}_{m_2,s}^r | \boldsymbol{\rho}_{m_2,s}^r, \bar{\mathbf{T}}_{m_1,s}^*)$  is equivalent to the method called ridge regression (Hoerl 1962, Hoerl and Kennard 1970), Tikhonov regularization (Tikhonov and Arsenin 1977), L2 regularization (Ng 2004) or joint maximum *a posteriori* (MAP; Bassett and Deride 2018) estimation of a neural network using a weight decay (Venables and Ripley 2002), depending on the applied setting.

Computationally, the regression for a chunk of time based on a realization of a Stage 2 model is fit using the multinom and nnet functions within the nnet package (v7.3-12; Venables and Ripley 2002) in R (R Core Team 2018), which formulate the regression as a single-hidden-layer neural network with skip-layer connections (Venables and Ripley 2002). The posterior distribution for the parameters is found with the gradient-based optimization routine known as the Broyden–Fletcher–Goldfarb–Shanno or BFGS Algorithm (Brayden 1970, Fletcher 1970, Goldfarb 1970, Shanno 1970), a quasi-Newtonian iterative searching method for non-linear models. Here, we use the negative log of the prior-penalized likelihood as the loss value ( $\mathcal{L}$ ):

$$\mathcal{L}_{m_2,s}^r = -\log\left(\mathcal{P}(\bar{\mathbf{T}}_{m_1,s}^* | \boldsymbol{\rho}_{m_2,s}^r, \hat{\mathbf{H}}_{m_2,s}^r) \mathcal{P}(\hat{\mathbf{H}}_{m_2,s}^r | \boldsymbol{\rho}_{m_2,s}^r)\right) \quad [19]$$

that is minimized to find the optimal parameter value set  $\hat{\mathbf{H}}_{m_2,s}^{*r}$

$$\hat{\mathbf{H}}_{m_2,s}^{*r} = \arg \min_{\hat{\mathbf{H}}_{m_2,s}^r} \mathcal{L}_{m_2,s}^r \quad [20]$$

which has the loss value  $\mathcal{L}_{m_2,s}^{*r}$  and corresponds to the mode of the posterior (Venables and Ripley 2002). The BFGS Algorithm works efficiently by not calculating the Hessian (matrix of partial second derivatives) of the prior-penalized loss function at every step in the optimization, but rather approximating it by comparing successive iterations of the Jacobian matrix (matrix of partial first derivatives), whose components correspond to the partial derivatives the loss with respect to the coefficients associated with each combination of covariate and topic. The full Jacobian of the loss equation has an extensive derivation (due to the multivariate application of the chain rule) that is based on the nuances of the data set being analyzed, but which collapses neatly because it is sparse (Appendix 4).

#### Combining Segment-Level Models: Inference About Change Point Locations

We then define the full time series model for the realization by collating the chunk-level models. Specifically, we row-wise stack the chunk-specific parameter matrices  $\hat{\mathbf{H}}_{m_2,1}^r$  to  $\hat{\mathbf{H}}_{m_2,S_{m_2}}^r$  creating  $\hat{\mathbf{H}}_{m_2}^r$  (dimension:  $(S_{m_2} C_{m_2}) \times k_{m_1}^*$ ) and we place the chunk-specific covariate matrices  $\mathbf{X}_{m_2,1}^r$  to  $\mathbf{X}_{m_2,S_{m_2}}^r$  into a diagonal square  $(S_{m_2} \times S_{m_2})$  matrix to produce a covariate matrix for the full time series  $\mathbf{X}_{m_2}^r$  (dimension:  $M \times (S_{m_2} C_{m_2})$ ). Recalling that the splitting of the corpus into segments and thus the fitting of chunk-level parameters is governed by the realized change point locations  $\boldsymbol{\rho}_{m_2}^r$  via the indication matrix  $\boldsymbol{\Xi}_{m_2}^r$ , we can then define the generalized linear equation for the full time series as

$$\mathbb{E}[\bar{\Gamma}_{m_1}^*]_{m_2}^r = f_{f(s)}(\text{softmax}, \mathbf{X}_{m_2}^r \hat{\mathbf{H}}_{m_2}^r, \boldsymbol{\Xi}_{m_2}^r) \quad [21]$$

where the segment function mapper ( $f_{f(s)}$ ) applies the **softmax** function to the segments of  $\mathbf{X}_{m_2}^r \hat{\mathbf{H}}_{m_2}^r$  as defined by  $\boldsymbol{\Xi}_{m_2}^r$  (Appendix 5). The use of the mapping function is needed here, even for a simple prediction, as it defines the values included in the **softmax** function, whose output depends on the full set of inputs (Appendix 3). Further, because the **softmax** is log-linear, entries in  $\mathbf{X}_{m_2}^r \hat{\mathbf{H}}_{m_2}^r$  that are 0 are actually meaningful (they are equal to the reference topic) when included in the regression, and so we must only include true 0s, and not the filler-entry 0s, from  $\mathbf{X}_{m_2}^r \hat{\mathbf{H}}_{m_2}^r$  in the calculation. Equation 21 can be used to calculate the optimal model-estimated topic proportions for all documents in the corpus ( $\mathbb{E}[\bar{\Gamma}_{m_1}^*]_{m_2}^r$ ) by simply substituting the MAP estimated parameter matrix into the segment function mapper, as in  $f_{f(s)}(\mathbf{X}_{m_2}^r \hat{\mathbf{H}}_{m_2}^r, \text{softmax}, \boldsymbol{\Xi}_{m_2}^r)$ .

Inferentially, we are interested in determining the (marginal posterior) probability distribution of the change point locations, given topic proportions  $\mathcal{P}(\boldsymbol{\rho}_{m_2} | \bar{\Gamma}_{m_1}^*)$ , which we relate to the probability of the topic proportions, given the change point locations (*i.e.*, the model likelihood,  $\mathcal{P}(\bar{\Gamma}_{m_1}^* | \boldsymbol{\rho}_{m_2})$ ) and the prior distribution of the change points ( $\mathcal{P}(\boldsymbol{\rho}_{m_2})$ ) via Bayes' theory:

$$\mathcal{P}(\boldsymbol{\rho}_{m_2} | \bar{\Gamma}_{m_1}^*) \propto \mathcal{P}(\bar{\Gamma}_{m_1}^* | \boldsymbol{\rho}_{m_2}) \mathcal{P}(\boldsymbol{\rho}_{m_2}) \quad [22]$$

We use the segment-level models' loss values to construct the likelihood of the topic proportions under the change point locations. Specifically, we sum the minimized loss (negative log prior-penalized likelihood) values across the chunks to calculate the total minimized loss for a realization of the model ( $\mathcal{L}_{m_2}^*$ ), which is equal to the negative log likelihood for the topic proportions given the change point locations

$$\mathcal{L}_{m_2}^* = \sum_{s=1}^{S_{m_2}} \mathcal{L}_{m_2,s}^* = -\ell(\bar{\Gamma}_{m_1}^* | \boldsymbol{\rho}_{m_2}^r) \quad [23]$$

Remembering that a realization of a model simply specifies the change point locations, we can define the distribution of minimized loss  $\mathcal{L}_{m_2}^*$  for the model as a function of change point locations based on the  $R_{m_2} = \binom{(t_M - 1) - t_1}{p_{m_2}}$  unique realizations of, each of which has its own  $\mathcal{L}_{m_2}^*$  (Western and Kleykamp 2004):

$$\mathcal{P}(\bar{\Gamma}_{m_1}^* | \boldsymbol{\rho}_{m_2} = \boldsymbol{\rho}_{m_2}^r) = \frac{e^{\mathcal{L}_{m_2}^*}}{\sum_{r=1}^{R_{m_2}} e^{\mathcal{L}_{m_2}^r}} \quad [24]$$

Although because of the sheer number of possible realizations for even a modest time series with multiple change points (Table 2), we use Markov Chain Monte Carlo (MCMC) methods to sample the probability distributions efficiently (rather than sample it systematically).

The prior probability distribution for the change point locations ( $\mathcal{P}(\boldsymbol{\rho}_{m_2})$ ) is a multivariate discrete distribution, which could take any number of specific formulations (Western and Kleykamp 2004, Ruggieri 2013). The original LDATS model (Christensen *et al.* 2018) allowed only a uniform prior, and that requirement is presently maintained in the coding of the package (Simonis *et al. In Development*), although relaxing this assumption in the software is planned (see **FUTURE DEVELOPMENTS**). The uniform prior allocates equal probability to each of the discrete time points from the time of the first document ( $\min \mathbf{t}$ ) to one time step before the last document ( $\max \mathbf{t} - 1$ ), and the selected times are then sorted chronologically

Because the probability distribution has a high potential for multiple modes (if multiple change points are reasonably likely), standard MCMC approaches may have difficulty fitting the model (Sambridge 2014). Thus, we employ parallel tempering MCMC (ptMCMC; also called Metropolis-coupled or replica-exchange MCMC), which endows a standard MCMC search with auxiliary chains that explore the distribution surface more rapidly than the focal chain (Swendsen and Wang 1986, Geyer 1991, Earl and Deem 2005). ptMCMC is a robust methodology that works well for generalized models (Guo *et al.* 2016) and is an efficient sampler of rough probability landscapes (Machata and Ellis 2011). Here, the posterior probability distribution’s surface is explored using  $H$  chains ( $h$  in  $1 \dots H$ ), each with its own temperature ( $\mathbf{a}_h$ ) defining its search ability: higher temperature chains have higher variances in their step sizes (they have flattened surfaces to search; Gupta *et al.* 2018) and are therefore more easily able to navigate the surface. This comes at a cost of instability, however, as higher temperature chains are less likely to settle in to stable distributions. We therefore use a range of temperatures grouped in a series ( $\mathbf{a} = \mathbf{a}_1 < \mathbf{a}_2 < \dots < \mathbf{a}_H$ ) to balance search breadth and depth, with the coolest ( $\mathbf{a}_1$ ) sampling the true surface (Earl and Deem 2005).

The specifics of the temperature regime are obviously then critical for fitting a model to data using ptMCMC (Kone and Kofke 2005, Rathore *et al.* 2005, Nagata and Watanabe 2008). Following the original LDATS model (Christensen *et al.* 2018), the current implementation allows for control over the temperature sequence control parameters to facilitate fitting a wide range of potential corpus time series (Simonis *et al. In Development*) and defines the temperatures as

$$\mathbf{a}_{1 \dots H-1} = 2^{\frac{(\text{seq}(0, \log_2 \mathbf{a}_{H-1}, H-1))^{1+q}}{(\log_2 \mathbf{a}_{H-1})^q}} \quad [25]$$

$$\mathbf{a}_H = \mathbf{a}_H$$

where  $\text{seq}(0, \log_2 \mathbf{a}_{H-1}, H-1)$  is a sequence of values from 0 to  $\log_2 \mathbf{a}_{H-1}$  of length  $H-1$ ,  $q$  is the exponent controlling the temperature series ( $q = 0$  produces a geometric sequence,  $q = 1$  implies squaring before exponentiating),  $\mathbf{a}_{H-1}$  is the penultimate temperature, and  $\mathbf{a}_H$  is the ultimate temperature. Currently, the control inputs ( $\mathbf{a}_H$ ,  $\mathbf{a}_{H-1}$ ,  $H$ , and  $q$ ) are available to the user, but are fixed for a given fit of the  $\mathcal{M}_2$  Stage 2 models. A target for future mathematical and coding development is to expand the inputs and allow for an adaptive approach to ptMCMC, which will facilitate a more plug-and-play approach to model fitting (where the user does not need to set any control parameters; see **FUTURE DEVELOPMENTS**).

The ptMCMC algorithm works by coupling the chains (which are taking their own walks on the distribution surface) through “swaps”, where neighboring chains exchange configurations in between steps (Geyer 1991, Falcioni and Deem 1999). Over the course of the search, each of the chains proceeds through  $G$  iterations (steps;  $g$  in  $1 \dots G$ ) of the Metropolis-Hastings (MH) algorithm (Metropolis *et al.* 1953, Hastings 1970), a classical MCMC method. Step  $g$  on chain  $h$  is some realization of the model, such that the values of the change points for  $r = h, g$  are  $\boldsymbol{\rho}_{m_2}^{h,g}$ . The chains are each initialized ( $g = 0$ ) with a draw from  $\mathcal{P}(\boldsymbol{\rho}_{m_2})$

and the best fit draw (determined by likelihood) is put in the focal chain with the next best fit draw is put in the next hottest chain until the worst fit draw is put in the hottest chain. Then, from each step until  $g = G$ , for each chain, a new set of change points is proposed ( $\rho_{m_2^{h,g}}$ ; the breve accent indicates a proposal (pre)-step), evaluated (loss is calculated), and then either accepted ( $\rho_{m_2^{h,g}} = \rho_{m_2^{h,g}}$ ) or rejected ( $\rho_{m_2^{h,g}} = \rho_{m_2^{h,g-1}}$ ). The proposed set of change point locations for step  $g$  is generated from the proposal distribution  $\mathcal{R}$  conditional on the previous change point locations ( $\mathcal{R}(\rho_{m_2^{h,g}} | \rho_{m_2^{h,g-1}})$ ). The proposal distribution describes the movement of a single change point via a symmetric geometric distribution with a user-controlled average step size. Functionally,  $\mathcal{R}$  is a joint distribution representing three steps: [1] selecting one of the  $P_{m_2}$  change points to move via a multinomial distribution with equal probabilities ( $\frac{1}{P_{m_2}}$ ), [2] determining the directionality of movement (earlier or later in the time series) using a binomial distribution with equal probability for the outcomes  $-1$  and  $1$ , and [3] calculating the magnitude of the movement (number of discrete time steps) with a geometric distribution (mean step size is  $\kappa$ ). The multiplication of the three distributions results in the  $P_{m_2}$ -length vector representing the proposal step from the current change point locations:

$$\rho_{m_2^{h,g}} \sim \rho_{m_2^{h,g-1}} + \text{Geom}\left(1, \frac{1}{\kappa}\right) \text{Binom}_{\{-1,1\}}(1, 0.5) \text{Multinom}_{P_{m_2}}\left(1, \frac{1}{P_{m_2}}\right) \quad [26]$$

where **Geom** is the version of the geometric distribution that has 1 as its minimum returned value. The only parameter that is available to the user to set is the average step size of the geometric distribution,  $\kappa$ . Ostensibly, this parameter could be adaptively set based on the data set, but that functionality is not presently included (although see **FUTURE DEVELOPMENTS**).

Because the proposal distribution is symmetric ( $\mathcal{R}(\rho_{m_2^{h,g}} | \rho_{m_2^{h,g-1}}) = \mathcal{R}(\rho_{m_2^{h,g-1}} | \rho_{m_2^{h,g}})$ ), we are able to use the simplified Metropolis acceptance rule for the proposal (Metropolis *et al.* 1953, Hastings 1960). Acceptance of the proposal is probabilistic and based on the difference in energy ( $\Delta\mathcal{E}$ ) between the current ( $m_2^{h,g}$ ) and proposed ( $m_2^{h,g}$ ) realizations of the model:

$$\Delta\mathcal{E}_{m_2^{h,g}} = \mathcal{E}_{m_2^{h,g}} - \mathcal{E}_{m_2^{h,g}} \quad [27]$$

where the energy of  $m_2^r$  is the log of the inverse posterior evaluated at  $\rho_{m_2} = \rho_{m_2^r}$ ,

$$\mathcal{E}_{m_2^r} = \log \frac{1}{\mathcal{P}(\bar{\Gamma}_{m_1^*} | \rho_{m_2} = \rho_{m_2^r}) \mathcal{P}(\rho_{m_2^r})} \quad [28]$$

which is equal to the difference between the total minimized loss ( $\mathcal{L}_{m_2^r}^*$ ; negative log likelihood), and the log of the prior:

$$\mathcal{E}_{m_2^r} = \mathcal{L}_{m_2^r}^* - \log \mathcal{P}(\rho_{m_2^r}) \quad [29]$$

(Metropolis *et al.* 1953, Gupta *et al.* 2018). The change in energy associated with a proposal is converted to an acceptance probability for the step:

$$u_{m_2^{h,g}} = \min\left(1, e^{-b_h \Delta\mathcal{E}_{m_2^{h,g}}}\right) \quad [30]$$

where  $b_h$  is the inverse temperature of chain  $h$  (i.e.,  $b_h = \frac{1}{a_h}$ ), such that higher temperature chains have higher acceptance probabilities for the same energy difference. A random number from the standard uniform distribution is then drawn ( $U(0,1)$ ), with the proposal being accepted if the random number is less than or equal to  $u_{m_2}^{h,g}$ , and rejected if the random number is larger than  $u_{m_2}^{h,g}$ .

After each Metropolis MCMC iteration, the chains are able to swap configurations with their nearest neighbors in the temperature series (following the Metropolis criterion; Metropolis *et al.* 1953), allowing them to share information and search the surface in combination (Earl and Deem 2005). Starting with the hottest pair of chains ( $H$  and  $H - 1$ ) and descending in temperature, the chains (generally,  $h$  and  $h - 1$ ) swap information, and similar to the within-chain steps, the temperatures scale the swap acceptance rates for the sharing of information between chains, such that hotter chains are more likely to accept swaps (Earl and Deem 2005, Gupta *et al.* 2018). The acceptance probability for the swap between the neighboring chains is the exponentiated product of the difference in the chains' inverse temperatures and energies:

$$u_{m_2}^{h:h-1,g} = \min \left( 1, e^{\Delta b_{m_2}^{h:h-1,g} \Delta \mathcal{E}_{m_2}^{h:h-1,g}} \right) \quad [31]$$

where

$$\Delta b_{m_2}^{h:h-1,g} = \frac{1}{a_h} - \frac{1}{a_{h-1}} \quad [32]$$

and

$$\Delta \mathcal{E}_{m_2}^{h:h-1,g} = \mathcal{E}_{m_2}^{h,g} - \mathcal{E}_{m_2}^{h-1,g} \quad [33]$$

Currently, the schedule of Metropolis iterations and chain swaps and the inclusion of the chains in each swap are set (as “swaps between each iteration” and “all neighbor pairs included in each swap”), although they could become adaptive or user-controlled to allow for more efficient sampling (this is the most computationally intensive step in the LDATS modeling process; see **FUTURE DEVELOPMENTS**).

We remove the first steps  $G_{\text{burn-in}}$  (set by the user) steps as burn-in and thin the resulting sample to a fraction of  $\tau$  ( $\tau = 1$  equates to no thinning and  $\tau \in (0,1]$ ) steps (Link and Eaton 2012, Hobbs and Hooten 2015), with the resulting change point location vectors we store as rows in the matrix  $\mathbf{P}_{m_2}^{h=1}$  ( $\mathbf{P}$  as in capital  $\rho$ ), which is dimension  $\text{floor}(\tau(G - G_{\text{burn-in}} + 1)) \times P_{m_2}$ , where  $\text{floor}$  is the round-down function.  $\mathcal{P}(\boldsymbol{\rho}_{m_2} | \bar{\mathbf{F}}_{m_1}^*)$  is then defined by the proportional representation of each of the  $R_{m_2}$  realizations of  $\boldsymbol{\rho}_{m_2}^r$  in the rows of  $\mathbf{P}_{m_2}^{h=1}$ . Throughout the ptMCMC algorithm, we track step and swap acceptances and count the trips made by particles (the bits of information that can move among replica chains during swaps) from the hottest to the coolest chain (Katzgraber *et al.* 2006), and calculate a final rate for each based on the full set of  $G$  iterations (Earl and Deem 2005). These rates, along with trace plots (Kruschke 2015), constitute the ptMCMC diagnostics presently provided in LDATS.

#### Combining Segment-Level Models: Inference About Within-Segment Parameters

Having defined  $\hat{\mathbf{H}}_{m_2}^r$  as the collated matrix of segment-level parameter matrices ( $\hat{\mathbf{H}}_{m_2,s}^r$ ) fit based on the realized change point locations and having estimated the probability of that specific set of change point locations using ptMCMC, we can now calculate the marginal posterior distribution of the regression parameters given the topic proportions ( $\mathcal{P}(\hat{\mathbf{H}}_{m_2}^r | \bar{\mathbf{F}}_{m_1}^*)$ ), acknowledging the uncertainty in the change point

locations (Western and Kleykamp 2004). In order to accomplish this, however, we must account for the variance-covariance structure among the parameters within  $\hat{\mathbf{H}}_{m_2}^r$ , which is facilitated by first unwinding the matrix into a vector  $\widetilde{\hat{\boldsymbol{\eta}}_{m_2}^r}$ , as indicated by the tilde accent, which constitutes the mean vector for the parameter distributions ( $\widetilde{\hat{\boldsymbol{\eta}}_{m_2}^r}$  is the unwound segment-specific matrix  $\hat{\mathbf{H}}_{m_2}^r$ ). Then, we place the precision (inverse variance-covariance) matrices for the segment-level model fits (generally,  $\mathbf{A}_{\widetilde{\hat{\boldsymbol{\eta}}_{m_2}^r}} = \boldsymbol{\Sigma}_{\widetilde{\hat{\boldsymbol{\eta}}_{m_2}^r}}^{-1}$ ) along the diagonal of a square ( $S_{m_2} \times S_{m_2}$ ) matrix, producing the block-diagonal precision matrix for the full time series model given the realized change points,  $\mathbf{A}_{\widetilde{\hat{\boldsymbol{\eta}}_{m_2}^r}}$ , whose off-diagonal entries are  $\infty$  (*i.e.*, no covariance). The probability distribution for the full set of regression parameters, given a realized set of change point locations and the topic proportions is then described by a multivariate normal distribution (MVN; Tiao and Zellner 1964) with mean vector  $\widetilde{\hat{\boldsymbol{\eta}}_{m_2}^r}$  and precision matrix  $\mathbf{A}_{\widetilde{\hat{\boldsymbol{\eta}}_{m_2}^r}$ :

$$\mathcal{P}(\hat{\mathbf{H}}_{m_2} | \boldsymbol{\rho}_{m_2}^r, \bar{\mathbf{F}}_{m_1}^*) = \mathcal{P}_{\text{MVN}_{S_{m_2} C_{m_2} k_{m_1}^*}}(\widetilde{\hat{\boldsymbol{\eta}}_{m_2}^r}, \mathbf{A}_{\widetilde{\hat{\boldsymbol{\eta}}_{m_2}^r}}) \quad [34]$$

where  $\mathcal{P}_{\text{MVN}}$  is the probability density function of an MVN defined for a set of specific parameter values  $\mathbf{x}$

$$\mathcal{P}_{\text{MVN}_{S_{m_2} C_{m_2} k_{m_1}^*}}(\widetilde{\hat{\boldsymbol{\eta}}_{m_2}^r}, \mathbf{A}_{\widetilde{\hat{\boldsymbol{\eta}}_{m_2}^r}}) = \frac{1}{\sqrt{|2\pi \mathbf{A}_{\widetilde{\hat{\boldsymbol{\eta}}_{m_2}^r}}^{-1}|}} e^{-\frac{1}{2}(\mathbf{x} - \widetilde{\hat{\boldsymbol{\eta}}_{m_2}^r})' \mathbf{A}_{\widetilde{\hat{\boldsymbol{\eta}}_{m_2}^r}} (\mathbf{x} - \widetilde{\hat{\boldsymbol{\eta}}_{m_2}^r})} \quad [35]$$

and  $S_{m_2} C_{m_2} k_{m_1}^*$  is the dimensionality of the specific MVN distribution here (one parameter per covariate per topic per segment). The value of  $\mathcal{P}(\hat{\mathbf{H}}_{m_2} | \boldsymbol{\rho}_{m_2}^r, \bar{\mathbf{F}}_{m_1}^*)$  is determined for each of the  $R_{m_2}$  realizations, which are then combined via Eq. 13 to estimate the marginal posterior probability of the regression coefficients, given just the topic proportions ( $\mathcal{P}(\hat{\mathbf{H}}_{m_2} | \bar{\mathbf{F}}_{m_1}^*)$ ). The original LDATS application did not calculate or return the posterior estimates of the regression parameters as the research was focused on the change point locations (Christensen *et al.* 2018), and so this constitutes a key update to the application.

### Multi-Model Inference

Given the fit of a specific Stage 2 model ( $m_2$ ) defined by the number of change points ( $P_{m_2}$ ) and the within-segment covariate model (the specific  $C_{m_2}$  predictors comprising  $\mathbf{X}_{m_2}$ ), we can then consider multiple Stage 2 models to determine the optimal configuration of the time series component.

As with Stage 1, we use AIC as our Stage 2 model selection criterion (Christensen *et al.* 2018, based on Gelman *et al.* 2014), defined for a specific time series model  $m_2$ :

$$\text{AIC}_{m_2} = -2\ell(\boldsymbol{\rho}_{m_2}, \hat{\mathbf{H}}_{m_2} | \bar{\mathbf{F}}_{m_1}^*) + 2k_{m_2} \quad [36]$$

where  $\ell$  is the log likelihood ( $\ell(\boldsymbol{\rho}_{m_2}, \hat{\mathbf{H}}_{m_2} | \bar{\mathbf{F}}_{m_1}^*) = \log \mathcal{P}(\bar{\mathbf{F}}_{m_1}^* | \boldsymbol{\rho}_{m_2}, \hat{\mathbf{H}}_{m_2})$ ) estimated from the ptMCMC samples retained after burn-in and thinning, and  $k_{m_2} = (S_{m_2} C_{m_2} (k_{m_1} - 1) + P_{m_2})$  is the number of parameters in the model: for each segment (in  $S_{m_2}$  total), there are  $C_{m_2} (k_{m_1} - 1)$  parameters fit for the multinomial regression (one for every covariate for every topic, save the first to account for the sum-to-1 constraint, Appendix 3), plus the  $P_{m_2}$  change point locations (Western and Kleykamp 2004, Christensen *et al.* 2018). The total number of models in Stage 2 ( $\mathcal{M}_2$ ) is

$$\mathcal{M}_2 = (P_{\max} + 1)\mathcal{U}_{\mathbf{X}} \quad [37]$$

where  $P_{\max}$  is the maximum number of change points included (and the  $+1$  to account for  $P_{m_2} = 0$ ) and  $\mathcal{U}_{\mathbf{X}}$  is the number of unique configurations of the covariate matrix ( $\mathbf{X}$ ). The optimal (according to AIC) time series model ( $m_2^*$ ) is determined by

$$m_2^* = \arg \min_{m_2 \in 1 \dots \mathcal{M}_2} \text{AIC}(m_2) \quad [38]$$

and has the corresponding set of parameters  $\{\boldsymbol{\rho}_{m_2^*}, \hat{\mathbf{H}}_{m_2^*}\}$ .

## MODEL EVALUATION

{This section will focus on evaluating the model under known parameters}

## MODEL APPLICATION

{This section will include one or two brief example applications of LDATS to real data sets}

## DISCUSSION

{This section will do the normal discussing of the present work in the context of the literature}

## FUTURE DEVELOPMENTS

The LDATS modeling framework (as presented here) and code package (Simonis *et al. In Development*) are stable and robust, but as with any methodology, improvements can be made through future developments. We have noted multiple components slated for developments throughout this manuscript; we add some additional components and briefly describe our current plans in this section.

Although the VEM approach to LDA model fitting is well developed in its application (Blei *et al.* 2003, Grün and Hornik 2011), there is one notable area of improvement that could be included: adaptive selection of the optimal number of topics  $k_{m_1^*}$ . The other common method to parameter estimation under LDA models, Gibbs sampling, provides a flexible approach to estimating the number of topics by considering it a free parameter to be fit (Griffiths and Steyvers 2004, Valle *et al.* 2018). However, there is also a method for adaptive selection of the number of topics under a VEM approach to LDA based on topic densities (Cao *et al.* 2009), which could be integrated into LDATS without needing to alter the underlying inference machinery.

The time series model options (based on Western and Kleykamp 2004) in LDATS do not presently include temporal autocorrelation, which is a feature of many time series (Cressie and Wikle 2011), and could be mistaken for (statistically speaking) change points, especially if the magnitude of the change is small (Jarušková 1997, Wang 2008, Beaulieu *et al.* 2012). Although there are methods for including temporal autocorrelation into change point models (Lund *et al.* 2007, Jandhyala *et al.* 2010, Beaulieu *et al.* 2012, Pandya *et al.* 2012), they all operate under the assumption of a univariate normal response variable, which is a much simpler segment-level regression model than the present softmax and offers the simple entrée of the error term for autocorrelation. Methods exist for including autocorrelation structure in multinomial regressions (Linderman *et al.* 2015), but have yet to be integrated with a change point model.

By splitting the corpus prior to fitting the segment-level models, LDATS makes the important assumption of no covariance among regressors across segments, although it is possible that there are substantially-non-0 covariances. It could therefore improve the predictive capacity of the model to fit the model across-segments in a fashion to allow covariances. Such a change to LDATS would require a shift in the underlying inferential machinery and is likely to considerably slow computation (Genz 1992), however.

LDATS also assumes mean-0, uncorrelated, vague priors for the regression parameters and a vague prior for the change point locations. Although the vagueness of the regression prior can be toggled by the user, no other control on the priors is available, and therefore the full strength of the Bayesian approach is not yet being leveraged in the time series component of LDATS (Chin Choy and Broemeling 1980, Western and Kleykamp 2004). Including more nuanced control of the priors necessitates substantial changes to the inferential code underlying the model, and in particular will require a shift in the fitting of the multinomial models away from an established code base (Venables and Ripley 2003).

Although the parameters controlling the ptMCMC algorithm ( $a_H$ ,  $a_{H-1}$ ,  $H$ ,  $q$ , and  $\kappa$ ) are available to the user to input, they are limited to fixed values, and so require the user to *a priori* know appropriate values or spend considerable time testing them. In addition, the swap schedule (in terms of frequency of swaps and the inclusion of chains within a swap) used for the ptMCMC is presently fixed and not available to control at all by the user. All of these controls have the potential to influence the convergence of the ptMCMC algorithm and the speed at which convergence is achieved. Therefore, the model's application could be improved by allowing flexibility in and control over the swap schedule and by taking an adaptive approach to ptMCMC that allows the controls to optimize themselves via the fitting process (Katzgraber *et al.* 2006, Trebst *et al.* 2006, Hasenbusch and Schaefer 2010, Miasojedow *et al.* 2013).

## ACKNOWLEDGEMENTS

The motivating study—known as the Portal Project—has been funded nearly continuously since 1977 by the National Science Foundation, most recently by DEB-1622425 to S. K. M. Ernest, which also supported (in part) E. Christensen's time. Much of the computational work (including time of J. Simonis, D. Harris, and H. Ye) was supported by the Gordon and Betty Moore Foundation's Data-Driven Discovery Initiative through Grant GBMF4563 to E. P. White. R. Diaz was supported in part by a National Science Foundation Graduate Research Fellowship (No. DGE-1315138 and DGE-1842473).

## AUTHOR CONTRIBUTIONS

J. L. Simonis provided insight on LDA applications and feedback on technical writing during development of the first version of the LDATS model and application, led the coding and mathematical development of the model into an R package, and led the writing on this manuscript. E. M. Christensen led the project during development of the first version of the LDATS model and its application to the Portal data, specifically conceiving the project, coding the pipeline wrappers of the analysis, and writing and editing the first description of the model and its application (Christensen *et al.* 2018). D. J. Harris was involved in developing and applying the first version of the LDATS model, specifically suggesting the LDA and change point approaches, coding the first version of the change point model, and writing and editing the first description of the model (Christensen *et al.* 2018). R. Diaz contributed code to the LDATS package, provided insight into model development, and conducted end-user code application testing. H. Ye contributed code to the LDATS package and insight into data structures and LDA algorithms. E. P. White helped design, troubleshoot, and supervise initial methods development and provided big-picture feedback on development of the R package. S. K. Morgan Ernest provided managerial oversight and feedback on the project in both the initial and second stages of LDATS development, tested applications of the code to data sets, and assisted with writing and editing of the first description of the model and its application (Christensen *et al.* 2018).

## LITERATURE CITED



- Andersen, T., J. Carstensen, E. Hernández-García, and C. M. Duarte. 2009. Ecological thresholds and regime shifts: approaches to identification. *Trends in Ecology and Evolution* **24**:49-57.
- Bartels, L. M. 1997. Specification uncertainty and model averaging. *American Journal of Political Science* **41**:641-674.
- Bassett, R. and J. Deride. 2018. Maximum *a posteriori* estimators as a limit of Bayes estimators. *Mathematical Programming, Series B*. <https://doi.org/10.1007/s10107-018-1241-0>
- Beauleiu, C., J. Chen, and J. L. Sarmiento. 2012. Change-point analysis as a tool to detect abrupt climate variations. *Philosophical Transactions of the Royal Society A* **370**:1228-1249.
- Berger, A. L., V. J. Della Pietra, and S. A. Della Pietra. 1996. A maximum entropy approach to natural language processing. *Computational Linguistics* **22**:39-71.
- Blei, D. M., A. Y. Ng, and M. I. Jordan. 2003. Latent Dirichlet Allocation. *Journal of Machine Learning Research* **3**: 993-1022.
- Bishop, C. M. 2006. Pattern Recognition and Machine Learning. Springer.
- Broyden, C. G. 1970. The convergence of a class of double-rank minimization algorithms. *Journal of the Institute of Mathematics and Its Applications* **6**:76-90.
- Cao, J., T. Xia, J. Li, Y. Zhang, and S. Tang. 2009. A density-based method for adaptive LDA model selection. *Neurocomputing* **72**:1775-1781.
- Chin Choy, J. H. and L. D. Broemeling. 1980. Some Bayesian inferences for a changing linear model. *Technometrics* **22**:71-78.
- Christensen, E. M., D. J. Harris, and S. K. M. Ernest. 2018. Long-term community change through multiple rapid transitions in a desert rodent community. *Ecology* **99**:1523-1529.
- Cressie, N. and C. K. Wikle. 2011. Statistics for Spatio-Temporal Data. Wiley, Hoboken, NJ, USA.
- Earl, D. J. and M. W. Deem. 2005. Parallel tempering: theory, applications, and new perspectives. *Physical Chemistry Chemical Physics* **7**: 3910-3916.
- Falcioni, M. and M. W. Deem. 1999. A biased Monte Carlo scheme for zeolite structure solution. *Journal of Chemical Physics* **110**:1754-1766.
- Fletcher, R. 1970. A new approach to variable metric algorithms. *Computer Journal* **13**:317-322.
- Gal, G. and W. Anderson. 2010. A novel approach to detecting a regime shift in a lake ecosystem. *Methods in Ecology and Evolution* **1**:45-52.
- Gelman, A. and J. Gill. 2007. Data Analysis Using Regression and Multilevel/Hierarchical Models. Cambridge University Press, Cambridge, UK.
- Gelman, A., J. Hwang, and A. Vehtari. 2014. Understanding predictive information criteria for Bayesian models. *Statistics and Computing* **24**:997-1016.

- Genz, A. 1992. Numerical computation of the multivariate normal probabilities. *Journal of Computational and Graphical Statistics* **1**:141-150.
- Geyer, C. J. 1991. Markov Chain Monte Carlo maximum likelihood. In *Computing Science and Statistics: Proceedings of the 23rd Symposium on the Interface*. pp 156-163. American Statistical Association, New York, USA.
- Goldfarb, D. 1970. A family of variable metric updates derived by variational means. *Mathematics of Computation* **24**:23-26.
- Griffiths, T. L. and M. Steyvers. 2004. Finding scientific topics. *Proceedings of the National Academy of Sciences* **101**:5228-5235.
- Grün, B. and K. Hornik. 2011. topicmodels: an R package for fitting topic models. *Journal of Statistical Software* **40**:13. <http://dx.doi.org/10.18637/jss.v040.i13>
- Guo, G., W. Shao, L. Lin, and X. Zhu. 2016. Parallel tempering for dynamic generalized linear models. *Communications in Statistics—Theory and Methods* **45**:6299-6310
- Gupta, S., L. Hainsworth, J. S. Hogg, R. E. C. Lee, and J. R. Faeder. 2018. Evaluation of parallel tempering to accelerate Bayesian parameter estimation in systems biology. arXiv:1801.09831 [q-bio.QM]
- Hasenbusch, M. and S. Schaefer. 2010. Speeding up parallel tempering simulations. *Physical Review* **82**:046707.
- Hastings, W. K. 1970. Monte Carlo sampling methods using Markov Chains and their applications. *Biometrika* **57**:97-109.
- Hobbs, N. T. and M. B. Hooten. 2015. Bayesian Models: A Statistical Primer for Ecologists. Princeton University Press, Princeton, NJ, USA.
- Hoerl, A. E. 1962. Application of ridge analysis to regression problems. *Chemical Engineering Progress*. **58**:54-59.
- Hoerl, A. E. and R. W. Kennard. 1970. Ridge regression: biased estimation for nonorthogonal problems. *Technometrics* **12**:55-67.
- Jandhyala, V. K., S. B. Fotopoulos, and J. You. 2010. Change-point analysis of mean annual rainfall data from Tucumán, Argentina. *Environmetrics* **21**:687-697.
- Jarušková, D. 1997. Some problems with application of change-point detection methods to environmental data. *Environmetrics* **8**:469–483.
- Jordan, M. Z. Ghahramani, T. Jaakkola, and L. Saul. 1999. Introduction to variational methods for graphical models. *Machine Learning* **37**:183-233.
- Karssenbergh, D., M. F. P. Bierkens, and M. Rietkerk. 2017. Catastrophic shifts in semiarid vegetation-soil systems may unfold rapidly or slowly. *The American Naturalist* **190**:e145-e155.
- Katzgraber, H. G., S. Trebst, D. A. Huse. And M. Troyer. 2006. Feedback-optimized parallel tempering Monte Carlo. *Journal of Statistical Mechanics: Theory and Experiment* **3**:P03018.
- Kone, A. and D. A. Kofke. 2005. Selection of temperature intervals for parallel-tempering simulations. *The Journal of Chemical Physics* **122**:206101.

Kruschke, J. K. 2015. Doing Bayesian Data Analysis: A Tutorial with R, JAGS, and Stan. 2<sup>nd</sup> Edition. Elsevier Academic Press, New York, NY, USA.

Linderman, S., M. Johnson, and R. P. Adams. 2015. Dependent multinomial models made easy: stick-breaking with the Pólya-gamma augmentation. *Advances in Neural Information Processing Systems* **28**: 3456-3464.

Link, W. A. and M. J. Eaton. 2012. On thinning of chains in MCMC. *Methods in Ecology and Evolution* **3**:112-115.

Lund, R., X. L. Wang, Q. Lu, J. Reeves, C. Gallagher, and Y. Feng. 2007. Changepoint detection in periodic and autocorrelated time series. *Journal of Climate* **20**:5178-5190.

Machata, J. and R. S. Ellis. 2011. Monte Carlo methods for rough free energy landscapes: population annealing and parallel tempering. *Journal of Statistical Physics* **144**:541-553.

Malouf, R. 2002. A comparison of algorithms for maximum entropy estimation. *Proceedings of the 6th Conference on Natural Language Learning* **20**:1-7.

Miasojedow, B., E. Moulines, and M. Vihola. 2013. An adaptive parallel tempering algorithm. *Journal of Computational and Graphical Statistics* **22**:649-664.

McCullagh, P. and J. A. Nelder. 1989. Generalized Linear Models. 2<sup>nd</sup> Edition. Chapman and Hall, New York, NY, USA.

McCune, B. and J.B. Grace. 2002. Analysis of Ecological Communities. MjM Software. Gleneden Beach, OR, USA.

Metropolis, N., A. W. Rosenbluth, M. N. Rosenbluth, A. H. Teller, and E. Teller. 1953. Equations of state calculations by fast computing machines. *Journal of Chemical Physics* **21**: 1087-1092

Nagata, K. and S. Watanabe. 2008. Asymptotic behavior of exchange ratio in exchange Monte Carlo method. *Neural Networks* **21**:980-988.

Ng, A. Y. and M. I. Jordan. 2002. On discriminative vs. generative classifiers: a comparison of logistic regression and naïve Bayes. In *Advances in Neural Information Processing Systems*, eds: T. G. Dietterich and S. Becker and Z. Ghahramani. pp 841-848. MIT Press, Boston, MA, USA.

Ng, A. 2004. Feature selection, L1 vs. L2 regularization, and rotation invariance. *Proceedings of the Twenty-first International Conference on Machine Learning*: 78.

Pandya, M., K. Bhatt, and H. C. Thakar. 2012. Bayesian estimation of change point in autoregressive process. *International Journal of Research and Reviews in Applied Sciences* **13**:41-52.

R Core Team. 2018. R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria.

Ratajczak, Z., S. R. Carpenter, A. R. Ives, C. J. Kucharik, T. Ramiadantsoa, M. A. Stegner, J. W. Williams, J. Zhang, and M. G. Turner. 2018. Abrupt change in ecological systems: inference and diagnosis. *Trends in Ecology & Evolution* **33**:513-526.

Rathore, N., M. Chopra., and J. J. de Pablo. 2005. Optimal allocation of replicas in parallel tempering simulations. *The Journal of Chemical Physics* **122**:024111.

- Ripley, B. D. 1993. Statistical aspects of neural networks. In *Networks and Chaos: Statistical and Probabilistic Aspects*, eds: O. E. Barndorff-Nielsen, J. L. Jensen, and W. S. Kendall. pp 40-123. Chapman Hall, London, UK.
- Ripley, B. D. 1996. Pattern Recognition and Neural Networks. Cambridge University Press, Cambridge, UK.
- Ronning, G. 1989. Maximum likelihood estimation of Dirichlet distributions. *Journal of Statistical Computation and Simulation* **34**:215-221.
- Ruggieri, E. 2013. A Bayesian approach to detecting change points in climatic records. *International Journal of Climatology* **33**:520-528.
- Sambridge, M. 2014. A Parallel Tempering algorithm for probabilistic sampling and multimodal optimization. *Geophysical Journal International* **196**:357-374.
- Scheffer, M. 2009. Critical transitions in nature and society. Princeton University Press, Princeton, NJ, USA.
- Shanno, D. F. July 1970. Conditioning of quasi-Newton methods for function minimization. *Mathematics of Computation* **24**:647-656.
- Simonis, J. L., E. M. Christensen, D. J. Harris, H. Ye, R. Diaz, E. P. White, and S. K. Morgan Ernest. *In Development*. LDATS: Latent Dirichlet Allocation coupled with Time Series analyses. R package v 0.0.6. <https://github.com/weecology/LDATS>
- Swendsen, R. H. and J. S. Wang. 1986. Replica Monte Carlo simulation of spin-glasses. *Physical Review Letters* **57**:2607-2609.
- Tiao, G. C. and A. Zellner. 1964. On the Bayesian estimation of multivariate regression. *Journal of the Royal Statistical Society: Series B (Methodological)* **26**:277-285 .
- Tikhonov, A. N.; V. Y. Arsenin. 1977. Solution of Ill-posed Problems. Winston & Sons, Washington DC, USA.
- Tingley, M.W., W.B. Monahan, S.R. Beissinger, and C. Moritz. 2009. Birds track their Grinnellian niche through a century of climate change. *Proceedings of the National Academy of Sciences* **106**:19637–19643.
- Trebst, S., M. Troyer, and U. H. E. Hansmann. 2006. Optimized parallel tempering simulations of proteins. *The Journal of Chemical Physics* **124**:174903.
- Valle, D., B. Balser, C. W. Woodall, and R. Chazdon. 2014. Decomposing biodiversity data using the Latent Dirichlet Allocation model, a probabilistic multivariate statistical model. *Ecology Letters* **17**: 1591-1601.
- Valle, D., P. Albuquerque, Q. Zhao, A. Barberan, and R. J. Fletcher Jr. 2018. Extending the Latent Dirichlet Allocation model to presence/absence data: a case study on North American breeding birds and biogeographic shifts expect from climate change. *Global Change Biology* **24**:5560-5572.
- Venables, W. N. and B. D. Ripley. 2002. Modern and Applied Statistics with S. Fourth Edition. Springer, New York, NY, USA.
- Wang, X. L. L. 2008. Accounting for autocorrelation in detecting mean shifts in climate data series using the penalized maximal t or F test. *Journal of Applied Meteorology and Climatology* **47**:2423–2444.

Western, B. and M. Kleykamp. 2004. A Bayesian change point model for historical time series analysis. *Political Analysis* **12**:354-374.

Williams, J.W., J.L. Blois, and B.N. Shuman. 2011. Extrinsic and intrinsic forcing of abrupt ecological change: case studies from the late Quaternary. *Journal of Ecology* **99**:664-677.

## TABLES

**Table 1.** Definitions of the notation used in LDATS.

Definition	Parameter	Expansions
Corpus (set of documents)	$D$	
Total documents	$M$	
Specific document	$d$	
Total words	$N$	
Total words in a document	$N_d$	
Specific word	$n$	
Total terms	$V$	
Specific term	$v$	
Weight of a document	$u_d$	$\mathbf{u}$
Time of a document	$t_d$	$\mathbf{t}$
Probability	$\mathcal{P}$	
Loss	$\mathcal{L}$	
Log-likelihood	$\ell$	
Expected value	$E$	
Number of fitted parameters	$k$	
Optimal value of a fit (parameter or model)	$*$	
Replicates of a model with different starting values	$\mathcal{N}$	
Derivative	$d$	$\mathcal{D}$
General index variables	$i, j$	
Total Stage 1 models	$\mathcal{M}_1$	
Specific Stage 1 model	$m_1$	
Total latent topics in a model	$k_{m_1}$	
Specific latent topic	$i$	
Observed word identity	$w_{d,n_d}$	$\mathbf{w}_d, \mathbf{w}$
Latent topic membership	$z_{m_1,d,n_d}$	$\mathbf{z}_{m_1,d}, \mathbf{z}_{m_1}$
Within-document topic probability	$\theta_{m_1,d,i}$	$\boldsymbol{\theta}_{m_1,d}, \boldsymbol{\theta}_{m_1}$
Document-level topic concentration	$\alpha_{m_1,i}$	$\alpha_{m_1}, \boldsymbol{\alpha}_{m_1,d}$
Word-level term probability	$\beta_{m_1,d,i,v}$	$\mathbf{B}_{m_1,d}, \mathbf{B}_{m_1}$
Random number generator seed	$\mathbf{f}$	
Variational distribution	$\mathcal{Q}$	
Variational lower bound	$L$	
Kullback-Leibler Divergence	$D_{KL}$	
Within-document topic concentration (variational)	$\gamma_{m_1,d,i}$	$\boldsymbol{\gamma}_{m_1,d}, \boldsymbol{\Gamma}_{m_1}$
Within-document topic probability (variational)	$\bar{\gamma}_{m_1,d,i_{m_1}}$	$\bar{\boldsymbol{\gamma}}_{m_1,d}, \bar{\boldsymbol{\Gamma}}_{m_1}$
Word-level topic probability (variational)	$\phi_{m_1,d,n,i}$	$\boldsymbol{\phi}_{m_1,d,n}, \boldsymbol{\Phi}_{m_1}$
Gamma function	$\Gamma$	
Digamma function	$\Psi$	
Indication of term identity for an observed word	$\omega_{d,n_d}^v$	
Total Stage 2 models	$\mathcal{M}_2$	

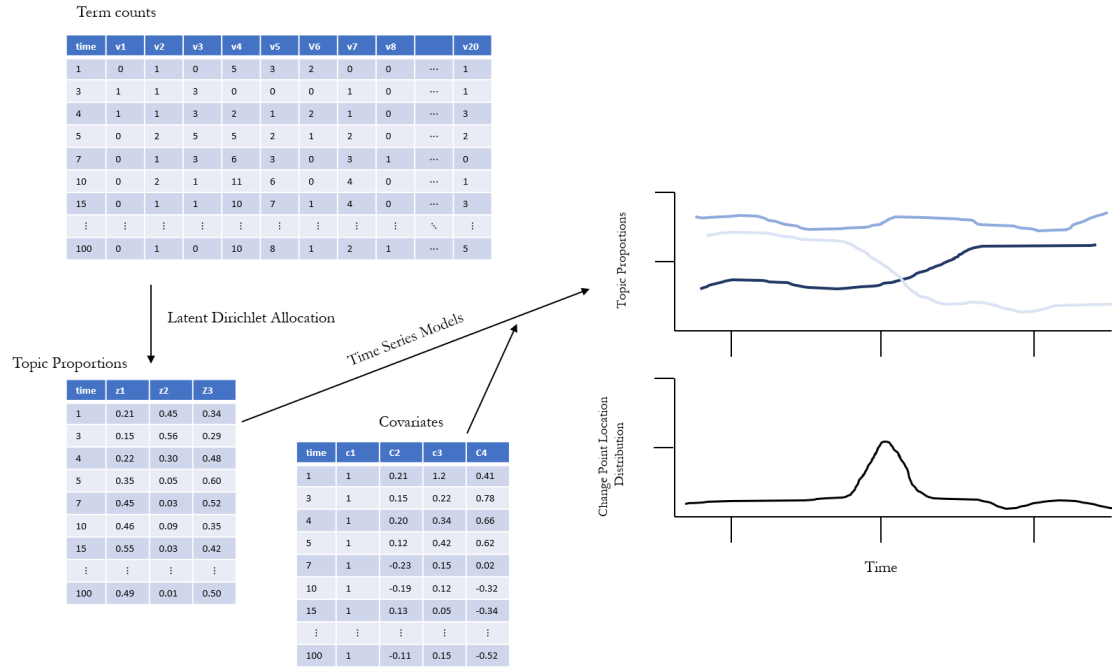
Specific Stage 2 model	$m_2$	
Total realizations of a specific Stage 2 model	$R_{m_2}$	
Specific realization of a specific Stage 2 model	$r$	
Total change points	$P_{m_2}$	
Specific change point	$p$	
Total segments of the time series	$S_{m_2}$	
Specific segment	$s$	
Change point location	$\rho_{m_2,p}$	$\rho_{m_2}$
Change point locations augmented with time range		$\tilde{\rho}_{m_2}$
Indication if a document is from a segment	$\xi_{m_2^r,d,s}$	$\xi_{m_2^r,s}, \Xi_{m_2^r}$
Total covariates in a model	$C_{m_2}$	
Specific covariate	$c$	
Value of a covariate	$x_{m_2,d,c_{m_2}}$	$\mathbf{x}_{m_2,d}, \mathbf{X}_{m_2,s}, \mathbf{X}_{m_2}$
Within-segment regressors	$\eta_{m_2^r,s,i,c}$	$\boldsymbol{\eta}_{m_2^r,s,i}, \mathbf{H}_{m_2^r,s}$
Augmented within-segment regressors	$\dot{\eta}_{m_2^r,s,i,c}$	$\dot{\eta}_{m_2^r,s,i}, \dot{\mathbf{H}}_{m_2^r,s}, \dot{\mathbf{H}}_{m_2^r}, \dot{\mathbf{H}}_{m_2}$
Precision on the prior for the regressors	$\lambda$	
General augmented within-segment regressor values	$\mathbf{x}$	$\mathbf{x}$
Total MCMC chains	$H$	
Specific MCMC chain	$h$	
Temperature of a chain	$a_h$	$\mathbf{a}$
Inverse temperature of a chain	$b_h$	$\mathbf{b}$
Exponent controlling the temperature sequence	$q$	
Final temperature	$a_H$	
Penultimate temperature	$a_{H-1}$	
Total MCMC iterations	$G$	
Specific MCMC iteration	$g$	
Proposal iteration for the step to $g$	$\tilde{g}$	
Statistical energy	$\mathcal{E}$	
Acceptance probability	$u$	
Proposal distribution	$\mathcal{R}$	
Average step size in the proposal distribution	$\kappa$	
Number of burn-in iterations	$G_{\text{burn-in}}$	
Fraction of samples remaining after thinning	$\tau$	
Document to segment mapping function	$f_{d \rightarrow s}$	
Cross entropy function	$f_{\text{CE}}$	
Softmax function	$\text{softmax}, f_S$	
Matrix multiplication function	$f_{\text{MM}}$	
Softmax of matrix multiplication function	$f_{S(\text{MM})}$	
Cross entropy of softmax of matrix multiplication	$f_{\text{CE}(S(\text{MM}))}$	
Penalty function	$f_P$	
Partition function	$f_{\text{part}}$	
Segment function mapper	$f_{f(s)}$	
Kronecker Delta	$\delta_{ij}$	

**Table 2.** Number of unique model realizations for a combination of time series length and number of change points.

Time steps → Change points ↓	50	100	200	500
1	49	99	199	499
2	1176	4851	1970	124251
3	18424	156849	1293699	10507399
4	211876	3764376	63391251	2552446876
5	1906884	71523144	2472258789	252692240724



# FIGURES



**Fig 1.** Schematic representation of the data-model relation of the LDATS framework. Documents are described by high-dimensional term counts, which are first reduced to low-dimension topic proportions. Then, the topic proportions are analyzed with time series models that include both covariate and change point dynamics.

## Appendix 1: Expanding the probability of a corpus given an LDA model

For a given Latent Dirichlet Allocation (LDA; Blei *et al.* 2003)  $m_1$ , the probability of the entire corpus, given the parameters, is the product of the probabilities of each document, given the parameters:

$$\mathcal{P}(\mathbf{w}|\alpha_{m_1}, \mathbf{B}_{m_1}) = \prod_{d=1}^M \mathcal{P}(\mathbf{w}_d|\alpha_{m_1}, \mathbf{B}_{m_1}) \quad [\text{A1.1}]$$

The probability of a document's data given the parameters is the product of [1] the word-level term-identity distributions, given the model parameters governing the within-document topic distribution and word-level term distributions ( $\mathcal{P}(\mathbf{w}_{d,n}|\boldsymbol{\theta}_{m_1,d}, \mathbf{B}_{m_1})$ ) and [2] the within-document topic distribution parameters given the document-level concentration parameter ( $\mathcal{P}(\boldsymbol{\theta}_{m_1,d}|\alpha_{m_1})$ ), integrated over the uncertainty in the within-document topic distribution ( $\boldsymbol{\theta}_{m_1,d}$ ):

$$\mathcal{P}(\mathbf{w}_d|\alpha_{m_1}, \mathbf{B}_{m_1}) = \int \mathcal{P}(\boldsymbol{\theta}_{m_1,d}|\alpha_{m_1}) \left( \prod_{n=1}^{N_d} \mathcal{P}(w_{d,n}|\boldsymbol{\theta}_{m_1,d}, \mathbf{B}_{m_1}) \right) d\boldsymbol{\theta}_{m_1,d} \quad [\text{A1.2}]$$

The word-level topic-identity distribution can be further decomposed into the product of [1] the term identity distribution given the topic identity and the word-level term distribution parameters ( $\mathcal{P}(w_{d,n}|z_{m_1,d,n}, \mathbf{B}_{m_1})$ ) and [2] the topic identity distribution given the parameters governing the within-document topic distribution ( $\mathcal{P}(z_{m_1,d,n}|\boldsymbol{\theta}_{m_1,d})$ ), integrated (summed due to discreteness) over the uncertainty in topic type  $z_{m_1,d,n}$

$$\mathcal{P}(w_{d,n}|\boldsymbol{\theta}_{m_1,d}, \mathbf{B}_{m_1}) = \sum_{z_{m_1,d,n}} \mathcal{P}(w_{d,n}|z_{m_1,d,n}, \mathbf{B}_{m_1}) \mathcal{P}(z_{m_1,d,n}|\boldsymbol{\theta}_{m_1,d}) \quad [\text{A1.3}]$$

Substituting Eq. A1.3 into Eq. A1.2,

$$\mathcal{P}(\mathbf{w}_d|\alpha_{m_1}, \mathbf{B}_{m_1}) = \int \mathcal{P}(\boldsymbol{\theta}_{m_1,d}|\alpha_{m_1}) \left( \prod_{n=1}^{N_d} \sum_{z_{m_1,d,n}} \mathcal{P}(w_{d,n}|z_{m_1,d,n}, \mathbf{B}_{m_1}) \mathcal{P}(z_{m_1,d,n}|\boldsymbol{\theta}_{m_1,d}) \right) d\boldsymbol{\theta}_{m_1,d} \quad [\text{A1.4}]$$

and then Eq. A1.4 into Eq 1.1,

$$\mathcal{P}(\mathbf{w}|\alpha_{m_1}, \mathbf{B}_{m_1}) = \prod_{d=1}^M \left[ \int \mathcal{P}(\boldsymbol{\theta}_{m_1,d}|\alpha_{m_1}) \left( \prod_{n=1}^{N_d} \sum_{z_{m_1,d,n}} \mathcal{P}(w_{d,n}|z_{m_1,d,n}, \mathbf{B}_{m_1}) \mathcal{P}(z_{m_1,d,n}|\boldsymbol{\theta}_{m_1,d}) \right) d\boldsymbol{\theta}_{m_1,d} \right] \quad [\text{A1.5}]$$

## Appendix 2: Variational Expectation Maximization estimation of a Latent Dirichlet Allocation

For the variational expansion of a Stage 1 LDA model  $\mathbf{m}_1$ ,  $\boldsymbol{\Gamma}_{\mathbf{m}_1}$  is an  $M \times k_{\mathbf{m}_1}$ -matrix akin to  $\boldsymbol{\theta}_{\mathbf{m}_1}$ , where row  $\mathbf{d}$  corresponds to document  $\mathbf{d}$ :  $\gamma_{\mathbf{m}_1,\mathbf{d},1} \dots \gamma_{\mathbf{m}_1,\mathbf{d},k_{\mathbf{m}_1}}$  (or  $\boldsymbol{\gamma}_{\mathbf{m}_1,\mathbf{d}}$ ), but contains the concentration parameters of a  $k_{\mathbf{m}_1}$ -dimension Dirichlet distribution and therefore are not constrained to sum to 1, and  $\boldsymbol{\Phi}_{\mathbf{m}_1}$  is an  $N \times k_{\mathbf{m}_1}$ -matrix whose rows correspond to the words across documents (indexed akin to  $\mathbf{w}_{\mathbf{d}}$  and  $\mathbf{z}_{\mathbf{d}}$ ) and whose columns correspond to the  $k_{\mathbf{m}_1}$  topics.  $\phi_{\mathbf{m}_1,\mathbf{d},n,i}$  is the probability that word  $n$  within document  $\mathbf{d}$  is from topic  $i$ , and  $\boldsymbol{\phi}_{\mathbf{m}_1,\mathbf{d},n;1\dots k_{\mathbf{m}_1}}$  (or  $\boldsymbol{\phi}_{\mathbf{m}_1,\mathbf{d},n}$ ) is a  $k_{\mathbf{m}_1}$ -length vector of probabilities defining the categorical distribution of that word's topic identity ( $\sum \phi_{\mathbf{m}_1,\mathbf{d},n} = 1$ ).  $\boldsymbol{\Phi}_{\mathbf{m}_1}$  contains  $M$  document-specific matrices (each is  $N_{\mathbf{d}} \times k_{\mathbf{m}_1}$  and notated  $\boldsymbol{\Phi}_{\mathbf{m}_1,\mathbf{d}}$ ). In comparison to  $\boldsymbol{\alpha}_{\mathbf{m}_1}$  and  $\mathbf{B}_{\mathbf{m}_1}$ , the variational parameters  $\boldsymbol{\Gamma}_{\mathbf{m}_1}$  and  $\boldsymbol{\Phi}_{\mathbf{m}_1}$  are document-specific and not coupled, thereby allowing estimation.

For a specific document  $\mathbf{d}$ , the variational distribution ( $\mathcal{Q}$ ) is

$$\mathcal{Q}(\boldsymbol{\theta}_{\mathbf{m}_1,\mathbf{d}}, \mathbf{z}_{\mathbf{m}_1,\mathbf{d}} | \boldsymbol{\gamma}_{\mathbf{m}_1,\mathbf{d}}, \boldsymbol{\Phi}_{\mathbf{m}_1,\mathbf{d}}) = \mathcal{Q}(\boldsymbol{\theta}_{\mathbf{m}_1,\mathbf{d}} | \boldsymbol{\gamma}_{\mathbf{m}_1,\mathbf{d}}) \prod_{n=1}^{N_{\mathbf{d}}} \mathcal{Q}(z_{\mathbf{m}_1,\mathbf{d},n} | \boldsymbol{\phi}_{\mathbf{m}_1,\mathbf{d},n}) \quad [\text{A2.1}]$$

In the Expectation (“E”) Step in the VEM algorithm, the distribution  $\mathcal{Q}(\boldsymbol{\theta}_{\mathbf{m}_1,\mathbf{d}}, \mathbf{z}_{\mathbf{m}_1,\mathbf{d}} | \boldsymbol{\gamma}_{\mathbf{m}_1,\mathbf{d}}, \boldsymbol{\Phi}_{\mathbf{m}_1,\mathbf{d}})$  is used to find a tight lower bound on  $\mathcal{P}(\mathbf{w}_{\mathbf{d}} | \boldsymbol{\alpha}_{\mathbf{m}_1}, \mathbf{B}_{\mathbf{m}_1})$  by optimizing the variational parameters  $\boldsymbol{\gamma}_{\mathbf{m}_1,\mathbf{d}}$  and  $\boldsymbol{\Phi}_{\mathbf{m}_1,\mathbf{d}}$  (*i.e.*, finding  $\boldsymbol{\gamma}_{\mathbf{m}_1,\mathbf{d}}^*$  and  $\boldsymbol{\Phi}_{\mathbf{m}_1,\mathbf{d}}^*$ , where the asterisks notate optimal values) with respect to minimizing the Kullback-Leibler Divergence ( $D_{\text{KL}}$ ) between  $\mathcal{Q}(\boldsymbol{\theta}_{\mathbf{m}_1,\mathbf{d}}, \mathbf{z}_{\mathbf{m}_1,\mathbf{d}} | \boldsymbol{\gamma}_{\mathbf{m}_1,\mathbf{d}}, \boldsymbol{\Phi}_{\mathbf{m}_1,\mathbf{d}})$  and  $\mathcal{P}(\boldsymbol{\theta}_{\mathbf{m}_1}, \mathbf{z}_{\mathbf{m}_1} | \mathbf{w}, \boldsymbol{\alpha}_{\mathbf{m}_1}, \mathbf{B}_{\mathbf{m}_1})$ :

$$(\boldsymbol{\gamma}_{\mathbf{m}_1,\mathbf{d}}^*, \boldsymbol{\Phi}_{\mathbf{m}_1,\mathbf{d}}^*) = \arg \min_{\boldsymbol{\gamma}_{\mathbf{m}_1,\mathbf{d}}, \boldsymbol{\Phi}_{\mathbf{m}_1,\mathbf{d}}} D_{\text{KL}} \left( \mathcal{Q}(\boldsymbol{\theta}_{\mathbf{m}_1,\mathbf{d}}, \mathbf{z}_{\mathbf{m}_1,\mathbf{d}} | \boldsymbol{\gamma}_{\mathbf{m}_1,\mathbf{d}}, \boldsymbol{\Phi}_{\mathbf{m}_1,\mathbf{d}}) \| \mathcal{P}(\boldsymbol{\theta}_{\mathbf{m}_1}, \mathbf{z}_{\mathbf{m}_1} | \mathbf{w}, \boldsymbol{\alpha}_{\mathbf{m}_1}, \mathbf{B}_{\mathbf{m}_1}) \right) \quad [\text{A2.2}]$$

Minimization of the distance is achieved through an iterative fixed-point method, where the derivative of  $D_{\text{KL}}$  is set to zero, producing a pair of update equations (Blei *et al.* 2003). First, the parameters describing the topic allocation of each word ( $\phi_{\mathbf{m}_1,\mathbf{d},n,i}$ ) are updated based on the topic distribution for the document ( $\boldsymbol{\gamma}_{\mathbf{m}_1,\mathbf{d}}$ ):

$$\phi_{\mathbf{m}_1,\mathbf{d},n,i} \propto \beta_{\mathbf{m}_1,\mathbf{d},w_{\mathbf{d},n},i} e^{\mathbb{E}_{\mathcal{Q}}[\log(\theta_{\mathbf{m}_1,\mathbf{d},i}) | \boldsymbol{\gamma}_{\mathbf{m}_1,\mathbf{d}}]} \quad [\text{A2.3}]$$

where  $\mathbb{E}_{\mathcal{Q}}$  is the expected value of the (log-scale) topic probability and is calculated using the digamma function ( $\Psi$ ), which is the logarithmic derivative of the gamma function ( $\Psi(a) = \frac{d}{da} \log(\Gamma(a))$ ), a quantity that is calculated through Taylor approximation:

$$\mathbb{E}_{\mathcal{Q}}[\log(\theta_{\mathbf{m}_1,\mathbf{d},i}) | \boldsymbol{\gamma}_{\mathbf{m}_1,\mathbf{d}}] = \Psi(\gamma_{\mathbf{m}_1,\mathbf{d},i}) - \Psi\left(\sum \gamma_{\mathbf{m}_1,\mathbf{d}}\right) \quad [\text{A2.4}]$$

Then, the parameters describing the topic distribution for the document ( $\gamma_{\mathbf{m}_1,\mathbf{d},i}$ ) are updated based on the word-level topic distributions for the sample ( $\boldsymbol{\Phi}_{\mathbf{m}_1,\mathbf{d}}$ ):

$$\gamma_{\mathbf{m}_1,\mathbf{d},i} = \alpha_{\mathbf{m}_1} + \sum_{n=1}^{N_{\mathbf{d}}} \phi_{\mathbf{m}_1,\mathbf{d},n,i} \quad [\text{A2.5}]$$

The update equations are alternated until the bound converges (*i.e.*, the updates do not yield changes to the

parameters), at which point the document-specific variation parameters have been optimized ( $\mathbf{V}_{m_1,d}^*$  and  $\Phi_{m_1,d}^*$  have been found) for the set of main parameters ( $\alpha_{m_1}, \mathbf{B}_{m_1}$ ).

The Maximization (“M”) Step in the VEM algorithm maximizes the overall lower bound with respect to the main model parameters  $\alpha_{m_1}$  and  $\mathbf{B}_{m_1}$  given the optimal variational parameters (Blei *et al.* 2003). This corresponds to obtaining maximum likelihood values of the model parameters using expected sufficient statistics for each sample under the approximate posterior calculated in the E Step (Blei *et al.* 2003). The update for the topic-level term distribution ( $\mathbf{B}_{m_1}$ ) is, analytically:

$$\beta_{m_1, i_{m_1} v} \propto \sum_{d=1}^M \sum_{n=1}^{N_d} \phi_{m_1, d, n}^* \omega_{d, n}^v \quad [\text{A2.7}]$$

where  $\omega_{d, n}^v$  is an indicator variable based on the term identity ( $v$ ) of the observed word ( $w_{d, n}$ ):

$$\omega_{d, n}^v = \begin{cases} 1, & w_{d, n} = v \\ 0, & w_{d, n} \neq v \end{cases} \quad [\text{A2.8}]$$

The update for the concentration parameter underlying the document-level topic distribution ( $\alpha_{m_1}$ ) requires an iterative approach to find a stationary point estimate. The optimization is conducted using the Newton-Raphson method (Ronning 1989), which repeats

$$\alpha_{m_1}^{\text{new}} = \alpha_{m_1}^{\text{old}} - \text{DD}(\alpha_{m_1}^{\text{old}})^{-1} \text{D}(\alpha_{m_1}^{\text{old}}) \quad [\text{A2.9}]$$

until convergence, where  $\text{D}$  represents the Jacobian matrix of first derivatives of a multivariate function and  $\text{DD}$  is the Hessian matrix of second derivatives of the function. Having updated the main model parameters (the M-Step), a new iteration of the E-Step followed by the M-Step is conducted, and the E-Step and M-Step are alternated until  $Q$  converges.

### Appendix 3: Softmax regression

We accommodate the multivariate proportional responses in our time series model by using multinomial logistic regression, also known as Softmax regression. The model is “log-linear” in that it relates the log of the expected proportion ( $\log E[\bar{\gamma}_{m_1^*,d_s,i}^*]_{m_2^r}$ ) to the linear predictors ( $\mathbf{x}_{m_2,d_s}^r \boldsymbol{\eta}_{m_2,s,i}^r$ ), although we formulate our model as the expected proportion ( $E[\bar{\gamma}_{m_1^*,d_s,i}^*]_{m_2^r}$ ) being a function of exponentiated predictors  $e^{\mathbf{x}_{m_2,d_s}^r \boldsymbol{\eta}_{m_2,s,i}^r}$ . And we handle the sum-to-1 constraint by normalizing  $e^{\mathbf{x}_{m_2,d_s}^r \boldsymbol{\eta}_{m_2,s,i}^r}$  with a document-specific partition function  $f_{\text{part}_{m_2^r,d_s}}$ :

$$E[\bar{\gamma}_{m_1^*,d_s,i}^*]_{m_2^r} = \frac{e^{\mathbf{x}_{m_2,d_s}^r \boldsymbol{\eta}_{m_2,s,i}^r}}{f_{\text{part}_{m_2^r,d_s}}} \quad [\text{A3.1}]$$

Given the sum-to-1 constraint ( $\sum_{i=1}^{k_{m_1^*}} \bar{\gamma}_{m_1^*,d_s,i}^* = \sum_{i=1}^{k_{m_1^*}} \frac{e^{\mathbf{x}_{m_2,d_s}^r \boldsymbol{\eta}_{m_2,s,i}^r}}{f_{\text{part}_{m_2^r,d_s}}} = 1$ ), we can simply define  $f_{\text{part}_{m_2^r,d_s}} = \sum_{j=1}^{k_{m_1^*}} e^{\mathbf{x}_{m_2,d_s}^r \boldsymbol{\eta}_{m_2,s,j}^r}$ , where we replace the  $i$  with  $j$  to avoid confusion with the focal topic  $i$ . This produces a generalized equation that is often referred to as the softmax function:

$$E[\bar{\gamma}_{m_1^*,d_s,i}^*]_{m_2^r} = \frac{e^{\mathbf{x}_{m_2,d_s}^r \boldsymbol{\eta}_{m_2,s,i}^r}}{\sum_{j=1}^{k_{m_1^*}} e^{\mathbf{x}_{m_2,d_s}^r \boldsymbol{\eta}_{m_2,s,j}^r}} = \text{softmax}\left(i, \mathbf{x}_{m_2,d_s}^r \boldsymbol{\eta}_{m_2,s,1}^r, \dots, \mathbf{x}_{m_2,d_s}^r \boldsymbol{\eta}_{m_2,s,k_{m_1^*}}^r\right) \quad [\text{A3.2}]$$

However, because of the sum-to-1 constraint, only  $k_{m_1^*} - 1$  of the proportions ( $\bar{\gamma}_{m_1^*,d_s,i}^*$  in  $\bar{\gamma}_{m_1^*,d_s,1}^* \dots \bar{\gamma}_{m_1^*,d_s,k_{m_1^*}}^*$ ), and by extension only  $k_{m_1^*} - 1$  of the parameter vectors ( $\boldsymbol{\eta}_{m_2,s,i}^r$  in  $\boldsymbol{\eta}_{m_2,s,1}^r \dots \boldsymbol{\eta}_{m_2,s,k_{m_1^*}}^r$ ), are uniquely identifiable. Thus, we define an augmented parameter vectors  $\hat{\boldsymbol{\eta}}_{m_2,s,i}^r$  (in  $\hat{\boldsymbol{\eta}}_{m_2,s,1}^r \dots \hat{\boldsymbol{\eta}}_{m_2,s,k_{m_1^*}}^r$ ), where

$$\hat{\boldsymbol{\eta}}_{m_2,s,i}^r = \boldsymbol{\eta}_{m_2,s,i}^r - \boldsymbol{\eta}_{m_2,s,1}^r \quad [\text{A3.3}]$$

setting the parameters associated with the first topic ( $i = 1$ ) to 0 ( $\hat{\boldsymbol{\eta}}_{m_2,s,1}^r = 0$ ,  $e^{\mathbf{x}_{m_2,d_s}^r \boldsymbol{\eta}_{m_2,s,1}^r} = 1$ ). This reduces the number of free parameter vectors (and number of proportions estimated) by 1 to the  $k_{m_1^*} - 1$  ( $\hat{\boldsymbol{\eta}}_{m_2,s,2}^r \dots \hat{\boldsymbol{\eta}}_{m_2,s,k_{m_1^*}}^r$ ) we are able to fit for this specific chunk of time  $s_{m_2^r}$ , resulting in the modified probability equation

$$E[\bar{\gamma}_{m_1^*,d_s,i}^*]_{m_2^r} = \frac{e^{\mathbf{x}_{m_2,d_s}^r \hat{\boldsymbol{\eta}}_{m_2,s,i}^r}}{\sum_{j=1}^{k_{m_1^*}} e^{\mathbf{x}_{m_2,d_s}^r \hat{\boldsymbol{\eta}}_{m_2,s,j}^r}} = \text{softmax}\left(i, \mathbf{x}_{m_2,d_s}^r \hat{\boldsymbol{\eta}}_{m_2,s,1}^r, \dots, \mathbf{x}_{m_2,d_s}^r \hat{\boldsymbol{\eta}}_{m_2,s,k_{m_1^*}}^r\right) \quad [\text{A3.4}]$$

We combine all of the  $k_{m_1^*}$  parameter vectors  $\hat{\boldsymbol{\eta}}_{m_2,s,1}^r$  to  $\hat{\boldsymbol{\eta}}_{m_2,s,k_{m_1^*}}^r$  (including the vector of 0s in  $\hat{\boldsymbol{\eta}}_{m_2,s,1}^r$ ) into a matrix  $\hat{\mathbf{H}}_{m_2,s}^r$ , which has  $k_{m_1^*}$  columns and a number of rows equal to the number of coefficients in model  $m_2$  ( $C_{m_2}$ ) including the intercept (*i.e.*, the length of  $\mathbf{x}_{m_2,d_s}^r$ ).

This allows us to further condense the expected probability equation to

$$\mathbb{E}[\bar{\gamma}_{m_1^*, d_s, i}]_{m_2^r} = \text{softmax}(i, \mathbf{x}_{m_2^r, d_s} \dot{\mathbf{H}}_{m_2^r, s}) \quad [\text{A3.5}]$$

thereby facilitating use of the generalized linear modeling framework. We expand the model to predict the proportions across all of the  $k_{m_1^*}$  topics within the document, which means we can drop the  $i$  input and produce the full set of values from the softmax function, which is a length-  $k_{m_1^*}$  row vector corresponding to the topic distribution of a single document:

$$\mathbb{E}[\bar{\gamma}_{m_1^*, d_s}]_{m_2^r} = \text{softmax}(\mathbf{x}_{m_2^r, d_s} \dot{\mathbf{H}}_{m_2^r, s}) \quad [\text{A3.6}]$$

We then expand the model across all documents within the chunk of time

$$\mathbb{E}[\bar{\mathbf{T}}_{m_1^*, s}]_{m_2^r} = \text{softmax}(\mathbf{X}_{m_2^r, s} \dot{\mathbf{H}}_{m_2^r, s}) \quad [\text{A3.7}]$$

where the covariates are held in a matrix  $(\mathbf{X}_{m_2^r, s})$  with the number of columns equal to the number of coefficients  $(C_{m_2})$  and the number of rows equal to the number of documents in the chunk  $(M_{m_2, s})$ . That is,  $\mathbf{X}_{m_2^r, s}$  is a series of  $\mathbf{x}_{m_2^r, d_s}$  row vectors. This equation relates directly to the generalized linear modeling equation that is typically written as  $\mathbf{g}(\mathbf{Y}) = \mathbf{X}\boldsymbol{\beta}$  or  $\mathbb{E}[\mathbf{Y}] = \mathbf{g}^{-1}(\mathbf{X}\boldsymbol{\beta})$ , where  $\mathbf{g}$  is the so-called link function and  $\mathbf{g}^{-1}$  is the inverse link function (McCullagh and Nelder 1989).

#### Appendix 4: Derivation of the Jacobian for the loss function

To construct the Jacobian of the loss function (which generates the negative prior-penalized log-likelihood ( $\mathcal{L}_{m_2,s} = -\log(\mathcal{P}(\bar{\Gamma}_{m_1,s}^* | \rho_{m_2,s}^r, \dot{\mathbf{H}}_{m_2,s}^r) \mathcal{P}(\dot{\mathbf{H}}_{m_2,s}^r | \rho_{m_2,s}^r))$ ), we first recognize that the loss for the segment is (weighted) additive across the documents within the chunk:

$$\mathcal{L}_{m_2,s} = - \sum_{d_s=1}^{M_{m_2,s}^r} u_{d_s} \left[ -\log \left( \mathcal{P}(\bar{\Gamma}_{m_1,s,d_s}^* | \rho_{m_2,s}^r, \dot{\mathbf{H}}_{m_2,s}^r) \mathcal{P}(\dot{\mathbf{H}}_{m_2,s}^r | \rho_{m_2,s}^r) \right) \right] \quad [\text{A4.1}]$$

We next acknowledge the order of operations of the component multivariate functions that comprise  $\mathcal{L}_{m_2,s}$ . To reduce notational clutter, we rename the cross entropy, softmax, matrix multiplication, and penalty functions  $f_{\text{CE}}$ ,  $f_{\text{S}}$ ,  $f_{\text{MM}}$ , and  $f_{\text{P}}$ , as well as rename the nested functions: the softmax of the matrix multiplication and the cross entropy of the softmax of the matrix multiplication become  $f_{\text{S}(\text{MM})}$  and  $f_{\text{CE}(\text{S}(\text{MM}))}$ . This allows us to write the loss equation for all documents within a time chunk as

$$\mathcal{L}_{m_2,s} = - \sum_{d_s=1}^{M_{m_2,s}^r} u_{d_s} \left[ \sum_{i=1}^{k_{m_1}^*} f_{\text{CE}(\text{S}(\text{MM}))}(\bar{\gamma}_{m_1,d_s,i}^*, \mathbf{x}_{m_2,d_s}^r, \dot{\mathbf{H}}_{m_2,s}^r) + f_{\text{P}}(\dot{\mathbf{H}}_{m_2,s}^r) \right] \quad [\text{A4.2}]$$

which can be condensed via the summation across topics to

$$\mathcal{L}_{m_2,s} = - \sum_{d_s=1}^{M_{m_2,s}^r} u_{d_s} \left( f_{\text{CE}(\text{S}(\text{MM}))}(\bar{\gamma}_{m_1,d_s}^*, \mathbf{x}_{m_2,d_s}^r, \dot{\mathbf{H}}_{m_2,s}^r) + f_{\text{P}}(\dot{\mathbf{H}}_{m_2,s}^r) \right) \quad [\text{A4.3}]$$

This highlights the chained (nested) aspect of the non-penalty functions (the cross entropy is calculated using the output of the softmax, which uses the output of the matrix multiplication), whose derivative can be expanded using the multivariate chain rule. For two general functions  $\mathbf{f}$  and  $\mathbf{g}$  chained as  $\mathbf{f}(\mathbf{g}(\mathbf{a}))$  (where  $\mathbf{a}$  contains the multivariate input values), we can write the function composite using the ring operator as  $(\mathbf{f} \circ \mathbf{g})(\mathbf{a})$ . We then take the multivariate derivative (denoted as function  $\mathcal{D}$ ) of the composite:

$$\begin{aligned} \mathcal{D}(\mathbf{f}(\mathbf{g}(\mathbf{a}))) &= \mathcal{D}((\mathbf{f} \circ \mathbf{g})(\mathbf{a})) \\ &= (\mathcal{D}(\mathbf{f}) \circ \mathbf{g})(\mathbf{a}) \cdot \mathcal{D}(\mathbf{g}(\mathbf{a})) \\ &= \mathcal{D}(\mathbf{f})(\mathbf{g}(\mathbf{a})) \cdot \mathcal{D}(\mathbf{g}(\mathbf{a})) \\ &= \mathcal{D}(\mathbf{f})(\mathbf{g}(\mathbf{a})) \cdot \mathcal{D}(\mathbf{g})(\mathbf{a}) \end{aligned} \quad [\text{A4.4}]$$

where  $\cdot$  is the dot product operator. Thus, the derivative of  $\mathbf{f}$  of  $\mathbf{g}$  of  $\mathbf{a}$  is the dot product of the derivative of  $\mathbf{f}$  evaluated at  $\mathbf{g}$  of  $\mathbf{a}$  and the derivative of  $\mathbf{g}$  evaluated at  $\mathbf{a}$ . Using the chain rule, we now expand the derivative of the loss function applied to a specific document within a specific chunk ( $d_s$ ). We start by expanding the outer layers ( $f_{\text{CE}}$  and  $f_{\text{S}}$ ):

$$\mathcal{D} \left( f_{\text{CE}(\text{S}(\text{MM}))}(\bar{\gamma}_{m_1,d_s,i}^*, \mathbf{x}_{m_2,d_s}^r, \dot{\mathbf{H}}_{m_2,s}^r) \right) = \mathcal{D}(f_{\text{CE}})(\bar{\gamma}_{m_1,d_s}^*, f_{\text{S}(\text{MM})}(\mathbf{x}_{m_2,d_s}^r, \dot{\mathbf{H}}_{m_2,s}^r)) \cdot \mathcal{D}(f_{\text{S}(\text{MM})}(\mathbf{x}_{m_2,d_s}^r, \dot{\mathbf{H}}_{m_2,s}^r)) \quad [\text{A4.5}]$$

We next expand the inner layers ( $f_{\text{S}}$  and  $f_{\text{MM}}$ ) by working with the right-hand-side of the dot product:

$$\mathcal{D}\left(f_{S(\text{MM})}(\mathbf{x}_{m_2^r, d_s}, \mathbf{H}_{m_2^r, s})\right) = \mathcal{D}(f_S)\left(f_{\text{MM}}(\mathbf{x}_{m_2^r, d_s}, \mathbf{H}_{m_2^r, s})\right) \cdot \mathcal{D}(f_{\text{MM}})(\mathbf{x}_{m_2^r, d_s}, \mathbf{H}_{m_2^r, s}) \quad [\text{A4.6}]$$

Combining these chained results gives the full Jacobian for the non-penalized component of the loss function applied to a single document:

$$\begin{aligned} \mathcal{D}\left(f_{\text{CE}(S(\text{MM}))}(\bar{\mathbf{y}}_{m_1^*, d_s, i}^*, \mathbf{x}_{m_2^r, d_s}, \mathbf{H}_{m_2^r, s})\right) \\ = \mathcal{D}(f_{\text{CE}})\left(\bar{\mathbf{y}}_{m_1^*, d_s}^*, f_{S(\text{MM})}(\mathbf{x}_{m_2^r, d_s}, \mathbf{H}_{m_2^r, s})\right) \cdot \mathcal{D}(f_S)\left(f_{\text{MM}}(\mathbf{x}_{m_2^r, d_s}, \mathbf{H}_{m_2^r, s})\right) \cdot \mathcal{D}(f_{\text{MM}})(\mathbf{x}_{m_2^r, d_s}, \mathbf{H}_{m_2^r, s}) \end{aligned} \quad [\text{A4.7}]$$

which is the dot product between [1] the dot product between [a] derivative of the cross-entropy function (evaluated at the softmax of the matrix multiplication of the coefficients) and [b] the derivative of the softmax (evaluated at the matrix multiplication of the coefficients), and [2] the derivative of the matrix multiplication evaluated at the coefficients.

We now define the derivative matrices (Jacobians) of each of the singular functions  $f_{\text{MM}}$ ,  $f_S$ , and  $f_{\text{CE}}$  and the composite functions  $f_{S(\text{MM})}$  and  $f_{\text{CE}(S(\text{MM}))}$ . To aid in this, we consider that the Jacobian of a function  $\mathbf{f}$  contains the partial derivatives of each output ( $f_i$ , a general component of  $\mathbf{f}$ ) with respect to each input ( $a_j$ , a general component of  $\mathbf{a}$ ), which can be written generally as  $\frac{\partial f_i}{\partial a_j}$  or  $\mathcal{D}_j f_i$ . The Jacobian for a given function then maps the input to the output, and so has dimensions equal to the number of output classes  $\times$  the number of input classes.

The function  $f_{\text{MM}}$  maps the  $\mathbf{H}_{m_2^r, s}$  matrix ( $C_{m_2} \times k_{m_1^*}$ ) to the dimensions of  $\mathbf{x}_{m_2^r, d_s} \mathbf{H}_{m_2^r, s}$  ( $1 \times k_{m_1^*}$ ) by left-multiplying  $\mathbf{H}_{m_2^r, s}$  by the covariate row matrix  $\mathbf{x}_{m_2^r, d_s}$  ( $1 \times C_{m_2}$ ). Thus, its Jacobian has  $k_{m_1^*}$  rows and  $C_{m_2} k_{m_1^*}$  columns:

$$\mathcal{D}(f_{\text{MM}}) = \begin{bmatrix} \mathcal{D}_1 f_{\text{MM}_1} & \cdots & \mathcal{D}_{C_{m_2} k_{m_1^*}} f_{\text{MM}_1} \\ \vdots & \ddots & \vdots \\ \mathcal{D}_1 f_{\text{MM}_{k_{m_1^*}}} & \cdots & \mathcal{D}_{C_{m_2} k_{m_1^*}} f_{\text{MM}_{k_{m_1^*}}} \end{bmatrix} \quad [\text{A4.8}]$$

For notekeeping purposes, entry  $c, i$  in the coefficient matrix  $\mathbf{H}_{m_2^r, s}$  ( $\dot{H}_{m_2^r, s, c, i}$ ) corresponds to the column  $(i-1)C_{m_2} + c$  in the Jacobian. In effect, the coefficient matrix  $\mathbf{H}_{m_2^r, s}$  is linearized in column-major order (iterating through all covariates within a given topic before progressing to the next topic). Recall that the matrix multiplication used to generate an output element (row)  $f_{\text{MM}_l}$  (for  $l$  in  $1 \dots k_{m_1^*}$ ) is just a linear combination of components

$$f_{\text{MM}_l} = x_{m_2^r, d_s, 1} \dot{H}_{m_2^r, s, 1, l} + x_{m_2^r, d_s, 2} \dot{H}_{m_2^r, s, 2, l} + \cdots + x_{m_2^r, d_s, C_{m_2}} \dot{H}_{m_2^r, s, C_{m_2}, l} \quad [\text{A4.9}]$$

and therefore, the partial derivative of the output element  $l_{m_1^*}$  with respect to an input element  $c_{m_2^r}, i_{m_1^*}$  is simply the relevant covariate or 0 (when beyond the relevant part of the Jacobian):

$$\mathcal{D}_{c, i} f_{\text{MM}_l} = \begin{cases} x_{m_2^r, d_s, c}, & i = l \\ 0, & i \neq l \end{cases} \quad [\text{A4.10}]$$

Moving to the softmax function,  $f_S$  maps  $\mathbf{x}_{m_2^r, d_s} \mathbf{H}_{m_2^r, s}$  to  $\text{E}[\bar{\mathbf{y}}_{m_1^*, d_s}^*]_{m_2^r}$ , both of which are of dimension  $1 \times k_{m_1^*}$ , because we are working within a single document. Thus, its Jacobian has  $k_{m_1^*}$  rows and  $k_{m_1^*}$  columns:



$$\mathcal{D}(f_s) = \begin{bmatrix} \mathcal{D}_1 f_{s_1} & \cdots & \mathcal{D}_{k_{m_1^*}} f_{s_1} \\ \vdots & \ddots & \vdots \\ \mathcal{D}_1 f_{s_{k_{m_1^*}}} & \cdots & \mathcal{D}_{k_{m_1^*}} f_{s_{k_{m_1^*}}} \end{bmatrix} \quad [\text{A4.11}]$$

We can write a generalized equation for the entries by describing the partial derivative of output  $j$  with respect to input  $i$ ,  $\mathcal{D}_i f_{s_j}$ :

$$\begin{aligned} \mathcal{D}_i f_{s_j} &= \frac{\partial f_{s_j}}{\partial (\mathbf{x}_{m_2, d_s}^r \dot{\mathbf{H}}_{m_2, s})_i} \\ &= \frac{\frac{\partial}{\partial \sum_{l=1}^{k_{m_1^*}} e^{\mathbf{x}_{m_2, d_s}^r \dot{\mathbf{H}}_{m_2, s, l}}} e^{\mathbf{x}_{m_2, d_s}^r \dot{\mathbf{H}}_{m_2, s, j}}}{\partial (\mathbf{x}_{m_2, d_s}^r \dot{\mathbf{H}}_{m_2, s})_i} \\ &= \frac{\partial}{\partial (\mathbf{x}_{m_2, d_s}^r \dot{\mathbf{H}}_{m_2, s})_i} \frac{e^{\mathbf{x}_{m_2, d_s}^r \dot{\mathbf{H}}_{m_2, s, j}}}{\sum_{l=1}^{k_{m_1^*}} e^{\mathbf{x}_{m_2, d_s}^r \dot{\mathbf{H}}_{m_2, s, l}}} \end{aligned} \quad [\text{A4.12}]$$

We decompose the generalized entry using the quotient rule, where for a function  $\mathbf{f}(\mathbf{a})$  that is equal to the ratio of two other functions:  $\mathbf{f}(\mathbf{a}) = \frac{\mathbf{g}(\mathbf{a})}{\mathbf{h}(\mathbf{a})}$ , the derivative of the function is

$$\mathcal{D}(\mathbf{f}(\mathbf{a})) = \frac{\mathcal{D}(\mathbf{g}(\mathbf{a}))\mathbf{h}(\mathbf{a}) - \mathcal{D}(\mathbf{h}(\mathbf{a}))\mathbf{g}(\mathbf{a})}{[\mathbf{h}(\mathbf{a})]^2} \quad [\text{A4.13}]$$

Here,  $\mathbf{g} = e^{\mathbf{x}_{m_2, d_s}^r \dot{\mathbf{H}}_{m_2, s, j}}$  and  $\mathbf{h} = \sum_{l=1}^{k_{m_1^*}} e^{\mathbf{x}_{m_2, d_s}^r \dot{\mathbf{H}}_{m_2, s, l}}$  and we differentiate each with respect to  $(\mathbf{x}_{m_2, d_s}^r \dot{\mathbf{H}}_{m_2, s})_i$ :

$$\mathcal{D}_i f_{s_j} = \frac{\frac{\partial e^{\mathbf{x}_{m_2, d_s}^r \dot{\mathbf{H}}_{m_2, s, j}}}{\partial (\mathbf{x}_{m_2, d_s}^r \dot{\mathbf{H}}_{m_2, s})_i} \sum_{l=1}^{k_{m_1^*}} e^{\mathbf{x}_{m_2, d_s}^r \dot{\mathbf{H}}_{m_2, s, l}} - \frac{\partial \sum_{l=1}^{k_{m_1^*}} e^{\mathbf{x}_{m_2, d_s}^r \dot{\mathbf{H}}_{m_2, s, l}}}{\partial (\mathbf{x}_{m_2, d_s}^r \dot{\mathbf{H}}_{m_2, s})_i} e^{\mathbf{x}_{m_2, d_s}^r \dot{\mathbf{H}}_{m_2, s, j}}}{\left[ \sum_{l=1}^{k_{m_1^*}} e^{\mathbf{x}_{m_2, d_s}^r \dot{\mathbf{H}}_{m_2, s, l}} \right]^2} \quad [\text{A4.14}]$$

Regardless of the specific input  $i$  that we are computing the partial derivative for  $\mathbf{h}$  with respect to, the value will always be  $e^{\mathbf{x}_{m_2, d_s}^r \dot{\mathbf{H}}_{m_2, s, i}}$ :

$$\frac{\sum_{l=1}^{k_{m_1^*}} e^{\mathbf{x}_{m_2, d_s}^r \dot{\mathbf{H}}_{m_2, s, l}}}{\partial (\mathbf{x}_{m_2, d_s}^r \dot{\mathbf{H}}_{m_2, s})_i} = e^{\mathbf{x}_{m_2, d_s}^r \dot{\mathbf{H}}_{m_2, s, i}} \quad [\text{A4.15}]$$

For  $\mathbf{g}$ , however, the value of the partial derivative is 0 unless  $i = j$ , in which case it is  $e^{\mathbf{x}_{m_2, d_s}^r \dot{\mathbf{H}}_{m_2, s, j}}$ :

$$\frac{\partial e^{\mathbf{x}_{m_2, d_s}^r \dot{\mathbf{H}}_{m_2, s, j}}}{\partial (\mathbf{x}_{m_2, d_s}^r \dot{\mathbf{H}}_{m_2, s})_i} = \begin{cases} e^{\mathbf{x}_{m_2, d_s}^r \dot{\mathbf{H}}_{m_2, s, j}} & i = j \\ 0 & i \neq j \end{cases} \quad [\text{A4.16}]$$

Thus, when  $i = j$ ,

$$\begin{aligned}
\mathcal{D}_i f_{S_j} &= \frac{e^{x_{m_2,d_s}^r \dot{\mathbf{h}}_{m_2,s,j}^r} \sum_{l=1}^{k_{m_1}^*} e^{x_{m_2,d_s}^r \dot{\mathbf{h}}_{m_2,s,l}^r} - e^{x_{m_2,d_s}^r \dot{\mathbf{h}}_{m_2,s,i}^r} e^{x_{m_2,d_s}^r \dot{\mathbf{h}}_{m_2,s,j}^r}}{\left[ \sum_{l=1}^{k_{m_1}^*} e^{x_{m_2,d_s}^r \dot{\mathbf{h}}_{m_2,s,l}^r} \right]^2} \\
&= \frac{e^{x_{m_2,d_s}^r \dot{\mathbf{h}}_{m_2,s,j}^r} \left( \sum_{l=1}^{k_{m_1}^*} e^{x_{m_2,d_s}^r \dot{\mathbf{h}}_{m_2,s,l}^r} - e^{x_{m_2,d_s}^r \dot{\mathbf{h}}_{m_2,s,i}^r} \right)}{\left[ \sum_{l=1}^{k_{m_1}^*} e^{x_{m_2,d_s}^r \dot{\mathbf{h}}_{m_2,s,l}^r} \right]^2} \\
&= \frac{e^{x_{m_2,d_s}^r \dot{\mathbf{h}}_{m_2,s,j}^r}}{\sum_{l=1}^{k_{m_1}^*} e^{x_{m_2,d_s}^r \dot{\mathbf{h}}_{m_2,s,l}^r}} \left( \frac{\sum_{l=1}^{k_{m_1}^*} e^{x_{m_2,d_s}^r \dot{\mathbf{h}}_{m_2,s,l}^r}}{\sum_{l=1}^{k_{m_1}^*} e^{x_{m_2,d_s}^r \dot{\mathbf{h}}_{m_2,s,l}^r}} - \frac{e^{x_{m_2,d_s}^r \dot{\mathbf{h}}_{m_2,s,i}^r}}{\sum_{l=1}^{k_{m_1}^*} e^{x_{m_2,d_s}^r \dot{\mathbf{h}}_{m_2,s,l}^r}} \right) \\
&= \frac{e^{x_{m_2,d_s}^r \dot{\mathbf{h}}_{m_2,s,j}^r}}{\sum_{l=1}^{k_{m_1}^*} e^{x_{m_2,d_s}^r \dot{\mathbf{h}}_{m_2,s,l}^r}} \left( 1 - \frac{e^{x_{m_2,d_s}^r \dot{\mathbf{h}}_{m_2,s,i}^r}}{\sum_{l=1}^{k_{m_1}^*} e^{x_{m_2,d_s}^r \dot{\mathbf{h}}_{m_2,s,l}^r}} \right) \\
&= f_{S_j} (1 - f_{S_i})
\end{aligned} \tag{A4.17}$$

Similarly, when  $i \neq j$ ,

$$\begin{aligned}
\mathcal{D}_i f_{S_j} &= \frac{0 - e^{x_{m_2,d_s}^r \dot{\mathbf{h}}_{m_2,s,i}^r} e^{x_{m_2,d_s}^r \dot{\mathbf{h}}_{m_2,s,j}^r}}{\left[ \sum_{l=1}^{k_{m_1}^*} e^{x_{m_2,d_s}^r \dot{\mathbf{h}}_{m_2,s,l}^r} \right]^2} \\
&= - \frac{e^{x_{m_2,d_s}^r \dot{\mathbf{h}}_{m_2,s,i}^r}}{\sum_{l=1}^{k_{m_1}^*} e^{x_{m_2,d_s}^r \dot{\mathbf{h}}_{m_2,s,l}^r}} \frac{e^{x_{m_2,d_s}^r \dot{\mathbf{h}}_{m_2,s,j}^r}}{\sum_{l=1}^{k_{m_1}^*} e^{x_{m_2,d_s}^r \dot{\mathbf{h}}_{m_2,s,l}^r}} \\
&= -f_{S_i} f_{S_j} \\
&= -f_{S_j} f_{S_i}
\end{aligned} \tag{A4.18}$$

Combining these conditions, we have

$$\mathcal{D}_i f_{S_j} = \begin{cases} f_{S_j} (1 - f_{S_i}) & i = j \\ -f_{S_j} f_{S_i} & i \neq j \end{cases} \tag{A4.19}$$

We can use the Kronecker delta function to condense the conditional equation to

$$\mathcal{D}_i f_{S_j} = f_{S_j} (\delta_{ij} - f_{S_i}) \tag{A4.20}$$

where

$$\delta_{ij} = \begin{cases} 1 & i = j \\ 0 & i \neq j \end{cases} \tag{A4.21}$$

The function set  $f_{S(\text{MM})}(\mathbf{x}_{m_2,d_s}^r, \dot{\mathbf{H}}_{m_2,s}^r)$  maps the  $\dot{\mathbf{H}}_{m_2,s}^r$  matrix ( $C_{m_2} \times k_{m_1}^*$ ) to the dimensions of

$\mathbb{E}[\bar{\mathbf{y}}_{m_1^*, d_s}^*]_{m_2^r}$  ( $1 \times k_{m_1^*}$ ) and so, like  $f_{\text{MM}}$ , its Jacobian has  $k_{m_1^*}$  rows and  $C_{m_2} k_{m_1^*}$  columns. We combine the Jacobians of  $f_S$  and  $f_{\text{MM}}$  to define the Jacobian of  $f_{S(\text{MM})}(\mathbf{x}_{m_2^r, d_s}, \dot{\mathbf{H}}_{m_2^r, s})$ :

$$\mathcal{D} \left( f_{S(\text{MM})}(\mathbf{x}_{m_2^r, d_s}, \dot{\mathbf{H}}_{m_2^r, s}) \right) = \begin{bmatrix} \mathcal{D}_1 f_{S(\text{MM})}(\mathbf{x}_{m_2^r, d_s}, \dot{\mathbf{H}}_{m_2^r, s})_1 & \cdots & \mathcal{D}_{C_{m_2} k_{m_1^*}} f_{S(\text{MM})}(\mathbf{x}_{m_2^r, d_s}, \dot{\mathbf{H}}_{m_2^r, s})_1 \\ \vdots & \ddots & \vdots \\ \mathcal{D}_1 f_{S(\text{MM})}(\mathbf{x}_{m_2^r, d_s}, \dot{\mathbf{H}}_{m_2^r, s})_{k_{m_1^*}} & \cdots & \mathcal{D}_{C_{m_2} k_{m_1^*}} f_{S(\text{MM})}(\mathbf{x}_{m_2^r, d_s}, \dot{\mathbf{H}}_{m_2^r, s})_{k_{m_1^*}} \end{bmatrix} \quad [\text{A4.22}]$$

For a general entry, the partial derivative is

$$\mathcal{D}_{c,i} f_{S(\text{MM})}(\mathbf{x}_{m_2^r, d_s}, \dot{\mathbf{H}}_{m_2^r, s})_j = \sum_{l=1}^{k_{m_1^*}} \mathcal{D}_l f_{S_j} \cdot \mathcal{D}_{c,i} f_{\text{MM}_l} \quad [\text{A4.23}]$$

Recalling that  $\mathcal{D}_l f_{S_j} \cdot \mathcal{D}_{c,i} f_{\text{MM}_l}$  is 0 except for when  $i = l$  (when it is  $x_{d_s, c, i}$ ), we can simplify this equation to

$$\mathcal{D}_{c,i} f_{S(\text{MM})}(\mathbf{x}_{m_2^r, d_s}, \dot{\mathbf{H}}_{m_2^r, s})_j = \mathcal{D}_i f_{S_j} x_{m_2^r, d_s, c, i} \quad [\text{A4.24}]$$

And recalling that  $\mathcal{D}_i f_{S_j} = f_{S_j}(\delta_{ij} - f_{S_i})$ , we can write this equation as

$$\mathcal{D}_{c,i} f_{S(\text{MM})}(\mathbf{x}_{m_2^r, d_s}, \dot{\mathbf{H}}_{m_2^r, s})_j = f_{S_j}(\delta_{ij} - f_{S_i}) x_{m_2^r, d_s, c, i} \quad [\text{A4.25}]$$

The cross-entropy function,  $f_{\text{CE}}$ , maps  $\mathbb{E}[\bar{\mathbf{y}}_{m_1^*, d_s}^*]_{m_2^r}$  (dimension  $1 \times k_{m_1^*}$ ) to the cross-entropy (loss) for the document, which is a scalar value. Thus, the Jacobian of  $f_{\text{CE}}$  is of dimension  $1 \times k_{m_1^*}$ :

$$\mathcal{D}(f_{\text{CE}}) = [\mathcal{D}_1 f_{\text{CE}} \quad \cdots \quad \mathcal{D}_{k_{m_1^*}} f_{\text{CE}}] \quad [\text{A4.26}]$$

A general entry in the Jacobian  $\mathcal{D}_l f_{\text{CE}}$  (*i.e.*, the partial derivative of the cross entropy loss with respect to topic  $l$ 's probability) is

$$\mathcal{D}_l f_{\text{CE}} = - \sum_{i=1}^{k_{m_1^*}} \frac{\partial \bar{y}_{m_1^*, d_s, i}^* \log \left( \mathbb{E}[\bar{y}_{m_1^*, d_s, i}^*]_{m_2^r} \right)}{\partial \mathbb{E}[\bar{y}_{m_1^*, d_s, l}^*]_{m_2^r}} \quad [\text{A4.27}]$$

Notably, the only instance where  $\mathbb{E}[\bar{y}_{m_1^*, d_s, l}^*]_{m_2^r}$  appears in the function being derived is when  $l = i$ , in which case, the derivative is

$$\mathcal{D}_{l=i} f_{\text{CE}} = - \frac{\bar{y}_{m_1^*, d_s, i}^*}{\mathbb{E}[\bar{y}_{m_1^*, d_s, l}^*]_{m_2^r}} \quad [\text{A4.28}]$$

Otherwise (*i.e.*, when  $l \neq i$ ), the function being derived is a constant and therefore has a derivative = 0.

Combining these conditions, we have

$$\mathcal{D}_l f_{\text{CE}} = \begin{cases} -\frac{\bar{\gamma}_{m_1^*, d_s, i}^*}{\mathbb{E}[\bar{\gamma}_{m_1^*, d_s, l}^*]_{m_2^r}} & l = i \\ 0 & l \neq i \end{cases} \quad [\text{A4.29}]$$

For notation, we identify the element  $l = i$  as  $\tilde{l}$ .

We can verify that the dimensionalities of the Jacobians are proper for combination via dot products:  $\mathcal{D}(f_{\text{CE}})(\bar{\gamma}_{m_1^*, d_s}^*, \mathbf{x}_{m_2^r, d_s}, f_{\text{S(MM)}}(\mathbf{x}_{m_2^r, d_s}, \hat{\mathbf{H}}_{m_2^r, s})) \cdot \mathcal{D}(f_{\text{S}})(f_{\text{MM}}(\mathbf{x}_{m_2^r, d_s}, \hat{\mathbf{H}}_{m_2^r, s})) \cdot \mathcal{D}(f_{\text{MM}})(\mathbf{x}_{m_2^r, d_s}, \hat{\mathbf{H}}_{m_2^r, s})$ .  $\mathcal{D}(f_{\text{CE}})$  is  $1 \times k_{m_1^*}$ ,  $\mathcal{D}(f_{\text{S}})$  is  $k_{m_1^*} \times k_{m_1^*}$ , and  $\mathcal{D}(f_{\text{MM}})$  is  $k_{m_1^*} \times C_{m_2} k_{m_1^*}$ . Thus, each of the two dot products has proper component matrices. In addition, the resulting matrix is  $1 \times C_{m_2} k_{m_1^*}$ , which heuristically matches the fact that the composite of the three functions maps the set of parameters ( $C_{m_2} k_{m_1^*}$  in total) to a single scalar value of the cross-entropy loss. Having verified the dimensions, we combine the elements across the three Jacobians to determine the derivative of  $f_{\text{CE}}$  of  $f_{\text{S}}$  of  $f_{\text{MM}}$ :  $\mathcal{D}(f_{\text{CE(S(MM))}})(\bar{\gamma}_{m_1^*, d_s}^*, \mathbf{x}_{m_2^r, d_s}, \hat{\mathbf{H}}_{m_2^r, s})$ .  $f_{\text{CE(S(MM))}}(\bar{\gamma}_{m_1^*, d_s}^*, \mathbf{x}_{m_2^r, d_s}, \hat{\mathbf{H}}_{m_2^r, s})$  maps the  $C_{m_2} k_{m_1^*}$  parameters (each  $\dot{\eta}_{s_{m_2^r, c_{m_2^r, i_{m_1^*}}}}$  entry in  $\hat{\mathbf{H}}_{m_2^r, s}$ ) to a scalar output (cross-entropy loss), so the resulting Jacobian is of dimensions  $1 \times C_{m_2} k_{m_1^*}$ :

$$\begin{aligned} & \mathcal{D}(f_{\text{CE(S(MM))}})(\bar{\gamma}_{m_1^*, d_s}^*, \mathbf{x}_{m_2^r, d_s}, \hat{\mathbf{H}}_{m_2^r, s}) \\ &= [\mathcal{D}_1 f_{\text{CE(S(MM))}}(\bar{\gamma}_{m_1^*, d_s}^*, \mathbf{x}_{m_2^r, d_s}, \hat{\mathbf{H}}_{m_2^r, s}) \quad \cdots \quad \mathcal{D}_{C_{m_2} k_{m_1^*}} f_{\text{CE(S(MM))}}(\bar{\gamma}_{m_1^*, d_s}^*, \mathbf{x}_{m_2^r, d_s}, \hat{\mathbf{H}}_{m_2^r, s})] \end{aligned} \quad [\text{A4.30}]$$

Similar to the rows in  $\mathcal{D}(f_{\text{MM}})$ , the single row of partial derivatives in  $\mathcal{D}(f_{\text{CE(S(MM))}})(\bar{\gamma}_{m_1^*, d_s}^*, \mathbf{x}_{m_2^r, d_s}, \hat{\mathbf{H}}_{m_2^r, s})$  corresponds to the column-major order linearized  $\hat{\mathbf{H}}_{m_2^r, s}$ . Following the indexing of  $\mathcal{D}(f_{\text{MM}})$ , we will index  $\mathcal{D}(f_{\text{CE(S(MM))}})(\bar{\gamma}_{m_1^*, d_s}^*, \mathbf{x}_{m_2^r, d_s}, \hat{\mathbf{H}}_{m_2^r, s})$  with  $c$  and  $i$  (as  $c, i$ ), where  $c, i$  refers to column (element)  $c C_{m_2} + i$  in  $\mathcal{D}(f_{\text{CE(S(MM))}})(\bar{\gamma}_{m_1^*, d_s}^*, \mathbf{x}_{m_2^r, d_s}, \hat{\mathbf{H}}_{m_2^r, s})$ . For a general entry  $c, i$ , then

$$\mathcal{D}_{c, i} f_{\text{CE(S(MM))}}(\bar{\gamma}_{m_1^*, d_s}^*, \mathbf{x}_{m_2^r, d_s}, \hat{\mathbf{H}}_{m_2^r, s}) = \sum_{j=1}^{k_{m_1^*}} \mathcal{D}_j(f_{\text{CE}})\left(\mathbb{E}[\bar{\gamma}_{m_1^*, d_s}^*]_{m_2^r}\right) \cdot \mathcal{D}_{c i} f_{\text{S(MM)}}(j, \mathbf{x}_{m_2^r, d_s}, \hat{\mathbf{H}}_{m_2^r, s}) \quad [\text{A4.31}]$$

Since only the  $\tilde{l}$  element of  $\mathcal{D}_j f_{\text{CE}}\left(\mathbb{E}[\bar{\gamma}_{m_1^*, d_s}^*]_{m_2^r}\right)$  is non-0, in which case it is  $-\frac{\bar{\gamma}_{m_1^*, d_s, \tilde{l}}^*}{\mathbb{E}[\bar{\gamma}_{m_1^*, d_s, \tilde{l}}^*]_{m_2^r}}$ , we can simplify this equation to

$$\mathcal{D}_{c, i} f_{\text{CE(S(MM))}}(\bar{\gamma}_{m_1^*, d_s}^*, \mathbf{x}_{m_2^r, d_s}, \hat{\mathbf{H}}_{m_2^r, s}) = -\frac{\bar{\gamma}_{m_1^*, d_s, \tilde{l}}^*}{\mathbb{E}[\bar{\gamma}_{m_1^*, d_s, \tilde{l}}^*]_{m_2^r}} \cdot \mathcal{D}_{c i} f_{\text{S(MM)}}(\tilde{l}, \mathbf{x}_{m_2^r, d_s}, \hat{\mathbf{H}}_{m_2^r, s}) \quad [\text{A4.32}]$$

Substituting in the derivative of the softmax (of the matrix multiplication) of the parameters,

$$\mathcal{D}_{c, i} f_{\text{CE(S(MM))}}(\bar{\gamma}_{m_1^*, d_s}^*, \mathbf{x}_{m_2^r, d_s}, \hat{\mathbf{H}}_{m_2^r, s}) = -\frac{\bar{\gamma}_{m_1^*, d_s, \tilde{l}}^*}{\mathbb{E}[\bar{\gamma}_{m_1^*, d_s, \tilde{l}}^*]_{m_2^r}} \cdot f_{\text{S} \tilde{i}}(\delta_{\tilde{i} \tilde{i}} - f_{\text{S} \tilde{i}}) x_{m_2^r, d_s, c, i} \quad [\text{A4.33}]$$

Noting that, by our definition,  $f_{S_{\bar{i}}} = \mathbb{E}[\bar{\gamma}_{m_1^*, d_s, \bar{i}}]_{m_2^r}$ , we can further simplify this equation:

$$\begin{aligned} \mathcal{D}_{c,i} f_{\text{CE(S(MM))}}(\bar{\gamma}_{m_1^*, d_s}^*, \mathbf{x}_{m_2^r, d_s}, \dot{\mathbf{H}}_{m_2^r, s}) &= -\frac{\bar{\gamma}_{m_1^*, d_s, i}^*}{f_{S_{\bar{i}}}} \cdot f_{S_{\bar{i}}}(\delta_{i\bar{i}} - f_{S_{\bar{i}}})x_{m_2^r, d_s, c, i} \\ &= -\bar{\gamma}_{m_1^*, d_s, i}^* \cdot (\delta_{i\bar{i}} - f_{S_{\bar{i}}})x_{m_2^r, d_s, c, i} \\ &= \bar{\gamma}_{m_1^*, d_s, i}^* \cdot (f_{S_{\bar{i}}} - \delta_{i\bar{i}})x_{m_2^r, d_s, c, i} \end{aligned} \quad [\text{A4.34}]$$

This row matrix  $\mathcal{D}(f_{\text{CE(S(MM))}}(\bar{\gamma}_{m_1^*, d_s}^*, \mathbf{x}_{m_2^r, d_s}, \dot{\mathbf{H}}_{m_2^r, s}))$  holds the  $C_{m_2} k_{m_1}^*$  partial derivatives associated with the  $C_{m_2}$  parameters and the  $k_{m_1}^*$  topics. With these functions, we can calculate the full set of partial derivatives required to evaluate the Jacobian for the loss equation without the penalty, the remaining component to be added.

Following the sum rule, given that we need to calculate the partial derivative of the loss with respect to each of the  $C_{m_2} k_{m_1}^*$  model parameters  $\dot{\mathbf{H}}_{m_2^r, s}$ , we start with

$$\mathcal{L}_{m_2^r, s}^* = - \sum_{d_s=1}^{M_{m_2^r, s}} u_{d_s} \left( f_{\text{CE(S(MM))}}(\bar{\gamma}_{m_1^*, d_s}^*, \mathbf{x}_{m_2^r, d_s}, \dot{\mathbf{H}}_{m_2^r, s}) + f_{\text{P}}(\dot{\mathbf{H}}_{m_2^r, s}) \right) \quad [\text{A4.35}]$$

focus in on a specific document  $d_s$ ,

$$\mathcal{L}_{m_2^r, d_s}^* = - \left( f_{\text{CE(S(MM))}}(\bar{\gamma}_{m_1^*, d_s}^*, \mathbf{x}_{m_2^r, d_s}, \dot{\mathbf{H}}_{m_2^r, s}) + f_{\text{P}}(\dot{\mathbf{H}}_{m_2^r, s}) \right) \quad [\text{A4.36}]$$

and take the derivative with respect to a particular parameter  $ci$  in the set,

$$\mathcal{D}_{c,i}(\mathcal{L}_{m_2^r, d_s}^*) = \mathcal{D}_{c,i} f_{\text{CE(S(MM))}}(\bar{\gamma}_{m_1^*, d_s}^*, \mathbf{x}_{m_2^r, d_s}, \dot{\mathbf{H}}_{m_2^r, s}) + \mathcal{D}_{c,i} f_{\text{P}}(\dot{\mathbf{H}}_{m_2^r, s}) \quad [\text{A4.37}]$$

Therefore, we need to calculate the partial derivative of the penalty function with respect to a particular parameter  $c, i$ . Remembering that the penalty function is,

$$f_{\text{P}}(\dot{\mathbf{H}}_{m_2^r, s}) = \sum_{i=\bar{i}}^{k_{m_1}^*} \sum_{c=1}^{C_{m_2}} \lambda_{m_2}(\dot{\eta}_{m_2^r, s, i, c})^2 \quad [\text{A4.38}]$$

the partial derivative of the penalty with respect to  $c, i$  is

$$\mathcal{D}_{c,i} f_{\text{P}}(\dot{\mathbf{H}}_{m_2^r, s}) = \sum_{i=\bar{i}}^{k_{m_1}^*} \sum_{c=1}^{C_{m_2}} \frac{\partial \lambda_{m_2}(\dot{\eta}_{m_2^r, s, i, c})^2}{\partial (\dot{\eta}_{m_2^r, s, i, c})} \quad [\text{A4.39}]$$

which evaluates to

$$\mathcal{D}_{c,i} f_{\text{P}}(\dot{\mathbf{H}}_{m_2^r, s}) = \sum_{i=\bar{i}}^{k_{m_1}^*} \sum_{c=1}^{C_{m_2}} 2\lambda_{m_2} \dot{\eta}_{m_2^r, s, i, c} \quad [\text{A4.40}]$$

Given that we also know that  $\mathcal{D}_{c,i}(\mathcal{L}_{m_2,d_s}^*) = \mathcal{D}_{c,i} f_{\text{CE(S(MM))}}(\bar{\mathbf{y}}_{m_1,d_s}^*, \mathbf{x}_{m_2,d_s}^r, \hat{\mathbf{H}}_{m_2,s}) + \mathcal{D}_{c,i} f_{\text{p}}(\hat{\mathbf{H}}_{m_2,s})$  and  $\mathcal{D}_{c,i} f_{\text{CE(S(MM))}}(\bar{\mathbf{y}}_{m_1,d_s}^*, \mathbf{x}_{m_2,d_s}^r, \hat{\mathbf{H}}_{m_2,s}) = \bar{y}_{m_1,d_s,i}^* \cdot (f_{s_i} - \delta_{ii}) x_{m_2,d_s,c,i}^r$ , the partial derivative of penalized loss with respect to input  $ci$  ( $\mathcal{D}_{c,i}(\mathcal{L}_{m_2,d_s}^*)$ ) is now fully defined:

$$\mathcal{D}_{c,i}(\mathcal{L}_{m_2,d_s}^*) = \bar{y}_{m_1,d_s,i}^* \cdot (f_{s_i} - \delta_{ii}) x_{m_2,d_s,c,i}^r + \sum_{i=i}^{k_{m_1}^*} \sum_{c=1}^{C_{m_2}} 2\lambda \dot{\eta}_{m_2,s,i,c}^r \quad [\text{A4.41}]$$

This equation determines the gradient of the loss across with respect to each input (parameter-topic combination  $c, i$ ) within a single document. Acknowledging multiple documents fall under the same parameter combination, we simply sum across all  $d_{s_{m_2}^r}$ :

$$\mathcal{D}_{c,i}(\mathcal{L}_{m_2,s}^*) = \sum_{d_s}^{M_{m_2,s}^r} \left[ \bar{y}_{m_1,d_s,i}^* \cdot (f_{s_i} - \delta_{ii}) x_{m_2,d_s,c,i}^r + \sum_{i=i}^{k_{m_1}^*} \sum_{c=1}^{C_{m_2}} 2\lambda \dot{\eta}_{m_2,s,i,c}^r \right] \quad [\text{A4.42}]$$

and distribute the summation

$$\mathcal{D}_{c,i}(\mathcal{L}_{m_2,s}^*) = \sum_{d_s}^{M_{m_2,s}^r} \bar{y}_{m_1,d_s,i}^* \cdot (f_{s_i} - \delta_{ii}) x_{m_2,d_s,c,i}^r + \sum_{d_s=1}^{M_{m_2,s}^r} \sum_{i=i}^{k_{m_1}^*} \sum_{c=1}^{C_{m_2}} 2\lambda \dot{\eta}_{m_2,s,i,c}^r \quad [\text{A4.43}]$$

to achieve the completely defined general entry to the Jacobian for the penalized loss equation used to fit a multinomial model to a chunk of documents' topic proportions.

## Appendix 5: Segment Function Mapper

The function  $f_{f(s)}$  is used to map a function to segments of data and has three inputs: [1] the function to map, [2] the matrix multiplication to map it onto, and [3] the indication matrix that defines the mapping. Specifically here, we use  $\Xi_{m_2}$  to map the **softmax** function to the matrix multiplication  $\mathbf{X}_{m_2} \hat{\mathbf{H}}_{m_2}$ :

$$\mathbb{E}[\bar{\Gamma}_{m_1}]_{m_2^r} = f_{f(s)}(\text{softmax}, \mathbf{X}_{m_2}^r \hat{\mathbf{H}}_{m_2}^r, \Xi_{m_2}^r) \quad [\text{A5.1}]$$

The operation of  $f_{f(s)}$  can be considered algorithmically:

For each column in  $\Xi_{m_2}^r$ , which corresponds to a segment  $s$  (in  $1 \dots S_{m_2}$ )

Select the rows from  $\mathbf{X}_{m_2}^r$  where  $\xi_{m_2,s}^r = 1$  (the document is in the segment):  $\mathbf{X}_{m_2,s}^r$

Select columns  $C_{m_2}(s-1)+1$  to  $C_{m_2}s$  from  $\hat{\mathbf{H}}_{m_2}^r$  (the regressors are for the segment):  $\hat{\mathbf{H}}_{m_2,s}^r$

Apply the **softmax** function to the selected entries:  $\mathbf{X}_{m_2,s}^r \hat{\mathbf{H}}_{m_2,s}^r$