
Bayesian Poisson Tucker Decomposition for Learning the Structure of International Relations

Aaron Schein

University of Massachusetts Amherst

ASCHEIN@CS.UMASS.EDU

Mingyuan Zhou

University of Texas at Austin

MINGYUAN.ZHOU@MCCOMBS.UTEXAS.EDU

David M. Blei

Columbia University

DAVID.BLEI@COLUMBIA.EDU

Hanna Wallach

Microsoft Research New York City

WALLACH@MICROSOFT.COM

Abstract

We introduce Bayesian Poisson Tucker decomposition (BPTD) for modeling country–country interaction event data. These data consist of interaction events of the form “country i took action a toward country j at time t .” BPTD discovers overlapping country–community memberships, including the number of latent communities. In addition, it discovers directed community–community interaction networks that are specific to “topics” of action types and temporal “regimes.” We show that BPTD yields an efficient MCMC inference algorithm and achieves better predictive performance than related models. We also demonstrate that it discovers interpretable latent structure that agrees with our knowledge of international relations.

1. Introduction

Like their inhabitants, countries interact with one another: they consult, negotiate, trade, threaten, and fight. These interactions are seldom uncoordinated. Rather, they are connected by a fabric of overlapping communities, such as security coalitions, treaties, trade cartels, and military alliances. For example, OPEC coordinates the petroleum export policies of its thirteen member countries, LAIA fosters trade among Latin American countries, and NATO guarantees collective defense against attacks by external parties.

A single country can belong to multiple communities, reflecting its different identities. For example, Venezuela—an oil-producing country and a Latin American country—is a member of both OPEC and LAIA. When Venezuela interacts with other countries, it sometimes does so as an OPEC member and sometimes does so as a LAIA member.

Countries engage in both within-community and between-community interactions. For example, when acting as an OPEC member, Venezuela consults with other OPEC countries, but trades with non-OPEC, oil-importing countries. Moreover, although Venezuela engages in between-community interactions when trading as an OPEC member, it engages in within-community interactions when trading as a LAIA member. To understand or predict how countries interact, we must account for their community memberships and how those memberships influence their actions.

In this paper, we take a new approach to learning overlapping communities from interaction events of the form “country i took action a toward country j at time t .” A data set of such interaction events can be represented as either 1) a set of event tokens, 2) a tensor of event type counts, or 3) a series of weighted multinetworks. Models that use the token representation naturally yield efficient inference algorithms, models that use the tensor representation exhibit good predictive performance, and models that use the network representation learn latent structure that aligns with well-known concepts such as communities. Previous models of interaction event data have each used a subset of these representations. Our approach—Bayesian Poisson Tucker decomposition (BPTD)—takes advantage of all three.

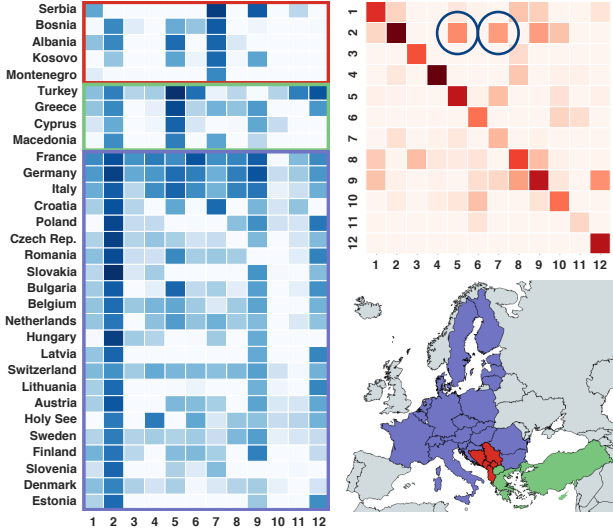


Figure 1. Latent structure learned by BPTD from country–country interaction events between 1995 and 2000. *Top right*: A community–community interaction network specific to a single topic of action types and temporal regime. The topic places most of its mass on the *Intend to Cooperate* and *Consult* actions, so this network represents cooperative community–community interactions. The two strongest between-community interactions (circled) are $2 \rightarrow 5$ and $2 \rightarrow 7$. *Left*: Each row depicts the overlapping community memberships for a single country. We show only those countries whose strongest community membership is to either community 2, 5, or 7. We ordered the countries accordingly. Countries strongly associated with community 7 are at highlighted in red; countries associated with community 5 are highlighted in green; and countries associated with community 2 are highlighted in purple. *Bottom right*: Each country is colored according to its strongest community membership. The latent communities have a very strong geographic interpretation.

BPTD builds on the classic Tucker decomposition (Tucker, 1964) to factorize a tensor of event type counts into three factor matrices and a four-dimensional core tensor (section 2). The factor matrices embed countries into communities, action types into “topics,” and time steps into “regimes.” The core tensor interacts communities, topics, and regimes. The country–community factors enable BPTD to learn overlapping community memberships, while the core tensor enables it to learn directed community–community interaction networks specific to topics of action types and temporal regimes. Figure 1 illustrates this structure. BPTD leads to an efficient MCMC inference algorithm (section 4) and achieves better predictive performance than related models (section 6). Finally, BPTD discovers interpretable latent structure that agrees with our knowledge of international relations (section 7).

2. Bayesian Poisson Tucker Decomposition

We can represent a data set of interaction events as a set of N event tokens, where a single token $e_n = (i \xrightarrow{a} j, t)$ indicates that sender country $i \in [V]$ took action $a \in [A]$ toward receiver country $j \in [V]$ during time step $t \in [T]$. Alternatively, we can aggregate these event tokens into a four-dimensional tensor \mathbf{Y} , where element $y_{i \xrightarrow{a} j}^{(t)}$ is a count of the number of events of type $(i \xrightarrow{a} j, t)$. This tensor will be sparse because most event types never actually occur in practice. Finally, we can equivalently view this count tensor as a series of T weighted multinet snapshots, where the weight on edge $i \xrightarrow{a} j$ in the t^{th} snapshot is $y_{i \xrightarrow{a} j}^{(t)}$.

BPTD models each element of count tensor \mathbf{Y} as

$$y_{i \xrightarrow{a} j}^{(t)} \sim \text{Po} \left(\sum_{c=1}^C \theta_{ic} \sum_{d=1}^C \theta_{jd} \sum_{k=1}^K \phi_{ak} \sum_{r=1}^R \psi_{tr} \lambda_{c \xrightarrow{k} d}^{(r)} \right), \quad (1)$$

where θ_{ic} , θ_{jd} , ϕ_{ak} , ψ_{tr} , and $\lambda_{c \xrightarrow{k} d}^{(r)}$ are positive real numbers. Factors θ_{ic} and θ_{jd} capture the rates at which countries i and j participate in communities c and d , respectively; factor ϕ_{ak} captures the strength of association between action a and topic k ; and ψ_{tr} captures how well regime r explains the events in time step t . We can collectively view the $V \times C$ country–community factors as a latent factor matrix Θ , where the i^{th} row represents country i ’s community memberships. Similarly, we can view the $A \times K$ action–topic factors and the $T \times R$ time-step–regime factors as latent factor matrices Φ and Ψ , respectively. Factor $\lambda_{c \xrightarrow{k} d}^{(r)}$ captures the rate at which community c takes actions associated with topic k toward community d during regime r . The $C \times C \times K \times R$ such factors form a core tensor Λ that interacts communities, topics, and regimes.

The country–community factors are gamma-distributed,

$$\theta_{ic} \sim \Gamma(\alpha_i, \beta_i), \quad (2)$$

where the shape and rate parameters α_i and β_i are specific to country i . We place an uninformative gamma prior over these shape and rate parameters: $\alpha_i, \beta_i \sim \Gamma(\epsilon_0, \epsilon_0)$. This hierarchical prior enables BPTD to express heterogeneity in the countries’ rates of activity. For example, we expect that the US will engage in more interactions than Burundi.

The action–topic and time-step–regime factors are also gamma-distributed; however, we assume that these factors are drawn directly from an uninformative gamma prior,

$$\phi_{ak}, \psi_{tr} \sim \Gamma(\epsilon_0, \epsilon_0). \quad (3)$$

Because BPTD learns a single embedding of countries into communities, it preserves the traditional network-based notion of community membership. Any sender–receiver asymmetry is captured by the core tensor Λ , which we can

view as a compression of count tensor Y . By allowing on-diagonal elements, which we denote by $\lambda_{c \circlearrowleft k}^{(r)}$ and off-diagonal elements to be non-zero, the core tensor can represent both within- and between-community interactions.

The elements of the core tensor are gamma-distributed,

$$\lambda_{c \circlearrowleft k}^{(r)} \sim \Gamma(\eta_c^\circ \eta_c^{\leftrightarrow} \nu_k \rho_r, \delta) \quad (4)$$

$$\lambda_{c \xrightarrow{k} d}^{(r)} \sim \Gamma(\eta_c^{\leftrightarrow} \eta_d^{\leftrightarrow} \nu_k \rho_r, \delta) \quad c \neq d. \quad (5)$$

Each community $c \in [C]$ has two positive weights η_c° and η_c^{\leftrightarrow} that capture its rates of within- and between-community interaction, respectively. Each topic $k \in [K]$ has a positive weight ν_k , while each regime $r \in [R]$ has a positive weight ρ_r . We place an uninformative prior over the within-community interaction rates and gamma shrinkage priors over the other weights: $\eta_c^\circ \sim \Gamma(\epsilon_0, \epsilon_0)$, $\eta_c^{\leftrightarrow} \sim \Gamma(\gamma_0 / C, \zeta)$, $\nu_k \sim \Gamma(\gamma_0 / K, \zeta)$, and $\rho_r \sim \Gamma(\gamma_0 / R, \zeta)$. These priors bias BPTD toward learning latent structure that is sparse. Finally, we assume that δ and ζ are drawn from an uninformative gamma prior: $\delta, \zeta \sim \Gamma(\epsilon_0, \epsilon_0)$.

As $K \rightarrow \infty$, the topic weights and their corresponding action–topic factors constitute a draw $G_K = \sum_{k=1}^{\infty} \nu_k \mathbb{1}_{\phi_k}$ from a gamma process (Ferguson, 1973). Similarly, as $R \rightarrow \infty$, the regime weights and their corresponding time-step–regime factors constitute a draw $G_R = \sum_{r=1}^{\infty} \rho_r \mathbb{1}_{\psi_r}$ from another gamma process. As $C \rightarrow \infty$, the within- and between-community interaction weights and their corresponding country–community factors constitute a draw $G_C = \sum_{c=1}^{\infty} \eta_c^{\leftrightarrow} \mathbb{1}_{\theta_c}$ from a marked gamma process (Kingman, 1972). The mark associated with atom $\theta_c = (\theta_{1c}, \dots, \theta_{Vc})$ is η_c° . We can view the elements of the core tensor and their corresponding factors as a draw $G = \sum_{c=1}^{\infty} \sum_{d=1}^{\infty} \sum_{k=1}^{\infty} \sum_{r=1}^{\infty} \lambda_{c \xrightarrow{k} d}^{(r)} \mathbb{1}_{\theta_c, \theta_d, \phi_k, \psi_r}$ from a gamma process, provided that the expected sum of the core tensor elements is finite. This multirelational gamma process extends the relational gamma process of Zhou (2015).

Proposition 1: *In the limit as $C, K, R \rightarrow \infty$, the expected sum of the core tensor elements is finite and equal to*

$$\mathbb{E} \left[\sum_{c=1}^{\infty} \sum_{k=1}^{\infty} \sum_{r=1}^{\infty} \left(\lambda_{c \circlearrowleft k}^{(r)} + \sum_{d \neq c} \lambda_{c \xrightarrow{k} d}^{(r)} \right) \right] = \frac{1}{\delta} \left(\frac{\gamma_0^3}{\zeta^3} + \frac{\gamma_0^4}{\zeta^4} \right).$$

We prove this proposition in the supplementary material.

3. Connections to Previous Work

Poisson CP decomposition: DuBois & Smyth (2010) developed a model that assigns each event token (ignoring time steps) to one of Q latent classes, where each class $q \in [Q]$ is characterized by three categorical distributions— θ_q^{\rightarrow}

over senders, θ_q^{\leftarrow} over receivers, and ϕ_q over actions—i.e.,

$$P(e_n = (i \xrightarrow{a} j, t) \mid z_n = q) = \theta_{i_q}^{\rightarrow} \theta_{j_q}^{\leftarrow} \phi_{aq}. \quad (6)$$

This model is closely related to the Poisson-based model of Schein et al. (2015), which explicitly uses the canonical polyadic (CP) tensor decomposition (Harshman, 1970) to factorize count tensor Y into four latent factor matrices. These factor matrices jointly embed senders, receivers, action types, and time steps into a Q -dimensional space,

$$y_{i \xrightarrow{a} j}^{(t)} \sim \text{Po} \left(\sum_{q=1}^Q \theta_{i_q}^{\rightarrow} \theta_{j_q}^{\leftarrow} \phi_{aq} \psi_{tq} \right), \quad (7)$$

where $\theta_{i_q}^{\rightarrow}$, $\theta_{j_q}^{\leftarrow}$, ϕ_{aq} , and ψ_{tq} are positive real numbers.

Schein et al.’s model generalizes Bayesian Poisson matrix factorization (Cemgil, 2009; Gopalan et al., 2014; 2015; Zhou & Carin, 2015) and non-Bayesian Poisson CP decomposition (Chi & Kolda, 2012; Welling & Weber, 2001).

Although Schein et al.’s model is expressed in terms of a tensor of event type counts, the relationship between the multinomial and Poisson distributions (Kingman, 1972) means that we can also express it in terms of a set of event tokens. This yields an equation that is similar to equation 6,

$$P(e_n = (i \xrightarrow{a} j, t) \mid z_n = q) \propto \theta_{i_q}^{\rightarrow} \theta_{j_q}^{\leftarrow} \phi_{aq} \psi_{tq}. \quad (8)$$

Conversely, DuBois & Smyth’s model can be expressed as a CP tensor decomposition. This equivalence is analogous to the relationship between Poisson matrix factorization and latent Dirichlet allocation (Blei et al., 2003).

We can make Schein et al.’s model nonparametric by adding a per-class positive weight $\lambda_q \sim \Gamma(\frac{\gamma_0}{Q}, \zeta)$, i.e.,

$$y_{i \xrightarrow{a} j}^{(t)} \sim \text{Po} \left(\sum_{q=1}^Q \theta_{i_q}^{\rightarrow} \theta_{j_q}^{\leftarrow} \phi_{aq} \psi_{tq} \lambda_q \right). \quad (9)$$

As $Q \rightarrow \infty$ the per-class weights and their corresponding latent factors constitute a draw from a gamma process.

Adding this per-class weight reveals that CP decomposition is a special case of Tucker decomposition where the cardinalities of the latent dimensions are equal and the off-diagonal elements of the core tensor are zero. DuBois & Smyth’s and Schein et al.’s models are therefore highly constrained special cases of BPTD that cannot capture dimension-specific structure, such as communities of countries or topics of action types. These models require each latent class to jointly summarize information about senders, receivers, action types, and time steps. This requirement conflates communities of countries and topics of action types, thus forcing each class to capture potentially redundant information. Moreover, by definition, CP decomposition models cannot express between-community interactions and cannot express sender–receiver asymmetry without learning completely separate latent factor matrices for

senders and receivers. These limitations make it hard to interpret these models as learning community memberships.

Infinite relational models: The infinite relational model (IRM) of Kemp et al. (2006) also learns latent structure specific to each dimension of an M -dimensional tensor; however, unlike BPTD, the elements of this tensor are binary, indicating the presence or absence of the corresponding event type. The IRM therefore uses a Bernoulli likelihood. Schmidt & Mørup (2013) extended the IRM to model a tensor of event counts by replacing the Bernoulli likelihood with a Poisson likelihood (and gamma priors):

$$y_{i \rightarrow j}^{(t)} \sim \text{Po} \left(\lambda_{z_i \xrightarrow{z_a} z_j}^{(z_t)} \right), \quad (10)$$

where $z_i, z_j \in [C]$ are the respective community assignments of countries i and j , $z_a \in [K]$ is the topic assignment of action a , and $z_t \in [R]$ is the regime assignment of time step t . This model, which we refer to as the gamma–Poisson IRM (GPIRM), allocates M -dimensional event types to M -dimensional latent classes—e.g., it allocates all tokens of type $(i \xrightarrow{a} j, t)$ to class $(z_i \xrightarrow{z_a} z_j, z_t)$.

The GPIRM is a special case of BPTD where the rows of the latent factor matrices are constrained to be “one-hot” binary vectors—i.e., $\theta_{ic} = \mathbb{1}(z_i = c)$, $\theta_{jd} = \mathbb{1}(z_j = d)$, $\phi_{ak} = \mathbb{1}(z_a = k)$, and $\psi_{tr} = \mathbb{1}(z_t = r)$. With this constraint, the Poisson rates in equations 1 and 10 are equal. Unlike BPTD, the GPIRM is a single-membership model. In addition, it cannot express heterogeneity in rates of activity of the countries, action types, and time steps. The latter limitation can be remedied by letting $\theta_{iz_i}, \theta_{jz_j}, \phi_{az_a}$, and ψ_{tz_t} be positive real numbers. We refer to this variant of the GPIRM as the degree-corrected GPIRM (DCGPIRM).

Stochastic block models: The IRM itself generalizes the stochastic block model (SBM) of Nowicki & Snijders (2001), which learns latent structure from binary networks. Although the SBM was originally specified using a Bernoulli likelihood, Karrer & Newman (2011) introduced an alternative specification that uses the Poisson likelihood:

$$y_{i \rightarrow j} \sim \text{Po} \left(\sum_{c=1}^C \theta_{ic} \sum_{d=1}^C \theta_{jd} \lambda_{c \rightarrow d} \right), \quad (11)$$

where $\theta_{ic} = \mathbb{1}(z_i = c)$, $\theta_{jd} = \mathbb{1}(z_j = d)$, and $\lambda_{c \rightarrow d}$ is a positive real number. Like the IRM and the GPIRM, the SBM is a single-membership model and cannot express heterogeneity in the countries’ rates of activity. Airolidi et al. (2008) addressed the former limitation by letting $\theta_{ic} \in [0, 1]$ such that $\sum_{c=1}^C \theta_{ic} = 1$. Meanwhile, Karrer & Newman (2011) addressed the latter limitation by allowing both θ_{iz_i} and θ_{jz_j} to be positive real numbers, much like the DCGPIRM. Ball et al. (2011) simultaneously addressed both limitations by letting $\theta_{ic}, \theta_{jd} \geq 0$, but constrained $\lambda_{c \rightarrow d} = \lambda_{d \rightarrow c}$. Finally, Zhou (2015) extended

Ball et al.’s model to be nonparametric and introduced the Poisson–Bernoulli distribution to link binary data to the Poisson likelihood in a principled fashion. In this model, the elements of the core matrix and their corresponding factors constitute a draw from a relational gamma process.

Non-Poisson Tucker decomposition: Researchers sometimes refer to the Poisson rate in equation 11 as being “bilinear” because it can equivalently be written as $\theta_j \Lambda \theta_i^\top$. Nickel et al. (2012) introduced RESCAL—a non-probabilistic bilinear model for binary data that achieves state-of-the-art performance at relation extraction. Nickel et al. (2015) then introduced several extensions for extracting relations of different types. Bilinear models, such as RESCAL and its extensions, are all special cases (albeit non-probabilistic ones) of Tucker decomposition.

Hoff (2015) recently developed a Gaussian-based Tucker decomposition model and multilinear tensor regression model (Hoff, 2014) for analyzing interaction event data.

Finally, there are many other Tucker decomposition methods (Kolda & Bader, 2009). Although these include non-parametric (Xu et al., 2012) and nonnegative variants (Kim & Choi, 20007; Mørup et al., 2008; Cichocki et al., 2009), BPTD is the first such model to use a Poisson likelihood.

4. Posterior Inference

Given an observed count tensor \mathbf{Y} , inference in BPTD involves “inverting” the generative process to obtain the posterior distribution over the parameters conditioned on \mathbf{Y} and hyperparameters ϵ_0 and γ_0 . The posterior distribution is analytically intractable; however, we can approximate it using a set of posterior samples. We draw these samples using Gibbs sampling, repeatedly resampling the value of each parameter from its conditional posterior given \mathbf{Y} , ϵ_0 , γ_0 , and the current values of the other parameters. We express each parameter’s conditional posterior in a closed form using gamma–Poisson conjugacy and the auxiliary variable techniques of Zhou & Carin (2012). We provide the conditional posteriors in the supplementary material.

The conditional posteriors depend on \mathbf{Y} via a set of “latent sources” (Cemgil, 2009) or subcounts. Because of the Poisson additivity theorem (Kingman, 1972), each latent source $y_{ic \xrightarrow{ak} jd}^{(tr)}$ is a Poisson-distributed random variable:

$$y_{ic \xrightarrow{ak} jd}^{(tr)} \sim \text{Po} \left(\theta_{ic} \theta_{jd} \phi_{ak} \psi_{tr} \lambda_{c \xrightarrow{k} d}^{(r)} \right) \quad (12)$$

$$y_{i \xrightarrow{a} j}^{(t)} = \sum_{c=1}^C \sum_{d=1}^D \sum_{k=1}^K \sum_{r=1}^R y_{ic \xrightarrow{ak} jd}^{(tr)}. \quad (13)$$

Together, equations 12 and 13 are equivalent to equation 1. In practice, we can equivalently view each latent source in

terms of the token representation described in section 2,

$$y_{ic \rightarrow ak \rightarrow jd}^{(tr)} = \sum_{n=1}^N \mathbb{1}(e_n = (i \xrightarrow{a} j, t)) \mathbb{1}(z_n = (c \xrightarrow{k} d, r)), \quad (14)$$

where each token’s class assignment z_n is an auxiliary latent variable. Using this representation, computing the latent sources (given the current values of the model parameters) simply involves allocating event tokens to classes, much like the inference algorithm for DuBois & Smyth’s model, and aggregating them using equation 14. The conditional posterior for each token’s class assignment is

$$\begin{aligned} P(z_n = (c \xrightarrow{k} d, r) | e_n = (i \xrightarrow{a} j, t), \mathbf{Y}, \epsilon_0, \gamma_0, \dots) \\ \propto \theta_{ic} \theta_{jd} \phi_{ak} \psi_{tr} \lambda_{c \xrightarrow{k} d}^{(r)}. \end{aligned} \quad (15)$$

Computation is dominated by the normalizing constant

$$Z_{i \xrightarrow{a} j}^{(t)} = \sum_{c=1}^C \sum_{d=1}^C \sum_{k=1}^K \sum_{r=1}^R \theta_{ic} \theta_{jd} \phi_{ak} \psi_{tr} \lambda_{c \xrightarrow{k} d}^{(r)}. \quad (16)$$

Computing this normalizing constant naïvely involves $O(C \times C \times K \times R)$ operations; however, because each latent class $(c \xrightarrow{k} d, r)$ is composed of four separate dimensions, we can improve efficiency. We instead compute

$$Z_{i \xrightarrow{a} j}^{(t)} = \sum_{c=1}^C \theta_{ic} \sum_{d=1}^C \theta_{jd} \sum_{k=1}^K \theta_{ak} \sum_{r=1}^R \psi_{tr} \lambda_{c \xrightarrow{k} d}^{(r)}, \quad (17)$$

which involves $O(C + C + K + R)$ operations.

Compositional allocation using equations 15 and 17 improves computational efficiency significantly over naïve non-compositional allocation using equations 15 and 16. In practice, we set C , K , and R to large values to approximate the nonparametric interpretation of BPTD. If, for example, $C = 50$, $K = 10$, and $R = 5$, computing the normalizing constant for equation 15 using equation 16 requires 2,753 times the number of operations implied by equation 17.

Proposition 2: *For an M -dimensional core tensor with $D_1 \times \dots \times D_M$ elements, computing the normalizing constant using non-compositional allocation requires $1 \leq \pi < \infty$ times the number of operations required to compute it using compositional allocation. When $D_1 = \dots = D_M = 1$, $\pi = 1$. As $D_m, D_{m'} \rightarrow \infty$ for any m and $m' \neq m$, $\pi \rightarrow \infty$.*

We prove this proposition in the supplementary material.

BPTD and other Poisson-based models yield allocation inference algorithms that take advantage of the inherent sparsity of the data and scale with the number of event tokens. In contrast, non-Poisson tensor decomposition models (including Hoff’s model) lead to algorithms that scale with the size of the count tensor. Allocation-based inference in BPTD is especially efficient because it *compositionally* allocates each M -dimensional event token to an

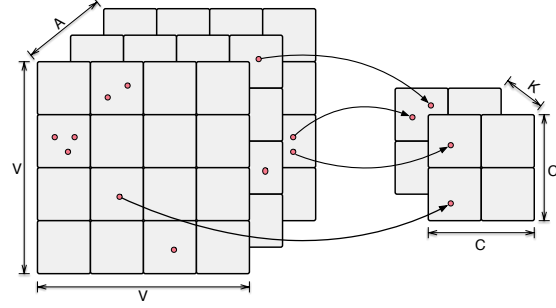


Figure 2. Compositional allocation. For clarity, we show the allocation process for a three-dimensional count tensor (ignoring time steps). Observed three-dimensional event tokens (left) are compositionally allocated to three-dimensional latent classes (right).

M -dimensional latent class. Figure 2 illustrates this process. CP decomposition models, such as those of DuBois & Smyth (2010) and Schein et al. (2015), only permit non-compositional allocation. For example, while BPTD allocates each token $e_n = (i \xrightarrow{a} j, t)$ to a four-dimensional latent class $(c \xrightarrow{k} d, r)$, Schein et al.’s model allocates e_n to a one-dimensional latent class q that cannot be decomposed. Therefore, when $Q = C \times C \times K \times R$, BPTD yields a faster allocation inference algorithm than Schein et al.’s model.

5. Country–Country Interaction Event Data

Our data come from the Integrated Crisis Early Warning System (ICEWS) of Boschee et al. and the Global Database of Events, Language, and Tone (GDELT) of Lee-taru & Schrodt (2013). ICEWS and GDELT both use the Conflict and Mediation Event Observations (CAMEO) hierarchy (Gerner et al.) for senders, receivers, and actions.

The top-level CAMEO coding for senders and receivers is their country affiliation, while lower levels in the hierarchy incorporate more specific attributes like their sectors (e.g., government or civilian) and their religious or ethnic affiliations. When studying international relations using CAMEO-coded event data, researchers usually consider only the senders’ and receivers’ countries. There are 249 countries represented in ICEWS, which include non-universally recognized states, such as *Occupied Palestinian Territory*, and former states, such as *Former Yugoslav Republic of Macedonia*; there are 233 countries in GDELT.

The top level for actions, which we use in our analyses, consists of twenty action classes, roughly ranked according to their overall sentiment. For example, the most negative is 20—*Use Unconventional Mass Violence*. CAMEO further divides these actions into the QuadClass scheme: Verbal Cooperation (actions 2–5), Material Cooperation (actions 6–7), Verbal Conflict (actions 8–16), and Material Conflict (16–20). The first action (1—*Make Statement*) is neutral.

6. Predictive Analysis

Baseline models: We compared BPTD’s predictive performance to that of three baseline models, described in section 3: 1) GPIRM, 2) DCGPIRM, and 3) the Bayesian Poisson tensor factorization (BPTF) model of [Schein et al. \(2015\)](#). All three models use a Poisson likelihood and have the same two hyperparameters as BPTD—i.e., ϵ_0 and γ_0 . We set ϵ_0 to 0.1, as recommended by [Gelman \(2006\)](#), and we set γ_0 so that $(\gamma_0 / C)^2 (\gamma_0 / K) (\gamma_0 / R) = 0.01$. This parameterization encourages the elements of the core tensor \mathbf{A} to be sparse. We implemented an MCMC inference algorithm for each model. We provide the full generative process for all three models in the supplementary material.

GPIRM and DCGPIRM are both Tucker decomposition models and thus allocate events to four-dimensional latent classes. The cardinalities of these latent dimensions are the same as BPTD’s—i.e., C , K , and R . In contrast, BPTF is a CP decomposition model and thus allocates events to one-dimensional latent classes. We set the cardinality of this dimension so that the total number of latent factors in BPTF’s likelihood was equal to the total number of latent factors in BPTD’s likelihood—i.e., $Q = \lceil \frac{(V \times C) + (A \times K) + (T \times R) + (C^2 \times K \times R)}{V + V + A + T + 1} \rceil$. We chose not to let BPTF and BPTD use the same number of latent classes—i.e., to set $Q = C^2 \times K \times R$. BPTF does not permit non-compositional allocation, so MCMC inference becomes very slow for even moderate values of C , K , and R . CP decomposition models also tend to overfit when Q is large ([Zhao et al., 2015](#)). Throughout our predictive experiments, we let $C = 20$, $K = 6$, and $R = 3$. These values were well-supported by the data, as we explain in section 7.

Experimental setup: We constructed twelve different observed tensors—six from ICEWS and six from GDELT. Five of the six tensors for each source (ICEWS or GDELT) correspond to one-year time spans with monthly time steps, starting with 2004 and ending with 2008; the sixth corresponds to a five-year time span with monthly time steps, spanning 1995–2000. We divided each tensor \mathbf{Y} into a training tensor $\mathbf{Y}_{\text{train}} = \mathbf{Y}^{(1)}, \dots, \mathbf{Y}^{(T-3)}$ and a test tensor $\mathbf{Y}_{\text{test}} = \mathbf{Y}^{(T-2)}, \dots, \mathbf{Y}^{(T)}$. We further divided each test tensor into a held-out portion and an observed portion via a binary mask. We experimented with two different masks: one that treats the elements involving the most active fifteen countries as the held-out portion and the remaining elements as the observed portion, and one that does the opposite. The first mask enabled us to evaluate the models’ reconstructions of the densest (and arguably most interesting) portion of each test tensor, while the second mask enabled us to evaluate their reconstructions of its complement. Across the entire GDELT database, for example, the elements involving the most active fifteen countries—i.e., 6% of all 233 countries—account for 30%

of the event tokens. Moreover, 40% of these elements are non-zero. These non-zero elements are highly dispersed, with a variance-to-mean ratio of 220. In contrast, only 0.7% of the elements involving the other countries are non-zero. These elements have a variance-to-mean ratio of 26.

For each combination of the four models, twelve tensors, and two masks, we ran 5,000 iterations of MCMC inference on the training tensor. We clamped the country–community factors, the action–topic factors, and the core tensor and then inferred the time-step–regime factors for the test tensor using its observed portion by running 1,000 iterations of MCMC inference. We saved every tenth sample after the first 500. We used each sample, along with the country–community factors, the action–topic factors, and the core tensor, to compute the Poisson rate for each element in the held-out portion of the test tensor. Finally, we averaged these rates across samples and used each element’s average rate to compute its probability. We combined the held-out elements’ probabilities by taking their geometric mean or, equivalently, by computing their inverse perplexity. We chose this combination strategy to ensure that the models were penalized heavily for making poor predictions on the non-zero elements and were not rewarded excessively for making good predictions on the zero elements. By clamping the country–community factors, the action–topic factors, and the core tensor after training, our experimental setup is analogous to that used to assess collaborative filtering models’ strong generalization ability ([Marlin, 2004](#)).

Results: Figure 3 illustrates the results for each combination of the four models, twelve tensors, and two masks. The top row contains the results from the twelve experiments involving the first mask, where the elements involving the most active fifteen countries were treated as the held-out portion. BPTD outperformed the baselines significantly. BPTF—itsself a state-of-the-art model—performed better than BPTD in only one experiment. In general, the Tucker decomposition allows BPTD to learn richer latent structure that generalizes better to held-out data. The bottom row contains the results from the experiments involving the second mask. The models’ performance was closer in these experiments, probably because of the large proportion of easy-to-predict zero elements. BPTD and BPTF performed indistinguishably in these experiments, and both models outperformed the GPIRM and the DCGPIRM. The single-membership nature of the GPIRM and the DCGPIRM prevents them from expressing high levels of heterogeneity in the countries’ rates of activity. When the held-out elements were highly dispersed, these models sometimes made extremely inaccurate predictions. In contrast, the mixed-membership nature of BPTD and BPTF allows them to better express heterogeneous rates of activity.

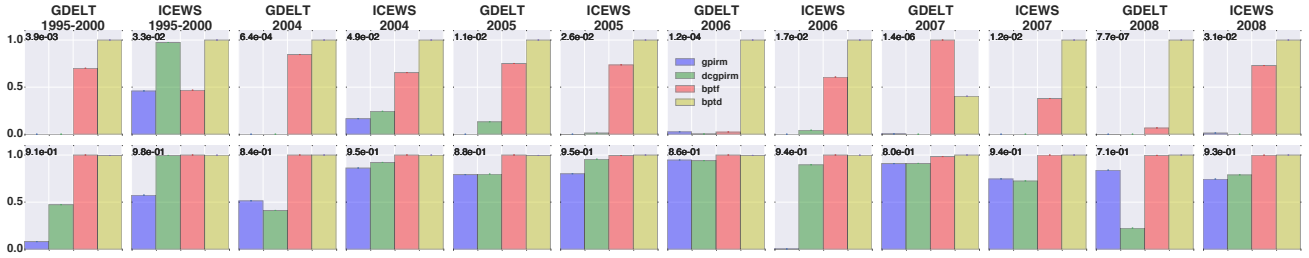


Figure 3. Predictive performance. Each plot shows the inverse perplexity (higher is better) for the four models: the GPIRM (blue), the DCGPIRM (green), BPTF (red), and BPTD (yellow). In the experiments depicted in the top row, we treated the elements involving the most active countries as the held-out portion; in the experiments depicted in the bottom row, we treated the remaining elements as the held-out portion. For ease of comparison, we scaled the inverse perplexities to lie between zero and one; we give the scales in the top-left corners of the plots. BPTD outperformed the baselines significantly when predicting the denser portion of each test tensor (top row).

7. Exploratory Analysis

We used a tensor of ICEWS events spanning 1995–2000, with monthly time steps, to explore the latent structure discovered by BPTD. We initially let $C = 50$, $K = 8$, and $R = 3$ —i.e., $C \times C \times K \times R = 60,000$ latent classes—and used the shrinkage priors to adaptively learn the most appropriate numbers of communities, topics, and regimes. We found $C = 20$ communities and $K = 6$ topics with weights that were significantly greater than zero. We provide a plot of the community weights in the supplementary material. Although all three regimes had non-zero weights, one had a much larger weight than the other two. For comparison, Schein et al. (2015) used fifty latent classes to model the same data, while Hoff (2015) used $C = 4$, $K = 4$, and $R = 4$ to model a similar tensor from GDELТ.

Topics of action types: We show the inferred action–topic factors as a heatmap in the left subplot of figure 4. We ordered the topics by their weights ν_1, \dots, ν_K , which are above the heatmap. The inferred topics correspond very closely to CAMEO’s QuadClass scheme. Moving from left to right, the topics place their mass on increasingly negative actions. Topics 1 and 2 place most of their mass on Verbal Cooperation actions; topic 3 places most of its mass on Material Cooperation actions and the neutral I —*Make Statement* action; topic 4 places most of its mass on Verbal Conflict actions and the I —*Make Statement* action; and topics 5 and 6 place their mass on Material Conflict actions.

Topic-partitioned community–community networks: In the right subplot of figure 4, we visualize the inferred community structure for topic $k = 1$ and the most active regime r . The bottom-left heatmap is the community–community interaction network $\Lambda_k^{(r)}$. The top-left heatmap depicts the rate at which each country i acts as a sender in each community c —i.e., $\theta_{ic} \sum_{j=1}^V \sum_{d=1}^C \theta_{jd} \lambda_{c \rightarrow d}^{(r)}$. Similarly, the bottom-right heatmap depicts the rate at which each country acts as a receiver in each community. The top-right heatmap depicts the number of times each country i took

an action associated with topic k toward each country j during regime r —i.e., $\sum_{c=1}^C \sum_{d=1}^C \sum_{a=1}^A \sum_{t=1}^T y_{ic \rightarrow jd}^{(tr)}$.

We grouped the countries by their strongest community memberships and ordered the communities by their within-community interaction weights $\eta_1^{\circ}, \dots, \eta_C^{\circ}$, from smallest to largest; the thin green lines separate the countries that are strongly associated with one community from the countries that are strongly associated with its adjacent communities.

Some communities contain only one or two strongly associated countries. For example, community 1 contains only the US, community 6 contains only China, and community 7 contains only Russia and Belarus. These communities mostly engage in between-community interaction. Other larger communities, such as communities 9 and 15, mostly engage in within-community interaction. Most communities have a strong geographic interpretation. Moving upward from the bottom, there are communities that correspond to Eastern Europe, East Africa, South-Central Africa, Latin America, Australasia, Central Europe, Central Asia, etc. The community–community interaction network summarizes the patterns in the top-right heatmap. This topic is dominated by the 4 —*Consult* action, so the network is symmetric; the more negative topics have asymmetric community–community interaction networks. We therefore hypothesize that cooperation is an inherently reciprocal type of interaction. We provide visualizations for the other five topics in the supplementary material.

8. Summary

We presented Bayesian Poisson Tucker decomposition (BPTD) for learning the latent structure of international relations from country–country interaction events of the form “country i took action a toward country j at time t .” Unlike previous models, BPTD takes advantage of all three representations of an interaction event data set: 1) a set of event tokens, 2) a tensor of event type counts, and 3) a series of weighted multinet network snapshots. BPTD uses a Poisson

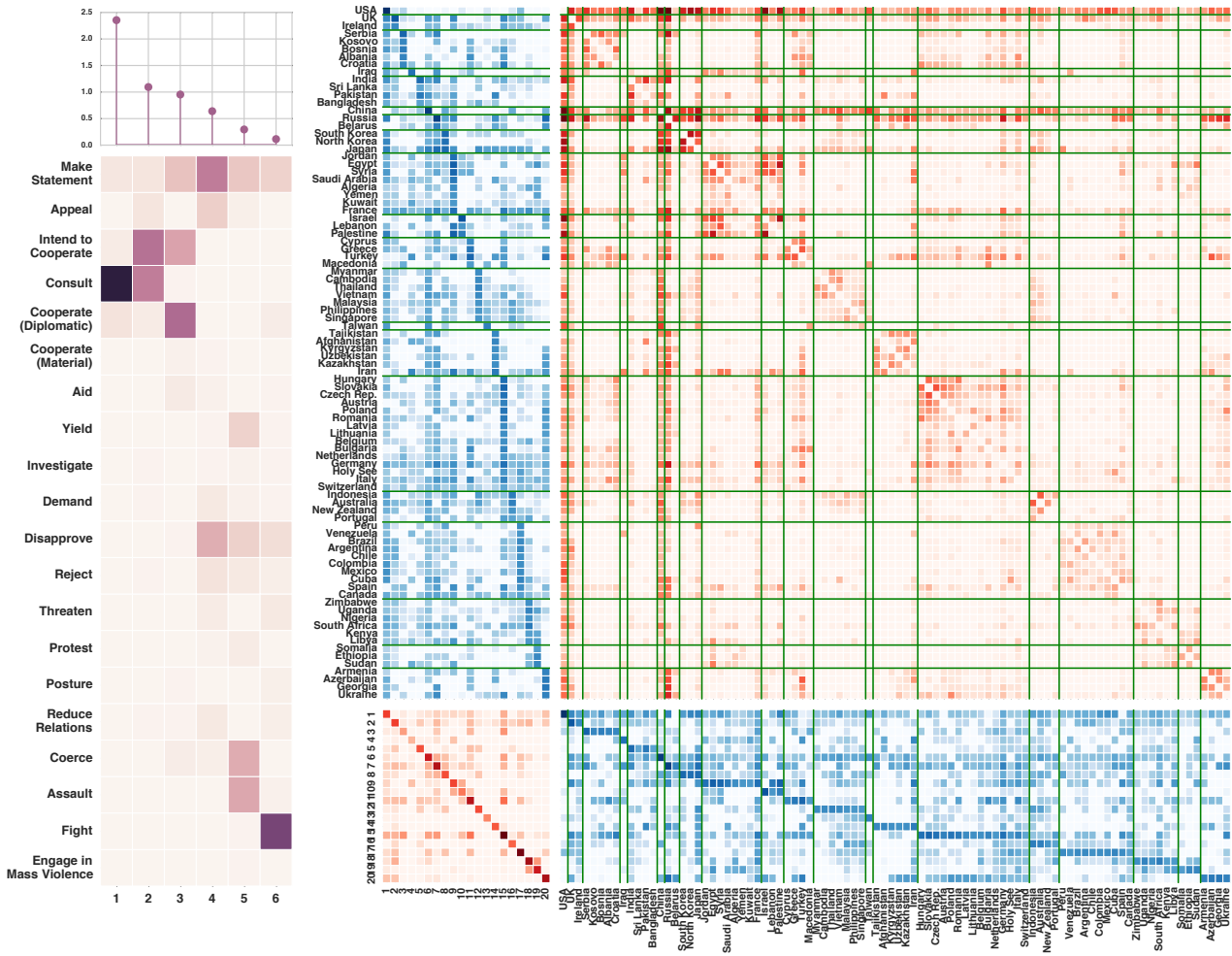


Figure 4. *Left*: Action–topic factors. The topics are ordered by ν_1, \dots, ν_K (above the heatmap). *Right*: Latent structure discovered by BPTD for topic $k = 1$ and the most active regime, including the community–community interaction network (bottom left), the rate at which each country acts as a sender (top left) and a receiver (bottom right) in each community, and the number of times each country i took an action associated with topic k toward each country j during regime r (top right). We show only the most active 100 countries.

likelihood, respecting the discrete nature of the data and its inherent sparsity. Moreover, BPTD yields a compositional allocation inference algorithm that is more efficient than non-compositional allocation algorithms. Because BPTD is a Tucker decomposition model, it shares parameters across latent classes. In contrast, CP decomposition models force each latent class to capture potentially redundant information. BPTD therefore “does more with less.” This efficiency is reflected in our predictive analysis: BPTD outperforms BPTF—a CP decomposition model—as well as two other baselines. BPTD learns interpretable latent structure that aligns with well-known concepts from the networks literature. Specifically, BPTD learns latent country–community memberships, including the number of communities, as well as directed community–community interaction networks that are specific to topics of action types and temporal regimes. This structure captures the complex-

ity of country–country interactions, while revealing patterns that agree with our knowledge of international relations. Finally, although we presented BPTD in the context of interaction events, BPTD is well suited to learning latent structure from other types of multidimensional count data.

Acknowledgements

We thank Abigail Jacobs and Brandon Stewart for helpful discussions. This work was supported by NSF #SBE-0965436, #IIS-1247664, #IIS-1320219; ONR #N00014-11-1-0651; DARPA #FA8750-14-2-0009, #N66001-15-C-4032; Adobe; the John Templeton Foundation; the Sloan Foundation; the UMass Amherst Center for Intelligent Information Retrieval. Any opinions, findings, conclusions, or recommendations expressed in this material are the authors’ and do not necessarily reflect those of the sponsors.

References

- Airoldi, E. M., Blei, D. M., Feinberg, S. E., and Xing, E. P. Mixed membership stochastic blockmodels. *Journal of Machine Learning Research*, 9:1981–2014, 2008.
- Ball, B., Karrer, B., and Newman, M. E. J. Efficient and principled method for detecting communities in networks. *Physical Review E*, 84(3), 2011.
- Blei, D., Ng, A., and Jordan, M. Latent Dirichlet allocation. *Journal of Machine Learning Research*, 3:993–1022, 2003.
- Boschee, E., Lautenschlager, J., O’Brien, S., Shellman, S., Starz, J., and Ward, M. ICEWS coded event data. Harvard Dataverse. V10.
- Cemgil, A. T. Bayesian inference for nonnegative matrix factorisation models. *Computational Intelligence and Neuroscience*, 2009.
- Chi, E. C. and Kolda, T. G. On tensors, sparsity, and nonnegative factorizations. *SIAM Journal on Matrix Analysis and Applications*, 33(4):1272–1299, 2012.
- Cichocki, A., Zdunek, R., Phan, A. H., and Amari, S. *Nonnegative Matrix and Tensor Factorizations: Applications to Exploratory Multi-Way Data Analysis and Blind Source Separation*. John Wiley & Sons, 2009.
- DuBois, C. and Smyth, P. Modeling relational events via latent classes. In *Proceedings of the Sixteenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 803–812, 2010.
- Ferguson, T. S. A Bayesian analysis of some nonparametric problems. *The Annals of Statistics*, 1(2):209–230, 1973.
- Gelman, A. Prior distributions for variance parameters in hierarchical models. *Bayesian Analysis*, 1(3):515–533, 2006.
- Gerner, D. J., Schrod, P. A., Abu-Jabr, R., and Yilmaz, Ö. Conflict and mediation event observations (CAMEO): A new event data framework for the analysis of foreign policy interactions. Working paper.
- Gopalan, P., Ruiz, F. J. R., Ranganath, R., and Blei, D. M. Bayesian nonparametric Poisson factorization for recommendation systems. In *Proceedings of the Seventeenth International Conference on Artificial Intelligence and Statistics*, volume 33, pp. 275–283, 2014.
- Gopalan, P., Hofman, J., and Blei, D. Scalable recommendation with Poisson factorization. In *Proceedings of the Thirty-First Conference on Uncertainty in Artificial Intelligence*, 2015.
- Harshman, R. Foundations of the PARAFAC procedure: Models and conditions for an “explanatory” multimodal factor analysis. *UCLA Working Papers in Phonetics*, 16: 1–84, 1970.
- Hoff, P. Multilinear tensor regression for longitudinal relational data. arXiv:1412.0048, 2014.
- Hoff, P. Equivariant and scale-free Tucker decomposition models. *Bayesian Analysis*, 2015.
- Karrer, B. and Newman, M. E. J. Stochastic blockmodels and community structure in networks. *Physical Review E*, 83(1), 2011.
- Kemp, C., Tenenbaum, J. B., Griffiths, T. L., Yamada, T., and Ueda, N. Learning systems of concepts with an infinite relational model. In *Proceedings of the Twenty-First National Conference on Artificial Intelligence*, 2006.
- Kim, Y.-D. and Choi, S. Nonnegative Tucker decomposition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 20007.
- Kingman, J. F. C. *Poisson Processes*. Oxford University Press, 1972.
- Kolda, T. G. and Bader, B. W. Tensor decompositions and applications. *SIAM Review*, 51(3):455–500, 2009.
- Leetaru, K. and Schrod, P. GDELT: Global data on events, location, and tone, 1979–2012. Working paper, 2013.
- Marlin, B. Collaborative filtering: A machine learning perspective. Master’s thesis, University of Toronto, 2004.
- Mørup, M., Hansen, L. K., and Arnfred, S. M. Algorithms for sparse nonnegative Tucker decompositions. *Neural Computation*, 20(8):2112–2131, 2008.
- Nickel, M., Tresp, V., and Kriegel, H.-P. Factorizing YAGO: Scalable machine learning for linked data. In *Proceedings of the Twenty-First International World Wide Web Conference*, pp. 271–280, 2012.
- Nickel, M., Murphy, K., Tresp, V., and Gabrilovich, E. A review of relational machine learning for knowledge graphs: From multi-relational link prediction to automated knowledge graph construction. arXiv:1503.00759, 2015.
- Nowicki, K. and Snijders, T. A. B. Estimation and prediction for stochastic blockstructures. *Journal of the American Statistical Association*, 96(455):1077–1087, 2001.
- Schein, A., Paisley, J., Blei, D. M., and Wallach, H. Bayesian Poisson tensor factorization for inferring multilateral relations from sparse dyadic event counts. In

Proceedings of the Twenty-First ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp. 1045–1054, 2015.

Schmidt, M. N. and Mørup, M. Nonparametric Bayesian modeling of complex networks: An introduction. *IEEE Signal Processing Magazine*, 30(3):110–128, 2013.

Tucker, L. R. The extension of factor analysis to three-dimensional matrices. In Frederiksen, N. and Gulliksen, H. (eds.), *Contributions to Mathematical Psychology*. Holt, Rinehart and Winston, 1964.

Welling, M. and Weber, M. Positive tensor factorization. *Pattern Recognition Letters*, 22(12):1255–1261, 2001.

Xu, Z., Yan, F., and Qi, Y. Infinite Tucker decomposition: Nonparametric Bayesian models for multiway data analysis. In *Proceedings of the Twenty-Ninth International Conference on Machine Learning*, pp. 1023–1030, 2012.

Zhao, Q., Zhang, L., and Cichocki, A. Bayesian CP factorization of incomplete tensors with automatic rank determination. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(9):1751–1763, 2015.

Zhou, M. Infinite edge partition models for overlapping community detection and link prediction. In *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics*, pp. 1135–1143, 2015.

Zhou, M. and Carin, L. Augment-and-conquer negative binomial processes. In *Advances in Neural Information Processing Systems Twenty-Five*, pp. 2546–2554, 2012.

Zhou, M. and Carin, L. Negative binomial process count and mixture modeling. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(2):307–320, 2015.

Supplementary Material for
“Bayesian Poisson Tucker Decomposition for Learning
the Structure of International Relations”

Proceedings of the 33rd International Conference on Machine Learning, New York, NY, USA, 2016.
JMLR: W&CP volume 48. Copyright 2016 by the author(s).

Aaron Schein

Mingyuan Zhou

David M. Blei

Hanna Wallach

1 Proposition 1

In the limit as $C, K, R \rightarrow \infty$, the expected sum of the core tensor elements is finite and equal to

$$\mathbb{E} \left[\sum_{c=1}^{\infty} \sum_{k=1}^{\infty} \sum_{r=1}^{\infty} \left(\lambda_{c \circlearrowleft^k}^{(r)} + \sum_{d \neq c} \lambda_{c \rightarrow d}^{(r)} \right) \right] = \frac{1}{\delta} \left(\frac{\gamma_0^3}{\zeta^3} + \frac{\gamma_0^4}{\zeta^4} \right).$$

The proof is very similar to that of Zhou (2015, Lemma 1). By the law of total expectation,

$$\begin{aligned} \mathbb{E} \left[\sum_{c=1}^{\infty} \sum_{k=1}^{\infty} \sum_{r=1}^{\infty} \left(\lambda_{c \circlearrowleft^k}^{(r)} + \sum_{d \neq c} \lambda_{c \rightarrow d}^{(r)} \right) \right] &= \sum_{c=1}^{\infty} \sum_{k=1}^{\infty} \sum_{r=1}^{\infty} \left(\mathbb{E} \left[\lambda_{c \circlearrowleft^k}^{(r)} \right] + \sum_{d \neq c} \mathbb{E} \left[\lambda_{c \rightarrow d}^{(r)} \right] \right) \\ &= \sum_{c=1}^{\infty} \sum_{k=1}^{\infty} \sum_{r=1}^{\infty} \left(\mathbb{E} \left[\frac{\eta_c^{\circlearrowleft} \eta_c^{\leftrightarrow} \nu_k \rho_r}{\delta} \right] + \sum_{d \neq c} \mathbb{E} \left[\frac{\eta_c^{\leftrightarrow} \eta_d^{\leftrightarrow} \nu_k \rho_r}{\delta} \right] \right) \\ &= \frac{1}{\delta} \sum_{c=1}^{\infty} \sum_{k=1}^{\infty} \sum_{r=1}^{\infty} \left(\mathbb{E} \left[\eta_c^{\circlearrowleft} \eta_c^{\leftrightarrow} \nu_k \rho_r \right] + \sum_{d \neq c} \mathbb{E} \left[\eta_c^{\leftrightarrow} \eta_d^{\leftrightarrow} \nu_k \rho_r \right] \right) \\ &= \frac{1}{\delta} \mathbb{E} \left[\sum_{k=1}^{\infty} \nu_k \right] \mathbb{E} \left[\sum_{r=1}^{\infty} \rho_r \right] \sum_{c=1}^{\infty} \left(\mathbb{E} \left[\eta_c^{\circlearrowleft} \eta_c^{\leftrightarrow} \right] + \sum_{d \neq c} \mathbb{E} \left[\eta_c^{\leftrightarrow} \eta_d^{\leftrightarrow} \right] \right) \\ &= \frac{1}{\delta} \left(\frac{\gamma_0}{\zeta} \right) \left(\frac{\gamma_0}{\zeta} \right) \sum_{c=1}^{\infty} \left(\mathbb{E} \left[\eta_c^{\circlearrowleft} \eta_c^{\leftrightarrow} \right] + \sum_{d \neq c} \mathbb{E} \left[\eta_c^{\leftrightarrow} \eta_d^{\leftrightarrow} \right] \right) \\ &= \frac{1}{\delta} \left(\frac{\gamma_0}{\zeta} \right)^2 \left(\sum_{c=1}^{\infty} \mathbb{E} \left[\eta_c^{\circlearrowleft} \right] \mathbb{E} \left[\eta_c^{\leftrightarrow} \right] + \mathbb{E} \left[\sum_{c=1}^{\infty} \sum_{d \neq c} \eta_c^{\leftrightarrow} \eta_d^{\leftrightarrow} \right] \right). \end{aligned}$$

The marks $\eta_c^{\circlearrowleft}$ are gamma distributed with mean 1, so

$$\begin{aligned} &= \frac{1}{\delta} \left(\frac{\gamma_0}{\zeta} \right)^2 \left(\mathbb{E} \left[\sum_{c=1}^{\infty} \eta_c^{\leftrightarrow} \right] + \mathbb{E} \left[\sum_{c=1}^{\infty} \sum_{d \neq c} \eta_c^{\leftrightarrow} \eta_d^{\leftrightarrow} \right] \right) \\ &= \frac{1}{\delta} \left(\frac{\gamma_0}{\zeta} \right)^2 \left(\frac{\gamma_0}{\zeta} + \mathbb{E} \left[\sum_{c=1}^{\infty} \sum_{d \neq c} \eta_c^{\leftrightarrow} \eta_d^{\leftrightarrow} \right] \right) \\ &= \frac{1}{\delta} \left(\frac{\gamma_0}{\zeta} \right)^2 \left(\frac{\gamma_0}{\zeta} + \mathbb{E} \left[\sum_{c=1}^{\infty} \sum_{d=1}^{\infty} \eta_c^{\leftrightarrow} \eta_d^{\leftrightarrow} \right] - \mathbb{E} \left[\sum_{c=1}^{\infty} \eta_c^{\leftrightarrow} \eta_c^{\leftrightarrow} \right] \right) \\ &= \frac{1}{\delta} \left(\frac{\gamma_0}{\zeta} \right)^2 \left(\frac{\gamma_0}{\zeta} + \mathbb{E} \left[\left(\sum_{c=1}^{\infty} \eta_c^{\leftrightarrow} \right) \left(\sum_{d=1}^{\infty} \eta_d^{\leftrightarrow} \right) \right] - \mathbb{E} \left[\sum_{c=1}^{\infty} \eta_c^{\leftrightarrow} \eta_c^{\leftrightarrow} \right] \right). \end{aligned}$$

Using $\mathbb{E}[(\sum_{c=1}^{\infty} \eta_c^{\leftrightarrow})(\sum_{d=1}^{\infty} \eta_d^{\leftrightarrow})] = \frac{\gamma_0^2}{\zeta^2} + \frac{\gamma_0}{\zeta^2}$, we can write

$$= \frac{1}{\delta} \left(\frac{\gamma_0}{\zeta} \right)^2 \left(\frac{\gamma_0}{\zeta} + \frac{\gamma_0^2}{\zeta^2} + \frac{\gamma_0}{\zeta^2} - \mathbb{E} \left[\sum_{c=1}^{\infty} \eta_c^{\leftrightarrow} \eta_c^{\leftrightarrow} \right] \right).$$

Finally, using Campbell's Theorem (Kingman, 1972), we know that $\mathbb{E}[\sum_{c=1}^{\infty} \eta_c^{\leftrightarrow} \eta_c^{\leftrightarrow}] = \frac{\gamma_0}{\zeta^2}$, so

$$\begin{aligned} &= \frac{1}{\delta} \left(\frac{\gamma_0}{\zeta} \right)^2 \left(\frac{\gamma_0}{\zeta} + \frac{\gamma_0^2}{\zeta^2} + \frac{\gamma_0}{\zeta^2} - \frac{\gamma_0}{\zeta^2} \right) \\ &= \frac{1}{\delta} \left(\frac{\gamma_0}{\zeta} \right)^2 \left(\frac{\gamma_0}{\zeta} + \frac{\gamma_0^2}{\zeta^2} \right) \\ &= \frac{1}{\delta} \left(\frac{\gamma_0^3}{\zeta^3} + \frac{\gamma_0^4}{\zeta^4} \right). \end{aligned}$$

2 Proposition 2

For an M -dimensional core tensor with $D_1 \times \dots \times D_M$ elements, computing the normalizing constant using non-compositional allocation requires $1 \leq \pi < \infty$ times the number of operations required by compositional allocation. When $D_1 = \dots = D_M = 1$, $\pi = 1$. As $D_m, D_{m'} \rightarrow \infty$ for any m and $m' \neq m$, $\pi \rightarrow \infty$.

Each event token occurs in an M -dimensional discrete coordinate space—i.e., $e_n = \mathbf{p}$, where $\mathbf{p} = (p_1, \dots, p_M)$ is a multi-index. Similarly, each event token's latent class assignment also occurs in an M -dimensional discrete coordinate space—i.e., $z_n = \mathbf{q}$, where $\mathbf{q} = (q_1, \dots, q_M)$ is a multi-index.

Assuming M factor matrices $\Theta^{(1)}, \dots, \Theta^{(M)}$ and an M -dimensional core tensor Λ ,

$$P(z_n = \mathbf{q} | e_n = \mathbf{p}) \propto \lambda_{\mathbf{q}} \prod_{m=1}^M \theta_{p_m q_m}^{(m)}.$$

The computational bottleneck in MCMC inference is computing the normalizing constant

$$Z_{\mathbf{p}} = \sum_{\mathbf{q}} \lambda_{\mathbf{q}} \prod_{m=1}^M \theta_{p_m q_m}^{(m)}.$$

If we use a naïve non-compositional approach, then (assuming each latent dimension m has cardinality D_m) the sum over \mathbf{q} involves $\prod_{m=1}^M D_m$ terms and each term requires M multiplications. Thus, computing $Z_{\mathbf{p}}$ requires a total of $M \prod_{m=1}^M D_m$ multiplications and $\prod_{m=1}^M D_m$ additions.¹

However, we can also compute $Z_{\mathbf{p}}$ using a compositional approach—i.e.,

$$Z_{\mathbf{p}} = \sum_{q_1=1}^{D_1} \theta_{p_1 q_1}^{(1)} \sum_{q_2=1}^{D_2} \theta_{p_2 q_2}^{(2)} \dots \sum_{q_M=1}^{D_M} \theta_{p_M q_M}^{(M)} \lambda_{\mathbf{q}}.$$

¹Computing a sum of N terms requires either N or $N - 1$ additions, depending on whether or not you add the first term to zero. We assume the former definition and say that computing a sum of N terms requires N additions.

This approach requires a total of $\sum_{m=1}^M D_m$ multiplications and $1 + \sum_{m=1}^M (D_m - 1)$ additions.

The ratio π of the number of operations (i.e., multiplications and additions) required by the non-compositional approach to the number of operations required by the compositional approach is

$$\begin{aligned}\pi &= \frac{\left(M \prod_{m=1}^M D_m\right) + \left(\prod_{m=1}^M D_m\right)}{\left(\sum_{m=1}^M D_m\right) + \left(1 + \sum_{m=1}^M (D_m - 1)\right)} \\ &= \frac{(M+1) \prod_{m=1}^M D_m}{\left(2 \sum_{m=1}^M D_m\right) - M + 1}.\end{aligned}$$

As the cardinalities D_1, \dots, D_M of the latent dimensions grow, the numerator grows at a faster rate than the denominator. Therefore π achieves its lower bound when $D_1 = \dots = D_M = 1$:

$$\Omega(\pi) = \frac{(M+1)}{(2M) - M + 1}.$$

Because the numerator grows at a faster rate than the denominator, we can find the upper bound by taking the limit as one or more cardinalities tend to infinity. We work with the inverse ratio

$$\begin{aligned}\pi^{-1} &= \frac{\left(2 \sum_{m=1}^M D_m\right) - M + 1}{(M+1) \prod_{m=1}^M D_m} \\ &= \frac{2}{M+1} \left(\sum_{m=1}^M \frac{D_m}{\prod_{m=1}^M D_m}\right) - \frac{M-1}{M+1} \left(\frac{1}{\prod_{m=1}^M D_m}\right) \\ &= \frac{2}{M+1} \left(\sum_{m=1}^M \frac{1}{\prod_{m' \neq m} D_{m'}}\right) - \frac{M-1}{M+1} \left(\frac{1}{\prod_{m=1}^M D_m}\right).\end{aligned}$$

First, we take the limit of π^{-1} as a single cardinality $D_m \rightarrow \infty$:

$$\begin{aligned}\lim_{D_m \rightarrow \infty} \pi^{-1} &= \lim_{D_m \rightarrow \infty} \frac{2}{M+1} \left(\sum_{m=1}^M \frac{1}{\prod_{n \neq m} D_n}\right) - \lim_{D_m \rightarrow \infty} \frac{M-1}{M+1} \left(\frac{1}{\prod_{m=1}^M D_m}\right) \\ &= \lim_{D_m \rightarrow \infty} \frac{2}{M+1} \left(\sum_{m=1}^M \frac{1}{\prod_{n \neq m} D_n}\right) \\ &= \frac{2}{M+1} \left(\frac{1}{\prod_{n \neq m} D_n}\right).\end{aligned}$$

However, as any second cardinality $D_{m'} \rightarrow \infty$,

$$\lim_{D_m, D_{m'} \rightarrow \infty} \pi^{-1} = \lim_{D_{m'} \rightarrow \infty} \frac{2}{M+1} \left(\frac{1}{\prod_{n \neq m} D_n}\right) \rightarrow 0.$$

Therefore, $\pi \rightarrow \infty$ as any two (or more) cardinalities tend to infinity.

3 Inference

Gibbs sampling repeatedly resamples the value of each latent variable from its conditional posterior. In this section, we provide the conditional posterior for each latent variable in BPTD.

We start by defining the Chinese restaurant table (CRT) distribution (Zhou & Carin, 2015): If $l \sim \text{CRT}(m, r)$ is a CRT-distributed random variable, then, we can equivalently say that

$$l \sim \sum_{n=1}^m \text{Bern} \left(\frac{r}{r+n-1} \right).$$

We also define $g(x) \equiv \ln(1+x)$.

Throughout this section, we use, e.g., $(\theta_{ic} | -)$ to denote θ_{ic} conditioned on \mathbf{Y} , ϵ_0 , γ_0 , and the current values of the other latent variables. We assume that \mathbf{Y} is partially observed and include a binary mask \mathbf{B} , where $b_{i \rightarrow j}^{(t)} = 0$ means that $y_{i \rightarrow j}^{(t)} = 0$ is unobserved, not an observed zero.

Action–Topic Factors:

$$\begin{aligned} y_{\cdot \leftrightarrow \cdot}^{(\cdot)} &\equiv \sum_{i=1}^V \sum_{c=1}^C \sum_{j \neq i}^C \sum_{d=1}^C \sum_{t=1}^T \sum_{r=1}^R y_{ic \rightarrow dj}^{(tr)} \\ \xi_{ak} &\equiv \sum_{i=1}^V \sum_{j \neq i}^C \sum_{t=1}^T b_{i \rightarrow j}^{(t)} \sum_{c=1}^C \theta_{ic} \sum_{d=1}^C \theta_{jd} \sum_{r=1}^R \psi_{tr} \lambda_{c \rightarrow d}^{(r)} \\ (\phi_{ak} | -) &\sim \Gamma \left(\epsilon_0 + y_{\cdot \leftrightarrow \cdot}^{(\cdot)}, \epsilon_0 + \xi_{ak} \right) \end{aligned}$$

Time-Step–Regime Factors:

$$\begin{aligned} y_{\cdot \rightarrow \cdot}^{(tr)} &\equiv \sum_{i=1}^V \sum_{c=1}^C \sum_{j \neq i}^C \sum_{d=1}^C \sum_{a=1}^A \sum_{k=1}^K y_{ic \rightarrow dj}^{(tr)} \\ \xi_{tr} &\equiv \sum_{i=1}^V \sum_{j \neq i}^C \sum_{a=1}^A b_{i \rightarrow j}^{(t)} \sum_{c=1}^C \theta_{ic} \sum_{d=1}^C \theta_{jd} \sum_{k=1}^K \phi_{ak} \lambda_{c \rightarrow d}^{(r)} \\ (\psi_{tr} | -) &\sim \Gamma \left(\epsilon_0 + y_{\cdot \rightarrow \cdot}^{(tr)}, \epsilon_0 + \xi_{tr} \right) \end{aligned}$$

Country–Community Factors:

$$\begin{aligned}
y_{ic\leftrightarrow}^{(\cdot)} &\equiv \sum_{j \neq i}^C \sum_{d=1}^A \sum_{a=1}^K \sum_{t=1}^T \sum_{r=1}^R \left(y_{ic \xrightarrow{ak} dj}^{(tr)} + y_{jd \xrightarrow{ak} ci}^{(tr)} \right) \\
\xi_{ic} &\equiv \sum_{j \neq i}^A \sum_{a=1}^T \left(b_{i \xrightarrow{a} j}^{(t)} \sum_{d=1}^C \theta_{jd} \sum_{k=1}^K \phi_{ak} \sum_{r=1}^R \psi_{tr} \lambda_{c \xrightarrow{k} d}^{(r)} + b_{j \xrightarrow{a} i}^{(t)} \sum_{d=1}^C \theta_{jd} \sum_{k=1}^K \phi_{ak} \sum_{r=1}^R \psi_{tr} \lambda_{d \xrightarrow{k} c}^{(r)} \right) \\
(\theta_{ic} | -) &\sim \Gamma \left(\alpha_i + y_{ic\leftrightarrow}^{(\cdot)}, \beta_i + \xi_{ic} \right)
\end{aligned}$$

Auxiliary Latent Country–Community Counts:

$$(\ell_{ic} | -) \sim \text{CRT} \left(y_{ic\leftrightarrow}^{(\cdot)}, \alpha_i \right)$$

Per-Country Shape Parameters:

$$(\alpha_i | -) \sim \Gamma \left(\epsilon_0 + \sum_{c=1}^C \ell_{ic}, \epsilon_0 + \sum_{c=1}^C g \left(\xi_{ic} \beta_i^{-1} \right) \right)$$

Per-Country Rate Parameters:

$$(\beta_i | -) \sim \Gamma \left(\epsilon_0 + C \alpha_i, \epsilon_0 + \sum_{c=1}^C \theta_{ic} \right)$$

Diagonal Elements of the Core Tensor:

$$\begin{aligned}
\omega_{c \circlearrowleft}^{(r)} &\equiv \eta_c^{\circlearrowleft} \eta_c^{\leftrightarrow} \nu_k \rho_r \\
y_{c \circlearrowleft}^{(r)} &\equiv \sum_{i=1}^V \sum_{j \neq i}^A \sum_{a=1}^T \sum_{t=1}^T y_{ic \xrightarrow{ak} cj}^{(tr)} \\
\xi_{c \circlearrowleft}^{(r)} &\equiv \sum_{i=1}^V \theta_{ic} \sum_{j \neq i}^A \theta_{jc} \sum_{a=1}^A \phi_{ak} \sum_{t=1}^T \psi_{tr} b_{i \xrightarrow{a} j}^{(t)} \\
(\lambda_{c \circlearrowleft}^{(r)} | -) &\sim \Gamma \left(\omega_{c \circlearrowleft}^{(r)} + y_{c \circlearrowleft}^{(r)}, \delta + \xi_{c \circlearrowleft}^{(r)} \right)
\end{aligned}$$

Off-Diagonal Elements of the Core Tensor:

$$\begin{aligned}
\omega_{c \rightarrow d}^{(r)} &\equiv \eta_c^{\leftrightarrow} \eta_d^{\leftrightarrow} \nu_k \rho_r & c \neq d \\
y_{c \rightarrow d}^{(r)} &\equiv \sum_{i=1}^V \sum_{j \neq i}^A \sum_{a=1}^A \sum_{t=1}^T y_{ic \rightarrow dj}^{(tr)} & c \neq d \\
\xi_{c \rightarrow d}^{(r)} &\equiv \sum_{i=1}^V \theta_{ic} \sum_{j \neq i}^A \theta_{jd} \sum_{a=1}^A \phi_{ak} \sum_{t=1}^T \psi_{tr} b_{i \rightarrow j}^{(t)} & c \neq d \\
\left(\lambda_{c \rightarrow d}^{(r)} \mid - \right) &\sim \Gamma \left(\omega_{c \rightarrow d}^{(r)} + y_{c \rightarrow d}^{(r)}, \delta + \xi_{c \rightarrow d}^{(r)} \right) & c \neq d
\end{aligned}$$

Core Rate Parameter:

$$\begin{aligned}
\omega_{\cdot \leftrightarrow \cdot}^{(\cdot)} &\equiv \sum_{c=1}^C \sum_{k=1}^K \sum_{r=1}^R \left(\omega_{c \circlearrowleft k}^{(r)} + \sum_{d \neq c} \omega_{c \rightarrow d}^{(r)} \right) \\
\lambda_{\cdot \leftrightarrow \cdot}^{(\cdot)} &\equiv \sum_{c=1}^C \sum_{k=1}^K \sum_{r=1}^R \left(\lambda_{c \circlearrowleft k}^{(r)} + \sum_{d \neq c} \lambda_{c \rightarrow d}^{(r)} \right) \\
(\delta \mid -) &\sim \Gamma \left(\epsilon_0 + \omega_{\cdot \leftrightarrow \cdot}^{(\cdot)}, \epsilon_0 + \lambda_{\cdot \leftrightarrow \cdot}^{(\cdot)} \right)
\end{aligned}$$

Diagonal Auxiliary Latent Core Counts:

$$\ell_{c \circlearrowleft k}^{(r)} \sim \text{CRT} \left(y_{c \circlearrowleft k}^{(r)}, \omega_{c \circlearrowleft k}^{(r)} \right)$$

Off-Diagonal Auxiliary Latent Core Counts:

$$\ell_{c \rightarrow d}^{(r)} \sim \text{CRT} \left(y_{c \rightarrow d}^{(r)}, \omega_{c \rightarrow d}^{(r)} \right) \quad c \neq d$$

Within-Community Weights:

$$\begin{aligned}
\ell_{c \circlearrowleft \cdot}^{(\cdot)} &\equiv \sum_{k=1}^K \sum_{r=1}^R \ell_{c \circlearrowleft k}^{(r)} \\
\xi_c^{\circlearrowleft} &\equiv \sum_{r=1}^R \rho_r \sum_{k=1}^K \nu_k \sum_{d \neq c} \eta_d^{\leftrightarrow} \left(g \left(\xi_{c \rightarrow d}^{(r)} \delta^{-1} \right) + g \left(\xi_{d \rightarrow c}^{(r)} \delta^{-1} \right) \right) \\
(\eta_c^{\circlearrowleft} \mid -) &\sim \Gamma \left(\frac{\gamma_0}{C} + \ell_{c \circlearrowleft \cdot}^{(\cdot)}, \zeta + \xi_c^{\circlearrowleft} \right)
\end{aligned}$$

Between-Community Weights:

$$\begin{aligned}\ell_{c \leftrightarrow \cdot}^{(\cdot)} &\equiv \ell_{c \odot \cdot}^{(\cdot)} + \sum_{d \neq c} \sum_{k=1}^K \sum_{r=1}^R \left(\ell_{c \rightarrow d}^{(r)} + \ell_{d \rightarrow c}^{(r)} \right) \\ \xi_c^{\leftrightarrow} &\equiv \sum_{r=1}^R \rho_r \sum_{k=1}^K \nu_k \left[\eta_c^{\odot} g \left(\xi_{c \odot^k}^{(r)} \delta^{-1} \right) + \sum_{d \neq c} \eta_d^{\leftrightarrow} \left(g \left(\xi_{c \rightarrow d}^{(r)} \delta^{-1} \right) + g \left(\xi_{d \rightarrow c}^{(r)} \delta^{-1} \right) \right) \right] \\ (\eta_c^{\leftrightarrow} | -) &\sim \Gamma \left(\frac{\gamma_0}{C} + \ell_{c \leftrightarrow \cdot}^{(\cdot)}, \zeta + \xi_c^{\leftrightarrow} \right)\end{aligned}$$

Topic Weights:

$$\begin{aligned}\ell_{\cdot \rightarrow k}^{(\cdot)} &\equiv \sum_{c=1}^C \sum_{d=1}^C \sum_{r=1}^R \ell_{c \rightarrow d}^{(r)} \\ \xi_k &\equiv \sum_{r=1}^R \rho_r \sum_{c=1}^C \eta_c^{\leftrightarrow} \left[\eta_c^{\odot} g \left(\xi_{c \odot^k}^{(r)} \delta^{-1} \right) + \sum_{d \neq c} \eta_d^{\leftrightarrow} \left(g \left(\xi_{c \rightarrow d}^{(r)} \delta^{-1} \right) + g \left(\xi_{d \rightarrow c}^{(r)} \delta^{-1} \right) \right) \right] \\ (\nu_k | -) &\sim \Gamma \left(\frac{\gamma_0}{K} + \ell_{\cdot \rightarrow k}^{(\cdot)}, \zeta + \xi_k \right)\end{aligned}$$

Regime Weights:

$$\begin{aligned}\ell_{\cdot \rightarrow \cdot}^{(r)} &\equiv \sum_{c=1}^C \sum_{d=1}^C \sum_{k=1}^K \ell_{c \rightarrow d}^{(r)} \\ \xi_r &\equiv \sum_{k=1}^K \nu_k \sum_{c=1}^C \eta_c^{\leftrightarrow} \left[\eta_c^{\odot} g \left(\xi_{c \odot^k}^{(r)} \delta^{-1} \right) + \sum_{d \neq c} \eta_d^{\leftrightarrow} \left(g \left(\xi_{c \rightarrow d}^{(r)} \delta^{-1} \right) + g \left(\xi_{d \rightarrow c}^{(r)} \delta^{-1} \right) \right) \right] \\ (\rho_r | -) &\sim \Gamma \left(\frac{\gamma_0}{R} + \ell_{\cdot \rightarrow \cdot}^{(r)}, \zeta + \xi_r \right)\end{aligned}$$

Weights Rate Parameter:

$$\begin{aligned}\omega &\equiv \sum_{c=1}^C \eta_c^{\odot} + \sum_{c=1}^C \eta_c^{\leftrightarrow} + \sum_{k=1}^K \nu_k + \sum_{r=1}^R \rho_r \\ (\zeta | -) &\sim \Gamma (\epsilon_0 + 4\gamma_0, \epsilon_0 + \omega)\end{aligned}$$

4 Baseline Models

BPTF (Schein et al., 2015):

$$\begin{aligned}
 y_{i \xrightarrow{a} j}^{(t)} &\sim \text{Po} \left(\sum_{q=1}^Q \theta_{iq}^{\rightarrow} \theta_{jq}^{\leftarrow} \phi_{aq} \psi_{tq} \lambda_q \right) \\
 \theta_{iq}^{\rightarrow} &\sim \Gamma(\epsilon_0, \beta_1) \\
 \theta_{jq}^{\leftarrow} &\sim \Gamma(\epsilon_0, \beta_2) \\
 \phi_{aq} &\sim \Gamma(\epsilon_0, \beta_3) \\
 \psi_{tq} &\sim \Gamma(\epsilon_0, \beta_4) \\
 \lambda_q &\sim \Gamma \left(\frac{\gamma_0}{Q}, \delta \right) \\
 \beta_1, \dots, \beta_4, \delta &\sim \Gamma(\epsilon_0, \epsilon_0)
 \end{aligned}$$

GPIRM (Schmidt & Mørup, 2013):

$$\begin{aligned}
 y_{i \xrightarrow{a} j}^{(t)} &\sim \text{Po} \left(\lambda_{z_i \xrightarrow{z_a} z_j}^{(z_t)} \right) \\
 z_i &\sim \text{Cat} \left(\frac{\eta_1}{\sum_c \eta_c}, \dots, \frac{\eta_C}{\sum_c \eta_c} \right) \\
 z_a &\sim \text{Cat} \left(\frac{\nu_1}{\sum_k \nu_k}, \dots, \frac{\nu_K}{\sum_k \nu_k} \right) \\
 z_t &\sim \text{Cat} \left(\frac{\rho_1}{\sum_r \rho_r}, \dots, \frac{\rho_R}{\sum_r \rho_r} \right) \\
 \eta_c &\sim \Gamma \left(\frac{\gamma_0}{C}, \zeta \right) \\
 \nu_k &\sim \Gamma \left(\frac{\gamma_0}{K}, \zeta \right) \\
 \rho_r &\sim \Gamma \left(\frac{\gamma_0}{R}, \zeta \right) \\
 \lambda_{c \xrightarrow{k} d}^{(r)}, \zeta &\sim \Gamma(\epsilon_0, \epsilon_0)
 \end{aligned}$$

DCGPIRM:

$$\begin{aligned}
 y_{i \xrightarrow{a} j}^{(t)} &\sim \text{Po} \left(\theta_i \theta_j \phi_a \psi_t \lambda_{z_i \xrightarrow{z_a} z_j}^{(z_t)} \right) \\
 \theta_i, \phi_a, \psi_t &\sim \Gamma(\epsilon_0, \epsilon_0)
 \end{aligned}$$

The rest of the generative process is the same as that of the GPIRM.

5 Supplementary Plots

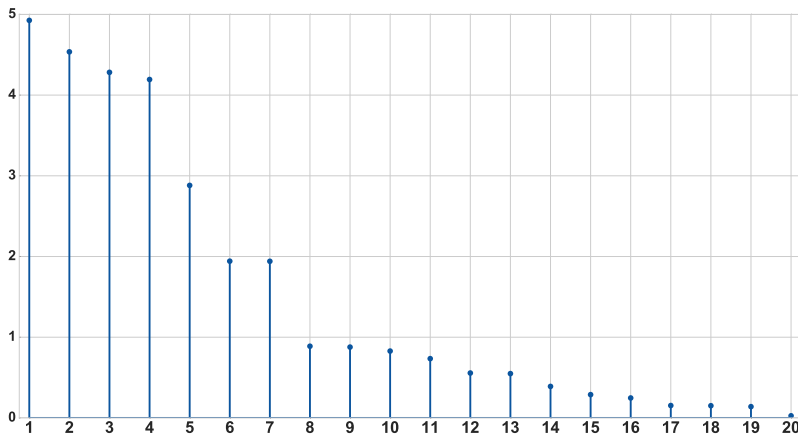


Figure 1: Inferred community weights $\eta_1^{\leftrightarrow}, \dots, \eta_C^{\leftrightarrow}$. We use the between-community weights to interpret shrinkage because they are used for the on- and off-diagonal elements of the core tensor.

References

- Kingman, J. F. C. *Poisson Processes*. Oxford University Press, 1972.
- Schein, A., Paisley, J., Blei, D. M., and Wallach, H. Bayesian Poisson tensor factorization for inferring multilateral relations from sparse dyadic event counts. In *Proceedings of the Twenty-First ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 1045–1054, 2015.
- Schmidt, M. N. and Mørup, M. Nonparametric Bayesian modeling of complex networks: An introduction. *IEEE Signal Processing Magazine*, 30(3):110–128, 2013.
- Zhou, M. Infinite edge partition models for overlapping community detection and link prediction. In *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics*, pp. 1135–1143, 2015.
- Zhou, M. and Carin, L. Negative binomial process count and mixture modeling. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(2):307–320, 2015.

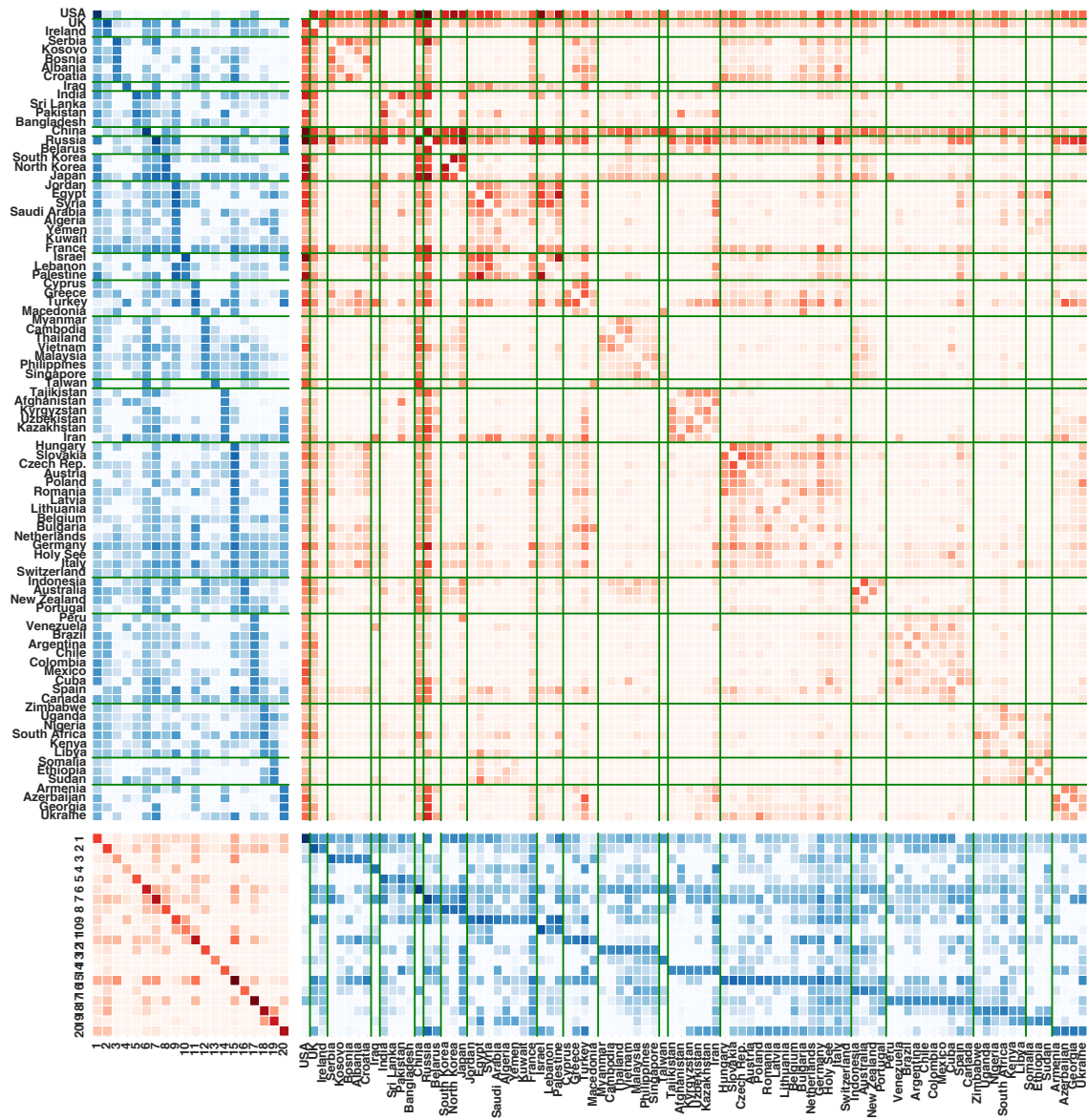


Figure 2: Latent structure discovered by BPTD for topic $k = 1$ (mostly Verbal Cooperation action types) and the most active regime, including the community–community interaction network (bottom left), the rate at which each country acts as a sender (top left) and a receiver (bottom right) in each community, and the number of times each country i took an action associated with topic k toward each country j during regime r (top right). We show only the most active 100 countries.

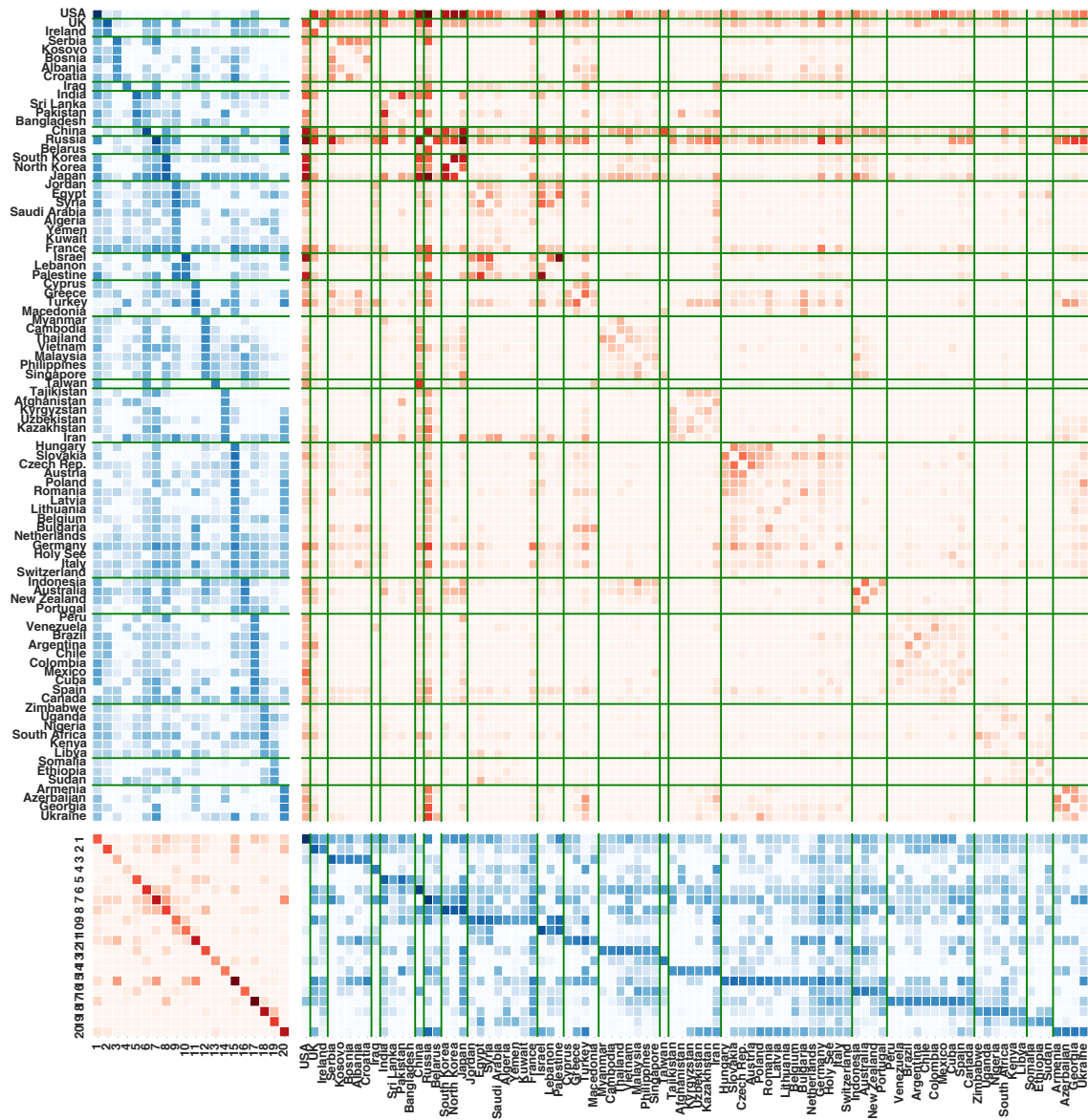


Figure 3: Latent structure discovered by BPTD for topic $k = 2$ (Verbal Cooperation).

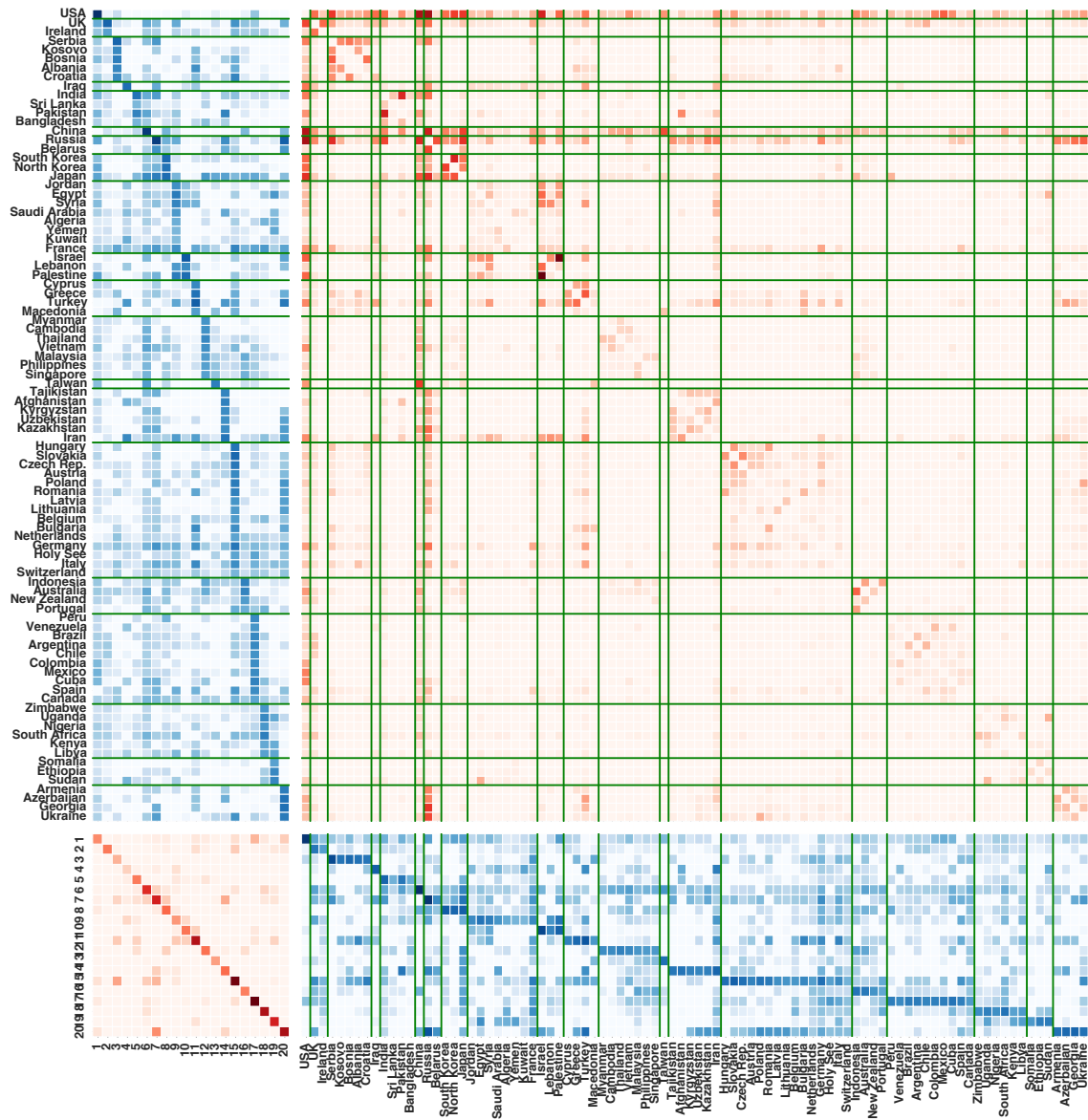


Figure 4: Latent structure discovered by BPTD for topic $k = 3$ (Material Cooperation).

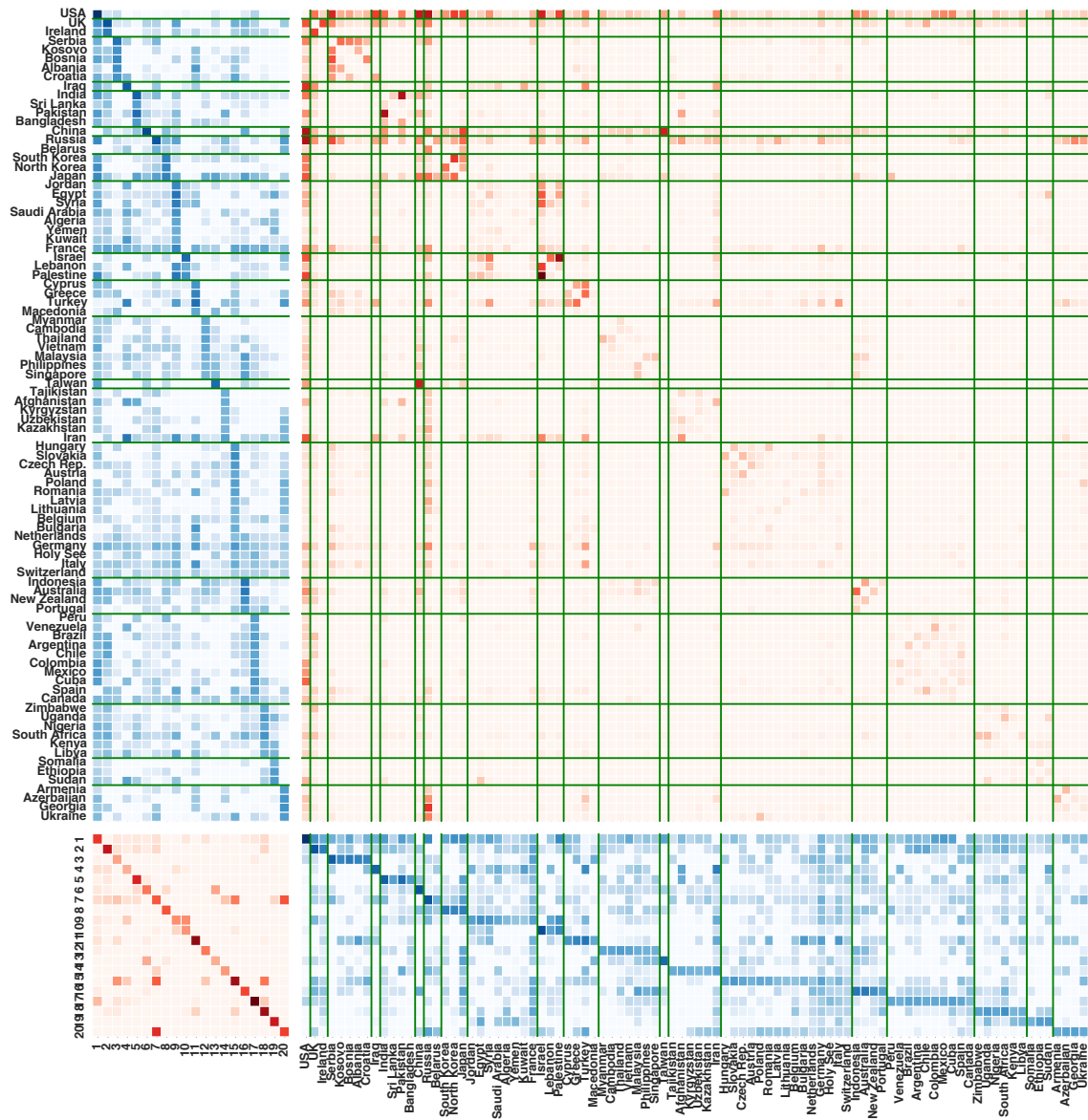


Figure 5: Latent structure discovered by BPTD for topic $k=4$ (Verbal Conflict).

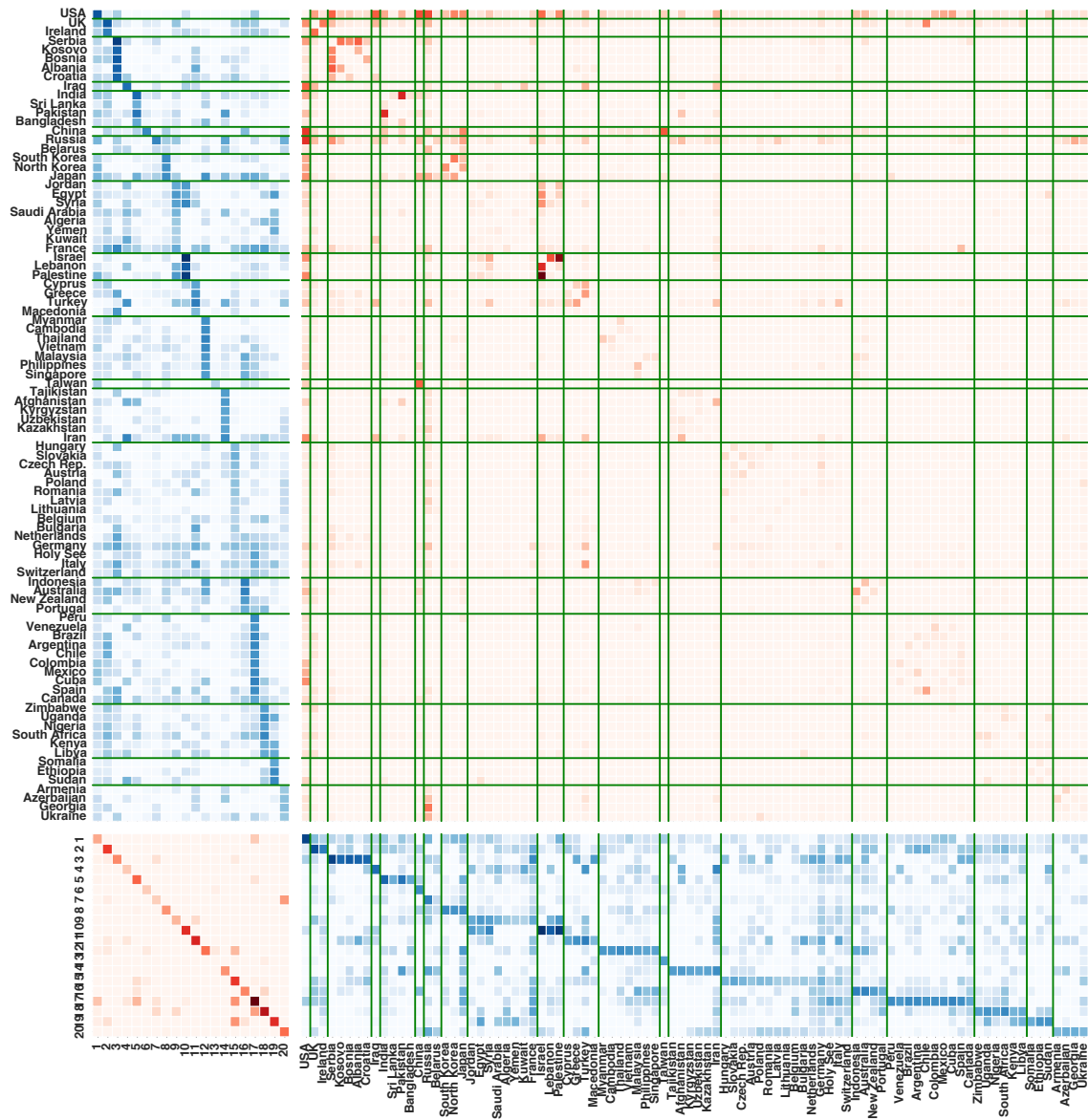


Figure 6: Latent structure discovered by BPTD for topic $k = 5$ (Material Conflict).

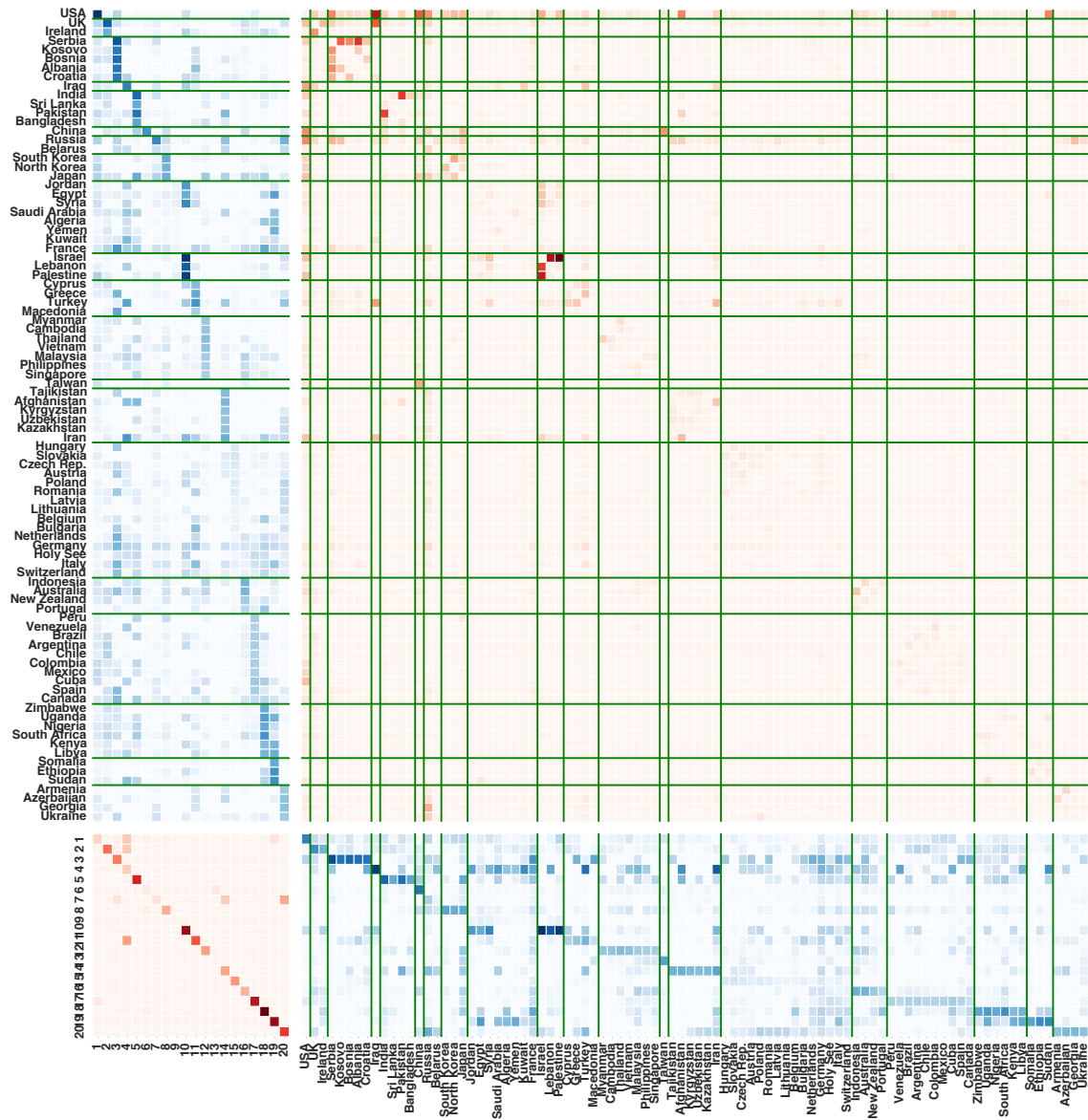


Figure 7: Latent structure discovered by BPTD for topic $k = 6$ (Material Conflict).