

A Network Model for Dynamic Textual Communications with Application to Government Email Corpora*

Bomin Kim¹, Aaron Schein³, Bruce Desmarais¹, and Hanna Wallach^{2,3}

¹Pennsylvania State University

²Microsoft Research NYC

³University of Massachusetts Amherst

November 2, 2017

Abstract

In this paper, we introduce the interaction-partitioned topic model (IPTM)—a probabilistic model for who communicates with whom about what, and when. Broadly speaking, the IPTM partitions time-stamped textual communications, such as emails, according to both the network dynamics that they reflect and their content. To define the IPTM, we integrate a dynamic version of the exponential random graph model—a generative model for ties that tend toward structural features such as triangles—and latent Dirichlet allocation—a generative model for topic-based content. The IPTM assigns each topic to an “interaction pattern”—a generative process for ties that is governed by a set of dynamic network features. Each communication is then modeled as a mixture of topics and their corresponding interaction patterns. We use the IPTM to analyze emails sent between department managers in Dare county government in North Carolina; these email corpora covers the Outer Banks during the time period surrounding Hurricane Sandy. Via this application, we demonstrate that the IPTM is effective at predicting and explaining continuous-time textual communications.

1 Introduction

In recent decades, real-time digitized textual communication has developed into a ubiquitous form of social and professional interaction (see, e.g., Kanungo and Jain, 2008; Szóstek, 2011; Burgess et al., 2004; Pew, 2016). From the perspective of the computational social scientist, this has lead to a growing need for methods of modeling interactions that manifest as text exchanged in continuous time (e.g., e-mail messages). A number of models that build upon topic modeling through Latent Dirichlet Allocation (Blei et al., 2003) to incorporate link data as well as textual content have been developed recently (McCallum et al., 2005; Lim et al., 2013; Krafft et al., 2012). These models are innovative in their extensions that incorporate network tie information. However, none of the models that are currently available in the literature integrate the rich random-graph structure offered by state of the art models for network structure—in particular, the exponential random graph model (ERGM) (Robins et al., 2007; Chatterjee et al., 2013; Hunter et al., 2008). The ERGM is the canonical model for network structure, as it is flexible enough to specify a generative model that accounts for nearly any pattern of tie formation (e.g., tie reciprocation, clustering, popularity effects) (Desmarais and Cranmer, 2017). We build upon recent extensions of ERGM that model time-stamped ties (Perry and Wolfe, 2013; Butts, 2008), and develop the interaction-partitioned topic model (IPTM) to simultaneously model the network structural patterns that govern tie formation, and the content in the communications.

*Prepared for presentation at the New Directions in Analyzing Text as Data (Text As Data 2017). This work was supported in part by the University of Massachusetts Amherst Center for Intelligent Information Retrieval and in part by National Science Foundation grants DGE-1144860, SES-1619644, and CISE-1320219. Any opinions, findings, and conclusions or recommendations are those of the authors and do not necessarily reflect those of the sponsors.

ERGM, and models based on ERGM, provide a framework for explaining or predicting ties between nodes using the network sub-structures in which the two nodes are embedded (e.g., an ERGM specification may predict ties between two nodes that have many shared partners). ERGM-style models have been used for many applications in which the ties between nodes are annotated with text. The text, despite providing rich information regarding the strength, scope, and character of the ties, has been largely excluded from these analyses, due to the inability of ERGM-style models to incorporate textual attributes of ties. These application domains include, among other applicaitons, the study of legislative networks in which networks reflect legislators' co-support of bills, but exclude bill text (Bratton and Rouse, 2011; Alemán and Calvo, 2013); the study of alliance networks in which networks reflect countries' co-signing of treaties, but exclude treaty text (Camber Warren, 2010; Cranmer et al., 2012b,a; Kinne, 2016); the study of scientific co-authorship networks that exclude the text of the co-authored papers (Kronegger et al., 2011; Liang, 2015; Fahmy and Young, 2016); and the study of text-based interaction on social media (e.g., users tied via 'mentions' on twitter) (Yoon and Park, 2014; Peng et al., 2016; Lai et al., 2017).

In defining and testing the IPTM we embed three core conceptual properties, in addition to modeling both text and network structure. First, we link the content component of the model, and network component of the model such that knowing who is communicating with whom at what time (i.e., the network component) provides information about the content of communication, and vice versa. Second, we fully specify the network dynamic component of the model such that, given the content of the communication and the history of tie formation, we can draw an exact, continuous-time prediction of when, by whom, and to whom the communication will be sent. Third, we formulate the network dynamic component of the model such that the model can represent, and be used to test hypotheses regarding, canonical processes relevant to network theory such as preferential attachment—the tendency for actors to prefer interacting with actors who have been popular in the past (Barabási and Albert, 1999; Vázquez, 2003; Jeong et al., 2003), reciprocity (Hammer, 1985; Rao and Bandyopadhyay, 1987), and transitivity—the tendency for the friends of friends to become friends (Louch, 2000; Burda et al., 2004). In what follows we (1) present the generative process for the IPTM, describing how it meets our theoretical criteria, (2) derive the sampling equations for Bayesian inference with the IPTM, and (3) illustrate the IPTM through application to email corpora of internal communications by government officials in Dare County, NC.

2 IPTM: Model Definition and Derivation

To define and derive the IPTM, we begin by describing a probabilistic process by which documents are generated, where documents include a sender, recipients, contents, and timing. We provide a fully parametric definition of each component of the generative process, which enables the model to be used to simulate distributions of who communicates with whom about what, and when. We take a Bayesian approach to inference for the parameters of the IPTM. In the next section, we derive equations for sampling from the posterior distributions of the IPTM parameters conditional on data generated by the generative process that we define in the current section.

The data generated under the IPTM consists of D unique documents. A single email, indexed by $d \in \{1, \dots, D\}$, is represented by the four components $(i^{(d)}, J^{(d)}, t^{(d)}, \mathbf{w}^{(d)})$. The first two are the sender and recipients of the email: an integer $i^{(d)} \in \{1, \dots, A\}$ indicates the identity of the sender out of A actors (or nodes) and a binary vector $J^{(d)} = \{j_r^{(d)}\}_{r=1}^{\|J^{(d)}\|_1}$, which indicates the identity of the receiver (or receivers) out of $A - 1$ actors, where $\|J_i^{(d)}\|_1 \in \{1, \dots, A - 1\}$ denotes the total number of possible receivers. Next, $t^{(d)}$ is the timestamp of the email d . Lastly, $\mathbf{w}^{(d)} = \{w_n^{(d)}\}_{n=1}^{N^{(d)}}$ is a set of tokens, or word type instances, that comprise the text of the email, where $N^{(d)}$ denotes the total number of tokens in a document.

In this section, we illustrate how the words $\mathbf{w}^{(d)}$ are generated according to latent Dirichlet allocation (Blei et al., 2003), and then how the other components, $(i^{(d)}, J^{(d)}, t^{(d)})$, are generated conditional on the document content. For simplicity, we assume that documents are ordered by time such that $t^{(d)} < t^{(d+1)}$ for all $d = 1, \dots, D$.

2.1 Content Generating Process

The content generating process follows from the generative process of Latent Dirichlet Allocation Blei et al. (2003). First we generate the global (corpus-wide) variables. Each topic k is associated with a cluster, or interaction pattern, assignment c_k , where c_k can take one of $c = \{1, 2, \dots, C\}$ values. There are two main sets of global variables—those that describe the content via topics and those that describe how people interact (interaction patterns). These variables are linked by a third set of variables that associate each topic with the pattern that best describes how people interact when talking about that topic. (Refer to Section 2.5 for Algorithm 1—Algorithm 5.)

There are K topics. Each topic k is a discrete distribution over V word types.

1. $\phi^{(k)} \sim \text{Dirichlet}(\beta, \mathbf{u})$ [**Algorithm 1**]

- A topic k is characterized by a discrete distribution over V word types with probability vector $\phi^{(k)}$. We specify a symmetric Dirichlet prior \mathbf{u} with the concentration parameter β for the probability vector $\phi^{(k)}$.

There are C interaction patterns. Each interaction pattern consists of a vector of coefficients $\mathbf{b}^{(c)}$ in \mathbb{R}^P and a vector of P -dimensional dynamic network statistics for directed edge (i, j) at time t $\mathbf{x}_t^{(c)}(i, j)$. The inner product of $\mathbf{b}^{(c)}$ and $\mathbf{x}_t^{(c)}(i, j)$ is used to generate both the recipient vector for a document and the timing of the document.

2. $\mathbf{b}^{(c)} \sim \text{Multivariate Normal}(\mu_{\mathbf{b}}, \Sigma_{\mathbf{b}})$ [**Algorithm 2**]:

- The vector of coefficients depends on the interaction pattern c . This means that there is variation across interaction patterns in the degree to which document timing and recipients depend upon the dynamic network statistics. The prior for $\mathbf{b}^{(c)}$ is a P -variate multivariate Normal with mean vector $\mu_{\mathbf{b}}$ and covariance matrix $\Sigma_{\mathbf{b}}$.

The topics and interaction patterns are tied together via a set of K categorical variables.

3. $c_k \sim \text{Uniform}(1, C)$ [**Algorithm 3**]:

- Each topic is associated with a single interaction pattern, and topics under same interaction pattern share the network properties via $\mathbf{b}^{(c)}$.

We have now defined all of the variables that make up the generative process of the IPTM. Here, given that the number of words $N^{(d)}$ is known, we assume the following generative process for the words in each document d in a corpus D [**Algorithm 4**]:

4-1. Choose document-topic distribution $\boldsymbol{\theta}^{(d)} \sim \text{Dir}(\alpha, \mathbf{m})$

4-2. For $n = 1$ to $N^{(d)}$:

- (a) Choose a topic $z_n^{(d)} \sim \text{Multinomial}(\boldsymbol{\theta}^{(d)})$
- (b) Choose a word $w_n^{(d)} \sim \text{Multinomial}(\phi^{(z_n^{(d)})})$

2.2 Stochastic Intensity

In this section, we illustrate how a set of dynamic network features and topic-interaction assignments jointly identify the stochastic intensity of a document, which plays a key role in the tie generating process in Section 2.4. Assume that each document $d \in \{1, \dots, D\}$ is associated with an $A \times A$ stochastic intensity matrix $\boldsymbol{\lambda}^{(d)}(t)$, where the $(i, j)^{th}$ element $\lambda_{ij}^{(d)}(t)$ can be interpreted as the likelihood of document d being sent from node i to node j at time t .

First, the content of a document is reflected in the stochastic intensity via the distribution of interaction patterns, $\{p_c^{(d)}\}_{c=1}^C$. To calculate the distribution of interaction patterns within a document, we estimate the proportion of words in document d which are assigned to the topics corresponding

to the interaction pattern c from Section 2.1:

$$p_c^{(d)} = \frac{\sum_{k:c_k=c} N^{(k|d)}}{N^{(d)}}, \quad (1)$$

where $N^{(k|d)}$ is the number of times topic k appears in the document d and $N^{(d)}$ is the total number of words, as defined earlier. By definition, $\sum_{c=1}^C p_c^{(d)} = 1$.

Now, we define the $(i, j)^{th}$ element of the stochastic intensity matrix $\lambda^{(d)}(t)$ in the form of the continuous-time ERGM:

$$\lambda_{ij}^{(d)}(t) = \sum_{c=1}^C p_c^{(d)} \cdot \exp\left\{ \mathbf{b}^{(c)T} \mathbf{x}_t^{(c)}(i, j) \right\}, \quad (2)$$

where $p_c^{(d)}$ is as defined in Equation (1), $\mathbf{b}^{(c)}$ is an unknown vector of coefficients in \mathbb{R}^P corresponding to the interaction pattern c , and $\mathbf{x}_t^{(c)}(i, j)$ is a vector of the P -dimensional dynamic network statistics for directed edge (i, j) at time t corresponding to the interaction pattern c . In other words, $\lambda^{(d)}$ can be seen as the weighted average of exponentiated terms over all interaction pattern, where $\{p_c^{(d)}\}_{c=1}^C$ is used as weights.

2.3 Dynamic Network Statistics

In this Section, we introduce detailed specifications of the dynamic network statsitics in our model. We develop a suite of nine different effects to be used as the components of $\mathbf{x}_t^{(c)}(i, j)$, (intercept, outdegree, indegree, send, receive, 2-send, 2-receive, sibling, and cosibling), which are incorporated as in Equation (2). These statistics capture common network properties such as popularity, centrality, reciprocity, and transitivity. Each network statistic is calculated for each interaction pattern $c = 1, \dots, C$, which means that each interaction pattern can be understood in terms of the ways that network dynamics shape tie formation within the interaction pattern. In addition to assigning interaction-pattern specific intercepts, we introduce the degree, dyadic, and triadic network statistics in this paper, and below are the detailed specifications.

We follow Perry and Wolfe (2013) and define each network feature to have potentially different effects within a number of intervals of recency in the formation of the ties that contribute to the network feature. We partition the interval $[-\infty, t)$ into $L = 4$ sub-intervals with equal length in the log-scale, by setting $\Delta_l = (6 \text{ hours}) \times 4^l$ for $l = 1, \dots, L - 1$ such that Δ_l takes the values 24 hours (=1 day), 96 hours (=4 days), 384 hours (=16 days):

$$\begin{aligned} [-\infty, t) &= [-\infty, t - \Delta_3) \cup [t - \Delta_3, t - \Delta_2) \cup [t - \Delta_2, t - \Delta_1) \cup [t - \Delta_1, t - \Delta_0) \\ &= [-\infty, t - 384h) \cup [t - 384h, t - 96h) \cup [t - 96h, t - 24h) \cup [t - 24h, t - 0) \\ &= I_t^{(4)} \cup I_t^{(3)} \cup I_t^{(2)} \cup I_t^{(1)}, \end{aligned}$$

where $\Delta_0 = 0$ and $I_t^{(l)}$ is the half-open interval $[t - \Delta_l, t - \Delta_{l-1})$.

In the application of the IPTM below, we do not include the last interval $I_t^{(4)}$, history before 16 days ago, since the time intervals covered by our datasets are only eight and twelve weeks in length. Although the specification of these dynamic network covariates could be reformulated based on the objectives of each study, in this paper, we define the degree and dyadic efffects for each $l = 1, \dots, L-1$ and $c = 1, \dots, C$ as

1. $\mathbf{outdegree}_{t,l}^{(c)}(i) = \sum_{d:t^{(d)} \in I_t^{(l)}} p_c^{(d)} \cdot I\{i \rightarrow \forall j\}$
2. $\mathbf{indegree}_{t,l}^{(c)}(j) = \sum_{d:t^{(d)} \in I_t^{(l)}} p_c^{(d)} \cdot I\{\forall i \rightarrow j\}$

3. $\text{send}_{t,l}^{(c)}(i,j) = \sum_{d:t^{(d)} \in I_t^{(l)}} p_c^{(d)} \cdot I\{i \rightarrow j\}$
4. $\text{receive}_{t,l}^{(c)}(i,j) = \sum_{d:t^{(d)} \in I_t^{(l)}} p_c^{(d)} \cdot I\{j \rightarrow i\}$

Next, we define four triadic statistics involving pairs of messages, which are analogous to 2-path statistics commonly used in the network science literature. While Perry and Wolfe (2013) adapted full sets of triadic statistics for each combination of time intervals (e.g. $3 \times 3 = 9$), we maintain 3 intervals per each statistic, by defining 3×3 time windows and sum the combination-specific statistics based on the interval where the triads are closed. (Refer to Figure 1.) As a result, our interval-adjusted definitions of triadic effects become

5. $\mathbf{2\text{-}send}_{t,l}^{(c)}(i,j) = \sum_{(l_1=l \text{ or } l_2=l)} \sum_{h \neq i,j} \left(\sum_{d:t^{(d)} \in I_t^{(l_1)}} p_c^{(d)} \cdot I\{i \rightarrow h\} \right) \cdot \left(\sum_{d':t^{(d')} \in I_t^{(l_2)}} p_c^{(d')} \cdot I\{h \rightarrow j\} \right)$
6. $\mathbf{2\text{-}receive}_{t,l}^{(c)}(i,j) = \sum_{(l_1=l \text{ or } l_2=l)} \sum_{h \neq i,j} \left(\sum_{d:t^{(d)} \in I_t^{(l_1)}} p_c^{(d)} \cdot I\{h \rightarrow i\} \right) \cdot \left(\sum_{d':t^{(d')} \in I_t^{(l_2)}} p_c^{(d')} \cdot I\{j \rightarrow h\} \right)$
7. $\mathbf{sibling}_{t,l}^{(c)}(i,j) = \sum_{(l_1=l \text{ or } l_2=l)} \sum_{h \neq i,j} \left(\sum_{d:t^{(d)} \in I_t^{(l_1)}} p_c^{(d)} \cdot I\{h \rightarrow i\} \right) \cdot \left(\sum_{d':t^{(d')} \in I_t^{(l_2)}} p_c^{(d')} \cdot I\{h \rightarrow j\} \right)$
8. $\mathbf{cosibling}_{t,l}^{(c)}(i,j) = \sum_{(l_1=l \text{ or } l_2=l)} \sum_{h \neq i,j} \left(\sum_{d:t^{(d)} \in I_t^{(l_1)}} p_c^{(d)} \cdot I\{i \rightarrow h\} \right) \cdot \left(\sum_{d':t^{(d')} \in I_t^{(l_2)}} p_c^{(d')} \cdot I\{j \rightarrow h\} \right),$

where $l_1 \in \{1, \dots, 3\}$ and $l_2 \in \{1, \dots, 3\}$.

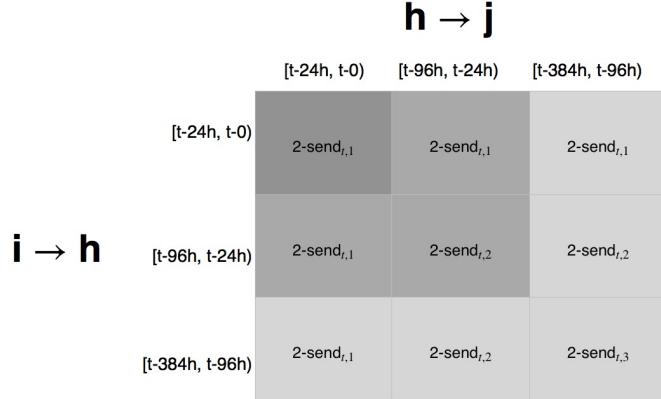


Figure 1: Example of 2-send statistic defined for each interval $l = 1, \dots, 3$. Cells with same shades sum up to one statistic, based on when the triads are “closed”.

2.4 Tie Generating Process

The tie generating process determines the sender, recipients, and timing $(i_o^{(d)}, J_o^{(d)}, t_o^{(d)})$ of the observed document, given the texts. We assume the following tie generating process for each document d in a corpus of D documents:

1. For each sender $i \in \{1, \dots, A\}$, we generate a binary receiver vector of length $A - 1$, $J_i^{(d)}$, from the non-empty Gibbs measure (Fellows and Handcock, 2017) for every $j \in \mathcal{A}_{\setminus i}$.

$$P(J_i^{(d)}) = \frac{1}{Z(\delta, \log(\lambda_i^{(d)}))} \exp \left\{ \log(I(\|J_i^{(d)}\|_1 > 0)) + \sum_{j \in \mathcal{A}_{\setminus i}} (\delta + \log(\lambda_{ij}^{(d)})) J_{ij}^{(d)} \right\}, \quad (3)$$

where δ is a real-valued intercept used to model the number of recipients—i.e., $\|J_i^{(d)}\|_1$, the ℓ_1 -norm (or sum) of the binary recipient vector. The prior distribution for δ is specified as $\text{Normal}(\mu_\delta, \sigma_\delta^2)$. As defined in Section 2.2, $\lambda_{ij}^{(d)}$ is a positive dyad-specific stochastic intensity included in the model, and we use $\lambda_i^{(d)} = \{\lambda_{ij}^{(d)}\}_{j \in \mathcal{A} \setminus i}$ to denote the vector of dyadic weights in which i is the sender. Note that we omitted the notation (t) from Equation (2) and used $\lambda_{ij}^{(d)}$ instead, since the stochastic intensity $\lambda_{ij}^{(d)}$ is always evaluated at time $t_+^{(d-1)}$. The λ_{ij} for d^{th} document is obtained using the history of interactions up to and including the time when the previous document was sent, $t^{(d-1)}$.

The normalizing constant for the non-empty Gibbs measure $Z(\delta, \log(\lambda_i^{(d)}))$, which is the sum of $P(J_i^{(d)})$ over the entire support, can be simplified as:

$$Z(\delta, \log(\lambda_i^{(d)})) = \left(\prod_{j \in \mathcal{A} \setminus i} \left(\exp\{\delta + \log(\lambda_{ij}^{(d)})\} + 1 \right) \right) - 1. \quad (4)$$

Derivation of the normalizing constant is provided in Appendix A.

2. For every sender $i \in \mathcal{A}$, generate the time increments given the latent ties from previous step:

$$\Delta T_{iJ_i} \sim \text{Log-normal}(\mu = \boldsymbol{\eta}^T \mathbf{y}_{iJ_i}^{(d)}, \sigma^2 = \sigma_T^2), \quad (5)$$

where $\boldsymbol{\eta}$ is an unknown vector of coefficients in \mathbb{R}^Q , $\mathbf{y}_{iJ_i}^{(d)}$ a vector of Q-dimensional covariates for given sender and receivers (i, J_i) for document d , and σ_T^2 is a variance parameter. We use $\boldsymbol{\eta} \sim \text{Multivariate Normal}(\mu_\eta, \Sigma_\eta)$ and $\sigma_T^2 \sim \text{Inverse-Gamma}(a_T, b_T)$ as priors. To be specific, $\mathbf{y}_{iJ_i}^{(d)}$ can include an intercept, as well as other features that could strongly affect “time to the next document” such as the day of the week and time of the day when the previous document was sent. Considering the fact that “who sends to whom, and how often they communicated before” is one major factor determining the time to next document, we include $\mathbf{x}_{iJ_i}^{(d)}$ as a default component of $\mathbf{y}_{iJ_i}^{(d)}$, where $\mathbf{x}_{iJ_i}^{(d)}$ is computed by taking the average of network history terms $\mathbf{x}_t^{(c)}(i, j)$ across the chosen receivers $J_i^{(d)}$:

$$\mathbf{x}_{iJ_i}^{(d)}(t) = \sum_{c=1}^C p_c^{(d)} \cdot \left(\frac{1}{\|J_i^{(d)}\|_1} \sum_{j: J_{ij}^{(d)}=1} \mathbf{x}_t^{(c)}(i, j) \right). \quad (6)$$

Note that when there are multiple chosen receivers (i.e. $\|J_i^{(d)}\|_1 > 1$), we call it as “multicast interaction”.

3. Set the observed sender, recipient, and time of the document simultaneously by choosing the sender who generated the minimum time in step 2 and the corresponding recipient and time increment (NOTE: $t_o^{(0)} = 0$):

$$\begin{aligned} i_o^{(d)} &= i_{\min(\Delta T_{iJ_i})}, \\ J_o^{(d)} &= J_{i^{(d)}}, \\ t_o^{(d)} &= t_o^{(d-1)} + \min(\Delta T_{iJ_i}). \end{aligned} \quad (7)$$

The intuition behind this choice is that all possible senders $i \in \mathcal{A}$ are competing against each other to send the next document, and correspondingly introduce the next modification to the history of interactions. Once that next document is sent, the actors in the network revise their plans for document sending, considering the new entry in the history of interactions.

2.5 Joint Generative Process

The algorithms we present in this section form the generative process for D documents. This generative process integrates Sections 2.1 through 2.4.

Algorithm 1 Parameters from priors

[Topic Word Distributions]
for $k=1$ to K do
| draw $\phi^{(k)} \sim \text{Dirichlet}(\beta, \mathbf{u})$
end

[Topic Interaction Pattern Assginments]
for $k=1$ to K do
| draw $c_k \sim \text{Uniform}(1, C)$
end

[Interaction Pattern Parameters]
for $c=1$ to C do
| draw $\mathbf{b}^{(c)} \sim \text{Multivariate Normal}(\mu_{\mathbf{b}}, \Sigma_{\mathbf{b}})$
end

[Time-related Parameters]
draw $\boldsymbol{\eta} \sim \text{Multivariate Normal}(\mu_{\boldsymbol{\eta}}, \Sigma_{\boldsymbol{\eta}})$
draw $\sigma_T^2 \sim \text{Inverse-Gamma}(a_T, b_T)$

[Recipient Size Parameter]
draw $\delta \sim \text{Normal}(\mu_{\delta}, \sigma_{\delta}^2)$

Algorithm 2 Document Generating Process

for $d=1$ to D do
| set $\bar{N}^{(d)} = \max(1, N^{(d)})$
| draw $\boldsymbol{\theta}^{(d)} \sim \text{Dirichlet}(\alpha, \mathbf{m})$
| for $n=1$ to $\bar{N}^{(d)}$ do
| | draw $z_n^{(d)} \sim \text{Multinomial}(\boldsymbol{\theta}^{(d)})$
| | if $N^{(d)} > 0$ then
| | | draw $w_n^{(d)} \sim \text{Multinomial}(\boldsymbol{\phi}^{(z_n^{(d)})})$
| | end
| end
| for $c=1$ to C do
| | set $p_c^{(d)} = \frac{\sum_{k:c_k=c} N^{(k|d)}}{N^{(d)}}$
| end
| for $i=1$ to A do
| | for $j=1$ to A do
| | | if $j \neq i$ then
| | | | calculate $\mathbf{x}_{t_{+}^{(d-1)}}^{(c)}(i, j)$
| | | | set $\lambda_{ij}^{(d)} = \sum_{c=1}^C p_c^{(d)} \cdot \exp\left\{ \mathbf{b}^{(c)T} \mathbf{x}_{t_{+}^{(d-1)}}^{(c)}(i, j) \right\}$
| | | end
| | end
| | draw $J_i^{(d)} \sim \text{Gibbs measure}(\{\lambda_{ij}^{(d)}\}_{j=1}^A, \delta)$
| | calculate $\mathbf{y}_{i, J_i}^{(d)}$ (including $\mathbf{x}_{i, J_i}^{(d)}$)
| | draw $\Delta T_{i, J_i} \sim \text{Lognormal}(\boldsymbol{\eta}^T \mathbf{y}_{i, J_i}^{(d)}, \sigma_T^2)$
| end
| set $i_o^{(d)} = i_{\min(\Delta T_{i, J_i})}$, $J_o^{(d)} = J_{i^{(d)}}$, and $t_o^{(d)} = t_o^{(d-1)} + \min(\Delta T_{i, J_i})$
end

3 Inference

We take a Bayesian approach to inferring the latent variables (i.e., parameters) in the IPTM. The likelihood function is implied by the generative process in Section 2.5. In this section, we derive the joint distribution over the variables introduced earlier $\Phi = \{\phi^{(k)}\}_{k=1}^K, \Theta = \{\theta^{(d)}\}_{d=1}^D, \mathcal{Z} = \{\mathbf{z}^{(d)}\}_{d=1}^D, \mathcal{C} = \{c_k\}_{k=1}^K, \mathcal{B} = \{\mathbf{b}^{(c)}\}_{c=1}^C, \delta, \boldsymbol{\eta}, \sigma_T^2$ as well as the augmented data $\mathcal{J}_a = \{\{J_i^{(d)}\}_{i \neq i_o^{(d)}}\}_{d=1}^D$, given the observed four components $\mathcal{W} = \{\mathbf{w}^{(d)}\}_{d=1}^D, \mathcal{I}_o = \{i_o^{(d)}\}_{d=1}^D, \mathcal{J}_o = \{J_o^{(d)}\}_{d=1}^D$, and $\mathcal{T}_o = \{t_o^{(d)}\}_{d=1}^D$, and the hyperparameters $(\beta, \mathbf{u}, \alpha, \mathbf{m}, \mu_{\mathbf{b}}, \Sigma_{\mathbf{b}}, \mu_{\delta}, \sigma_{\delta}^2, \boldsymbol{\mu}_{\boldsymbol{\eta}}, \Sigma_{\boldsymbol{\eta}}, a_T, b_T)$.

After integrating out Φ and Θ using Dirichlet-multinomial conjugacy (Griffiths and Steyvers, 2004) we sample the remaining unobserved variables from their joint posterior distribution using Markov chain Monte Carlo (MCMC) methods.

Our inference goal is to draw samples from the posterior distribution

$$\begin{aligned} & P(\mathcal{Z}, \mathcal{C}, \mathcal{B}, \delta, \boldsymbol{\eta}, \sigma_T^2, \mathcal{J}_a | \mathcal{W}, \mathcal{I}_o, \mathcal{T}_o, \beta, \mathbf{u}, \alpha, \mathbf{m}, \mu_{\mathbf{b}}, \Sigma_{\mathbf{b}}, \mu_{\delta}, \sigma_{\delta}^2, \boldsymbol{\mu}_{\boldsymbol{\eta}}, \Sigma_{\boldsymbol{\eta}}, a_T, b_T) \\ & \propto P(\mathcal{Z}, \mathcal{C}, \mathcal{B}, \boldsymbol{\eta}, \sigma_T^2, \mathcal{W}, \mathcal{J}_a, \mathcal{I}_o, \mathcal{T}_o | \beta, \mathbf{u}, \alpha, \mathbf{m}, \mu_{\mathbf{b}}, \Sigma_{\mathbf{b}}, \mu_{\delta}, \sigma_{\delta}^2, \boldsymbol{\mu}_{\boldsymbol{\eta}}, \Sigma_{\boldsymbol{\eta}}, a_T, b_T) \\ & = P(\mathcal{Z} | \alpha, \mathbf{m}) P(\mathcal{C} | \mathcal{C}) P(\mathcal{B} | \mathcal{C}, \mu_{\mathbf{b}}, \Sigma_{\mathbf{b}}) P(\delta | \mu_{\delta}, \sigma_{\delta}^2) P(\boldsymbol{\eta} | \boldsymbol{\mu}_{\boldsymbol{\eta}}, \Sigma_{\boldsymbol{\eta}}) P(\sigma_T^2 | a_T, b_T) \\ & \quad \times P(\mathcal{W} | \mathcal{Z}, \beta, \mathbf{u}) P(\mathcal{J}_a, \mathcal{I}_o, \mathcal{T}_o | \mathcal{Z}, \mathcal{C}, \mathcal{B}, \delta, \boldsymbol{\eta}, \sigma_T^2). \end{aligned} \tag{8}$$

To summarize the inference procedure outlined above, we provide pseudocode for Markov Chain Monte Carlo (MCMC) sampling. For better performance and interpretability of the topics we infer, we run n_1 iterations of the hyperparameter optimization technique called “new fixed-point iterations using the Digamma recurrence relation” in Wallach (2008), for every outer iteration o . In addition, while we update the categorical variables \mathcal{Z} , \mathcal{C} , and \mathcal{J}_a once per outer iteration using Gibbs sampling method, we specify a larger number of inner iterations (n_2) for the continuous variables that require Metropolis-Hastings \mathcal{B} and δ , respectively, considering slower mixing. Time-related parameters $\boldsymbol{\eta}$ and σ_T^2 also use Metropolis-Hastings update with n_3 number of inner iterations. When summarizing model results, we only use the samples from the last (i.e., O^{th}) outer loop. The detailed derivation of sampling equations can be found in Appendix B.

Algorithm 3 MCMC

```
set initial values of  $\mathcal{Z}^{(0)}, \mathcal{C}^{(0)}, \mathcal{B}^{(0)}, \delta^{(0)}, \boldsymbol{\eta}^{(0)}, \sigma_T^2$ 
for  $o=1$  to  $O$  do
    for  $n=1$  to  $n_1$  do
        | optimize  $\alpha$  and  $\mathbf{m}$  using hyperparameter optimization in Wallach (2008)
    end
    for  $d=1$  to  $D$  do
        for  $i \in \mathcal{A}_{\setminus i_o^{(d)}}$  do
            | sample the augmented data  $J_i^{(d)}$  following Section B.2
        end
        for  $n=1$  to  $N^{(d)}$  do
            | draw of  $z_n^{(d)} \sim \text{Multinomial}(p^{\mathcal{Z}})$  following Section B.3
        end
    end
    for  $k=1$  to  $K$  do
        | draw  $c_k \sim \text{Multinomial}(p^C)$  following Section B.4
    end
    for  $n=1$  to  $n_2$  do
        | sample  $\mathcal{B}$  and  $\delta$  using Metropolis-Hastings following Section B.5
    end
    for  $n=1$  to  $n_3$  do
        | sample  $\boldsymbol{\eta}$  and  $\sigma_T^2$  using Metropolis-Hastings following Section B.6
    end
end
```

Summarize the results with:

last sample of \mathcal{C} , last sample of \mathcal{Z} , last chain of \mathcal{B} and δ , and last chain of $\boldsymbol{\eta}$ and σ_T^2

4 Getting It Right (GiR) Test

Software development is integral to the objective of applying IPTM to real world data. Code review is a valuable process in any research computing context, and the prevalence of software bugs in statistical software is well documented (e.g., Altman et al., 2004; McCullough, 2009). With highly complex models such as IPTM, there are many ways in which software bugs can be introduced and go unnoticed. As such, we present a joint analysis of the integrity of our generative model, sampling equations, and software implementation.

Geweke (2004) introduced the “Getting it Right” (GiR) test—a joint distribution test of posterior simulators which can detect errors in sampling equations as well as coding errors. The test involves comparing the distributions of variables simulated from two joint distribution samplers, which we call “forward” and “backward” samples. The forward sampler draws unobservable variables from the prior and then generates the observable data conditional on the unobservables. The backward sampler alternates between the inference and an observables simulator, by running the inference code on observable data to obtain posterior estimates of the unobservable variables and then re-generating the observables given the inferred unobservables. The backward sampler is initialized by running an iteration of inference on observables drawn directly from the prior. Since the only information on which both the forward and backward samplers are based is the prior, if the sampling equations are correct and the code is implemented without bugs, each variable should have the same distribution in the forward and backward samples.

In the forward samples, both observable and unobservable variables are generated using Algorithm 2. In the backward samples, unobservable variables are generated using the sampling equations for inference, which we derived in Section 3. In order to generate observable variables in the backward samples, we use the collapsed-time generative process, which we presented in Section ???. For each forward and backward sample that consists of D number of documents, we save these statistics:

1. Mean of network effect parameters $(\mathbf{b}_p^{(1)}, \dots, \mathbf{b}_p^{(C)})$ for every $p = 1, \dots, P$,
2. Network statistic ‘send’ calculated for the last D^{th} document for every $l = 1, \dots, 3$
3. δ value used to generate the samples
4. Mean of the observed recipient size $\|J_o^{(d)}\|_1$ across $d = 1, \dots, D$,
5. Mean of time-increments $t^{(d)} - t^{(d-1)}$ across $d = 1, \dots, D$,
6. Mean topic-interaction pattern assignment c_k across $k = 1, \dots, K$,
7. Number of tokens in topics assigned to each interaction pattern $c = 1, \dots, C$,
8. Number of tokens assigned to each topic $k = 1, \dots, K$,
9. Number of tokens assigned to each unique word type $w = 1, \dots, W$,
10. Time-increment parameter η used to generate the samples, for every $q = 1, \dots, Q$
11. Log-normal variance parameter σ_T^2 used to generate the samples

To keep the computational burden of re-running thousands of rounds of inference manageable, we run GiR using a relatively small artificial sample, consisting of 5 documents, 4 tokens per document, 4 actors, 5 unique word types, 2 interaction patterns, and 4 topics per each forward or backward samples. For detailed settings including the prior specifications, see Appendix C.3. We generated 10^5 sets of forward and backward samples, and then calculated 1,000 quantiles for each of the network effect parameters, and 50 quantiles for the rest of the statistics. We also calculated t-test and Mann-Whitney test p-values in order to test for differences in the distributions generated in the forward and backward samples. Before we calculated these statistics, we thinned our samples by taking every 9th sample starting at the 10,000th sample for a resulting sample size of 10,000, in order to reduce the autocorrelation in the Markov chains. In each case, if we observe a large p-value, this gives us evidence that the distributions generated under forward and backward sampling have the same locations. We depict the GiR results using probability-probability (PP) plots. To compare two samples with a PP-plot we calculate the empirical quantile in each sample of a set of values observed across the two samples, then plot the sets of quantiles in the two samples against each other. If the two samples are from equivalent distributions, the quantiles should line up on a line with zero y -intercept, and unit slope (i.e., a 45-degree line). The GiR test results are depicted in Figure 2, which show that we pass the test on every statistic.

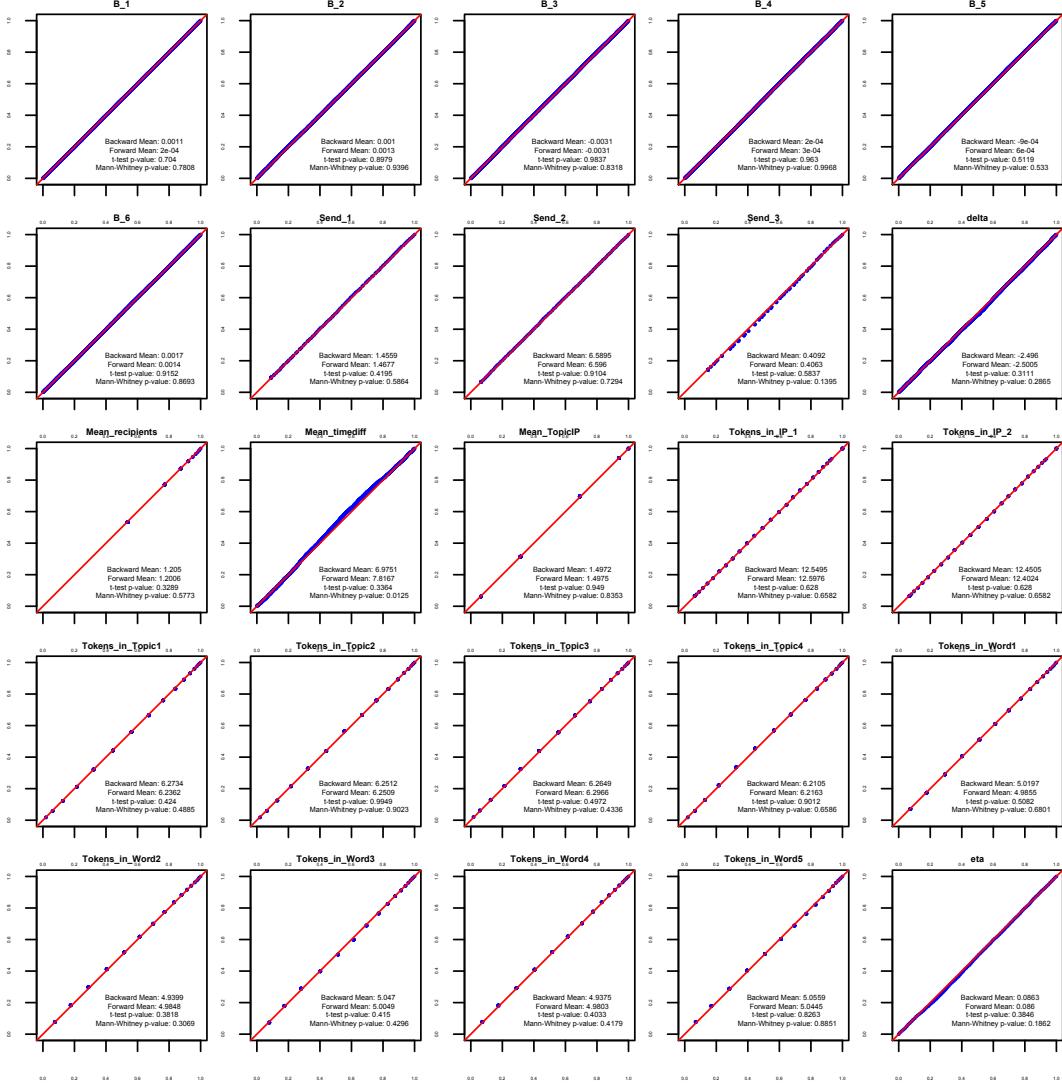


Figure 2: Probability-Probability plot for the 25 GiR test statistics.

5 Application: North Carolina County Government Email Communication During Hurricane Sandy

In our application of the IPTM, we use a subset of the North Carolina county government e-mail dataset collected by ben Aaron et al. (2017). This dataset includes internal e-mail corpora covering the inboxes and outboxes of managerial-level employees of North Carolina county governments. Each county corpus covers a three-month span in 2012. The full dataset covers over twenty counties, but we focus on one county, Dare County, for which the time span included a notable national emergency—Hurricane Sandy (October 26, 2012—October 30, 2012). We chose Dare County, (1) in order to see whether and how communication networks surrounding Hurricane Sandy differed from those surrounding other governmental functions, and (2) to limit the scope of this initial application.



Figure 3: Geographical location of Dare County in North Carolina

5.1 Exploratory Data Analysis

Before discussing the IPTM results, we present a set of exploratory analyses in which we examine characteristics of the data that are relevant to the prevalence of Hurricane Sandy in Dare County government email networks. Dare County data spans October 1st to November 30th containing $D = 1,456$ emails between $A = 27$ actors from 22 departments, and with a vocabulary of size $W = 2907$. We ran the exploratory analysis by looking at three different plots—sending and receiving counts, networks, and word counts—to visualize the networks and content of the email data, with emphasis on the changes during hurricane Sandy.

In Dare county, we saw considerable change in email sending/receiving behaviors during hurricane Sandy (Figure 4). As the hurricane approaches on October 26th, the manager from the emergency services department sent significantly more emails than before, and at the same time there was dramatic rise in the receiving counts for almost every department. Further analysis demonstrated that emergency services department sent a lot of “multicast” emails with a large number of receivers during the Sandy period.

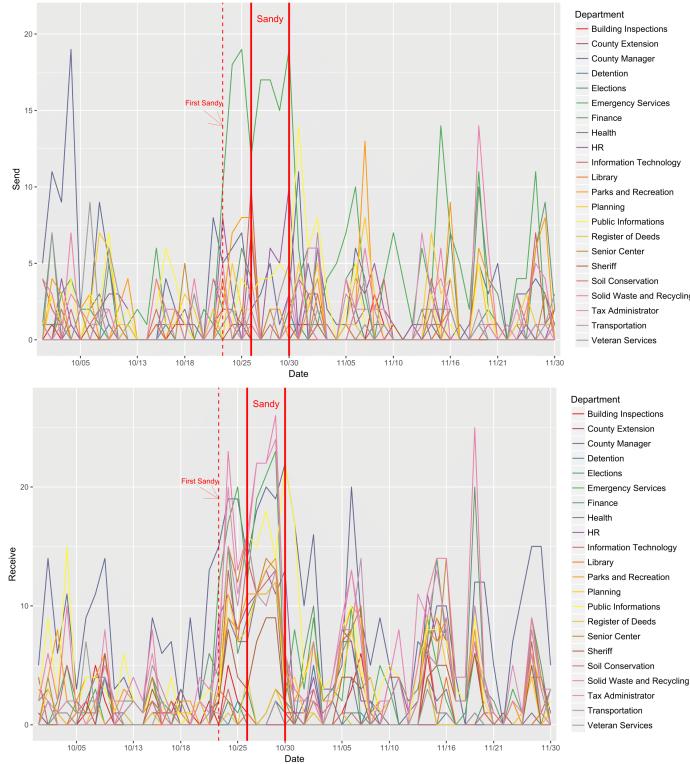


Figure 4: The number of emails sent to (upper) and received by (lower) department in Dare County

The network plots in Figure 5 illustrated same patterns we found in Figure 4. Again, the manager

from the emergency services department became highly central in the network during the Sandy period, and it maintained this pattern after Sandy. Hurricane related conversations continued after Sandy passed the county, since there remained post-hurricane issues from the damage.

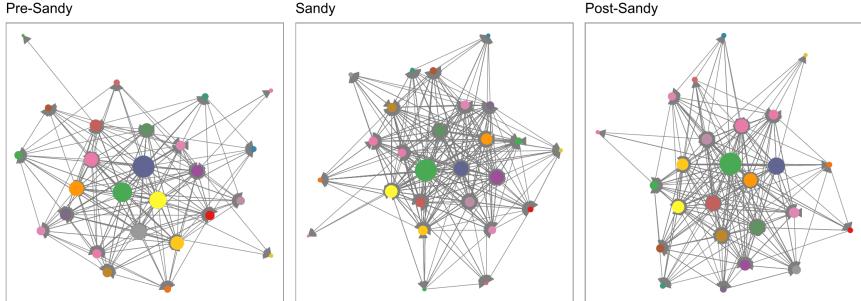


Figure 5: Network plot for three time windows: before Sandy (October 1—October 18), during Sandy (October 19—November 2), and after Sandy (November 3—November 30), in Dare County

Figure 6 reflects the hurricane’s effects on email exchanges as well, and it matches our interpretations from the network aspects. Usage of the two words, ‘hurricane’ and ‘Sandy’, exploded starting a few days before Sandy arrived in Dare County, and multiple emails used the words again in November, implying the continuous discussions on hurricane-related topics.

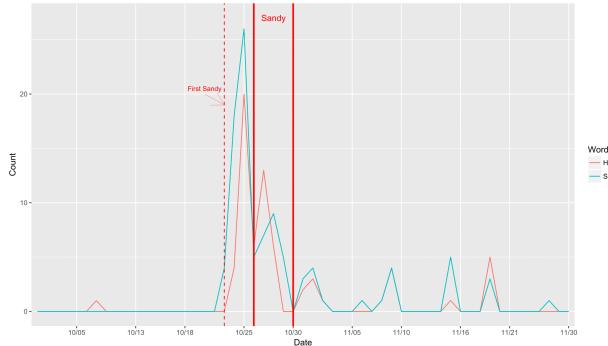


Figure 6: Frequency plot counting how many times the word ‘hurricane’ and ‘sandy’ appeared

Based on this preliminary exploration, we see that Sandy received significant attention in Dare County. However, this basic descriptive finding does not shed light on whether the network structure of communication surrounding Sandy differed from the typical communication patterns.

5.2 IPTM Results

In this section we present the IPTM results. Researchers who use the IPTM can test hypotheses regarding network structure, which is a common—perhaps the most common—use of ERGM-style models for networks. However, with the IPTM those hypotheses can be conditioned on the content area of communication. To provide an example, we articulate expectations regarding the content-conditional structural properties of the county government email networks. Information diffuses more efficiently in networks characterized by a lack of loops (Lin et al., 2010; Iribarren and Moro, 2011) and closed triangles (Roca et al., 2010; Tadić and Thurner, 2004). Assuming that the county governments’ internal communication networks are characterized by efficient communication structure, we expect to see communication regarding the everyday business of the county characterized by negative reciprocity (i.e., 2-loop) effects, and negative triadic effects. The “receive” term captures reciprocity, and triadic effects are captured by 2-send, 2-receive, sibling and cosibling. Reciprocity

and closed triangles are, however, common structural properties of social networks. We expect to see communication regarding personal/social matters to be characterized by positive reciprocity and triadic effects. Lastly, we have limited expectations regarding how communication surrounding Hurricane Sandy will be structured. We take an exploratory approach to the question of whether or not discussion surrounding Sandy forms an efficient communication network structure.

For the IPTM application to the Dare County email data we used $C = 2$ and $K = 20$. We ran the model for $O = 100$ outer iterations with hyperparameter optimization steps for $n_1 = 5$, and the inner iterations for \mathcal{B} and δ were set as $n_2 = 5500$ and $n_3 = 550$, respectively. This time, first 500 and 50 iterations were discarded as a burn-in for inference on \mathcal{B} and δ , and every 10th samples were taken as a thinning process for \mathcal{B} . To present the result, each interaction pattern is summarized with (a) Table 1: the top 15 most likely words to be generated in the topics within the interaction pattern, and (b) Table 2: boxplots visualizing posterior estimates of dynamic network effects $\mathbf{b}^{(c)}$ within each interaction pattern.

First, we see significant differences in the contents related to each interaction pattern. Overall, 55.2% of the words were assigned to the topics in interaction pattern 1, and 44.8% were assigned to the topics in interaction pattern 2. Table 1 demonstrates 5 examples of topics from each interaction pattern, where each topic is summarized by the top 15 words. Interaction pattern 1 seems to represent topics commonly used by government managers. On the other hand, it is interesting that interaction pattern 2 included several topics (e.g. topic 2 and topic 18) with hurricane related words, such as ‘storm’, ‘impacts’ ‘damage’ and ‘ocean’. In addition, one more impressive point is that topic 12 contained words related politics or election, e.g. ‘survey’, ‘parties’, and ‘elections’. In summary, interaction pattern 1 represents usual administrative communications between managers in county government, and interaction pattern 2 to reflects temporary conversations driven by events or emergencies occurring in the county.

IP	1	1	1	1	1
Topic	3 (0.078)	5 (0.065)	19 (0.064)	13 (0.057)	15 (0.053)
Word	water relocation location hills utilities mustian hydrant department skyco kill devil road lane tank map	planning meter room asked needed sure afternoon cheryl johnson issues case letter antennas inspection keep	phone collins drive marshall director human resources manteo phr fax box timesheets -lsb- wanted touch	questions board december call sheets agenda nov hope item weekly management internet told care comp	contact info problem release check weather priority readings rodanthe top collection located health heads ahead
IP	2	2	2	2	2
Topic	14 (0.058)	12 (0.047)	2 (0.045)	6 (0.044)	18 (0.036)
Word	time hours leave monday administrative employees employee work day friday october storm tomorrow hour question	survey voice copy discovery regional parties disclosed elections pin sending prior editor students cost residents	road mirlo storm beach high coastal impacts saturday dot night winds hold bridge expressed normal	library week working place best start visit year albemarle librarian web learning east holiday system	status system area south forecast track pay move assessment opens damage well ocean operation addition

Table 1: Summary of topic-token assignments from Dare County data: top 15 words assigned to each topic, corresponding to interaction pattern assignments

In the Dare County analysis the two interaction patterns exhibited quite different network effects. The effects in interaction pattern 1 were generally greater in magnitude than those in interaction pattern 2, implying that topics related to interaction pattern 2 are less affected by previous email exchanges than those in interaction pattern 1. Most effects are greater in magnitude for the time interval $[t-24h, t)$ and the effect is diminishing as we move to older time intervals $[t-96h, t-24h)$ and $[t-384h, t-96h)$. On the contrary, the statistic ‘send’ had strongest effect in the time interval farthest in the past $[t-384h, t-96h)$. Furthermore, negative reciprocity, which we noted is expected in an efficient communication network, are found in the ‘receive’ statistics effect (except $[t-96h, t-24h)$), in interaction pattern 1. Most of the posterior distributions of network effects in interaction pattern 2 are centered close to zero, providing little evidence of complex network dynamics. The differences between interaction patterns can be explained by the content differences conveyed through Table 1. Since interaction pattern 2 consists of highly time-sensitive topics, ties are less likely to be effected by previous interactions.

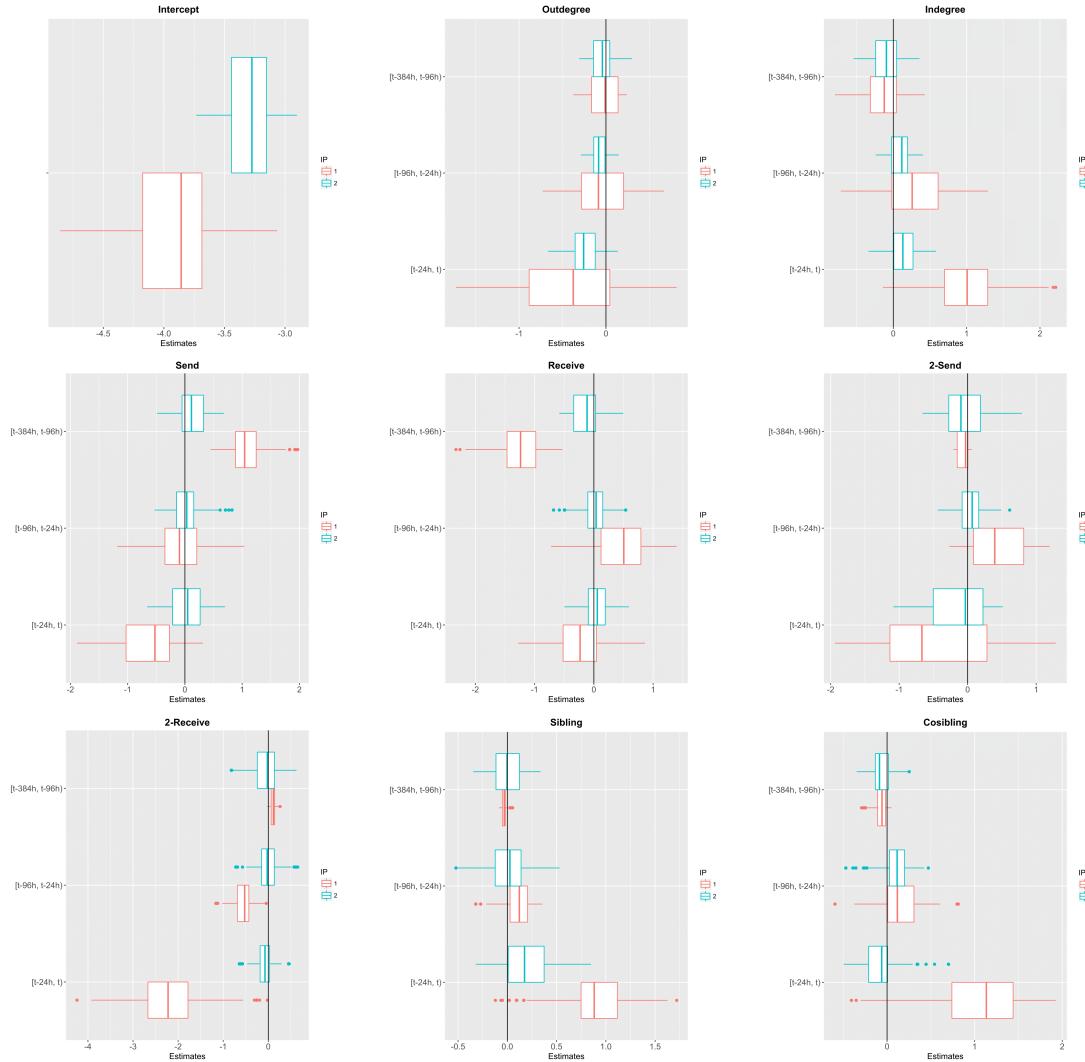


Table 2: 95% credible intervals of posterior estimates of the network effects $b^{(c)}$: $c = 1$ (red) and $c = 2$ (green), using Dare County data

6 Document Prediction Experiments

We use a set of posterior predictive experiments, or document-ahead forecasting, to evaluate the performance of the IPTM as compared to alternative parameterizations of the IPTM and alternative modeling approaches via regressions. For a randomly chosen document $D \in \{M, M+1, \dots, D\}$, we fit the IPTM to the corpus consisting of the first $d = \{1, \dots, D-1\}$ documents, then use the inferred posterior distributions to generate a distribution of predicted tie data for document D conditional on the content in document D , $\mathbf{w}^{(D)}$. A reasonable choice for M would be $D/2$, to assure a sufficient size training set for the first document in the test set. The variables that need to be sampled are the tie data $(i_o^{(D)}, J_o^{(D)}, t_o^{(D)})$, and we would compare the simulated ones to the observed data. Detailed pseudocode for generating predicted tie data using the IPTM is demonstrated in Appendix D.1.

6.1 Comparison Models

We first compare the ability of the IPTM in predicting tie data to that of the reduced parameterization of the IPTM. Specifically, by setting the number of interaction pattern $C = 1$ (i.e. $\{p_c^d\}_{d=1}^D = 1$), the IPTM reduces to pure network model which disregards the content reflection on network dynamics. However, this reduced version of the model still is innovative network model in that we jointly make inference on sender, receiver, and timestamp of the document. There exists numerous network modeling approaches in the literature, however, we are unaware of a model that can be used to jointly predict the sender, receivers, and timing of an e-mail/document conditional upon the interaction history. Therefore, the two main goals of this ablation study are 1) to test whether the reduced model itself serves as a novel network model with good predictive performance, and 2) to test whether including content information via interaction pattern provides additional advantage on document-ahead forecasting for tie data.

Next, we compare the IPTM to the comparative model that is built upon two regression models, in order to test if the IPTM or the ablated version of IPTM have any benefit over other existing models. To build an alternative modeling approach, we combine established existing models to make comparable predictions using the same inputs. The objective underlying our selection and use of comparison models is to evaluate the performance gains from jointly inferring the parameters that govern the generation of tie data—senders, receivers, and timing. Several models exist that could be used to model any of these three data types individually, but, to our knowledge, the literature does not offer any models that can be used to jointly generate all three types of tie data integrated into the IPTM. We train two separate models to use in predicting document D recipients, sender, and timing, respectively. These models are selected to closely mirror the structure of the corresponding components of the IPTM. In each model the documents used for training include documents 1 through $D-1$. The model of recipients will be a logistic regression model in which the dependent variable is the observed value of $J_{ij}^{(d)}$, where i indexes the observed sender, $i_o^{(d)}$. The network statistics, $\mathbf{x}_t^{(c)}(i, j)$, will be used as the covariates, with all $c = 1$ (i.e., all past interactions are within the same single interaction pattern in the comparison models). Specifically, the model used for predicting recipients of document D will have the form

$$P(J_{ij}^{(D)} = 1) = \frac{1}{1 + \exp(-(δ + \mathbf{b}^T \mathbf{x}_t(i, j)))}. \quad (9)$$

The model of when an email was sent is an exponential regression model in which the dependent variable is $t_i^{(D)}$. The single covariate in this model, assumed to have a coefficient of 1, is $\lambda_{iJ_i}^{(D)}(t)$, which is defined in Equation (6). We estimate an intercept term δ in this regression model in order to calibrate the scale of the model. Specifically, the probability density function of the model for the timing of the email is given by

$$f(t_i^{(D)}, \lambda_{iJ_i}^{(D)}(t), \eta) = \eta \lambda_{iJ_i}^{(D)}(t) \exp(-\eta \lambda_{iJ_i}^{(D)}(t) t_i^{(D)}), \quad (10)$$

where $\ln(\eta)$ is the intercept that calibrates the scale of the exponential distribution. The predicted sender of document D is determined by simulating $t_i^{(D)}$ for each i , and selecting the i that corresponds to the minimum value of $t_i^{(D)}$. The minimum value of $t_i^{(D)}$ is the predicted timing for document D , i.e. $t_o^{(D)}$, and similarly as the IPTM, the minimum time chooses the sender and recipient of the document D , i.e. $i_o^{(D)}$ and $J_o^{(D)}$. See Appendix D.2 for pseudocode.

6.2 Predictive Experiment Results

To verify that our model is applicable beyond the Dare County email data, we also performed the posterior predictive experiments using the Enron email data set —the most widely studied email data set. For this predictive experiment, we took a subset of the original data such that we only include emails between actors who sent over 100 emails, and actors who received over 100 emails from the chosen senders. E-mails that were not sent to at least one other active actor were discarded, and also preprocessed to remove any stop words, URLs, quoted text, and signatures. These steps resulted in a total of 3,925 emails involving 33 actors.

For both data sets, Dare County and Enron, we randomly selected 200 documents from the later half of the corpus and generated 100 samples of predicted tie data for every document D by running Algorithm 5 and Algorithm 6 ($O = 10$ and $R = 50$) with some burn-ins and thinnings ($b = 10$ and $q = 4$). Due to computational burden, we used the reduced set of network statistics (“intercept” and “dyadic”) for Enron, while we used full set of network statistics (“intercept” and “dyadic”, “degree”, “triadic”). Moreover, we ran the same predictive experiments with 21 unique combinations of the number of interaction patterns ($C = 1, 2, 3$) and the number of topics ($K = 2, 5, 10, 25, 50, 75, 100$) as a grid-search based hyperparameter selection process. We compare the predictions in terms of classification accuracy in predicting senders and receivers, as well as prediction error in document timing. Figure 7 presents the F_1 scores on sender predictions, multiclass version of the area under the ROC curve (AUC) measure (Hand and Till, 2001) on receiver predictions, and median absolute error (MAE) on time predictions for each document we predicted, all averaged over the entire samples.¹

Figure 7 demonstrates the ability IPTM in predicting the sender, receiver(s), and timetsamps of email. We currently do not have solid conclusion from this prediction experiment, however, we plan to further improve our model to achieve better predictability.

7 Topic Coherence

Topic coherence metrics Mimmo et al. (2011) are often used to evaluate the semantic coherence in topic models. In order to test whether the IPTM’s incorporation of network properties improves or impairs the ability of modeling text, we compared the coherence of topics inferred using our model with the coherence of topics inferred using the latent dirichlet allocation (LDA). Instead of re-fit the data using standard LDA algorithms, we used the topic assignments from the IPTM with $C = 1$, which simply makes the IPTM reduced to LDA in terms of topic assignments by unlinking the text and networks (same analogy as ablation study in Section 6.1). For each model, we varied the number of topics from 1 to 100 and drew five samples from the joint posterior distribution over the unobserved random variables in that model. We evaluated the topics resulting from each sample and averaged over the five samples, where the result is shown in Figure 8.² For Dare County data, our model works significantly better than LDA in terms of topic coherence, for any choice of the number of topics greater than or equal to 10. For Enron data, although the IPTM did not outperform LDA, still the topic coherence scores using the IPTM were not significantly different from LDA for all three cases, $C = 3$, $C = 2$, and $C = 1$. This result, when combined with the results in Section 6, demonstrates that our model can achieve goo predictive performance while producing coherent topics.

¹Number of documents currently used: Dare = 35 and Enron = 200.

²Currently, Enron data used two samples, varying the number of topics from 1 to 75.

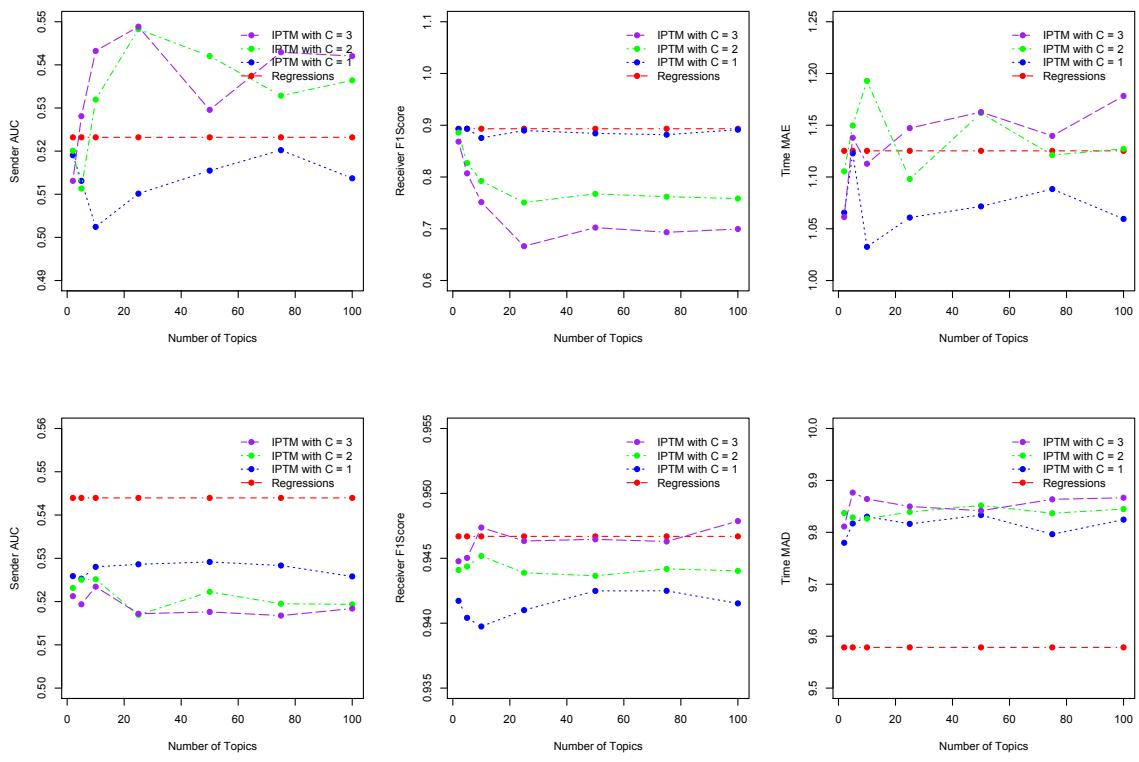


Figure 7: Average AUC, F1 score, MSE for the Dare County email network (upper) and the Enron data set (lower).

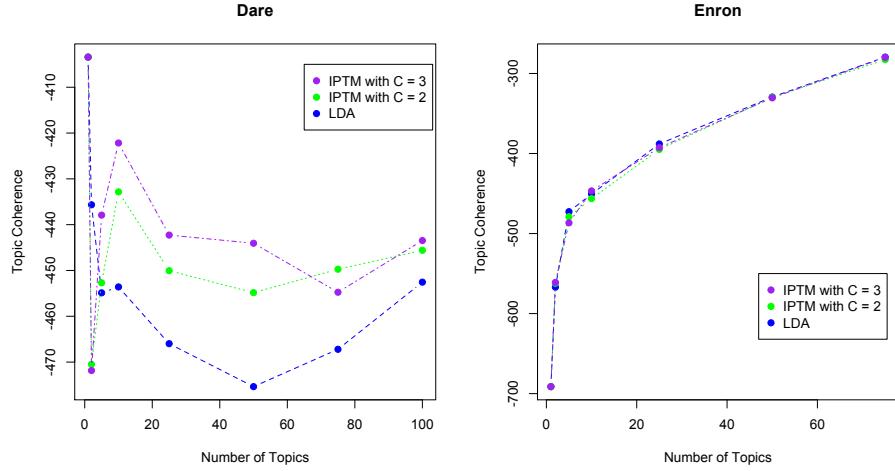


Figure 8: Average topic coherence scores for the Dare County email network (left) and the Enron dataset (right).

8 Posterior Predictive Checks

We finally perform posterior predictive checks Rubin et al. (1984) in order to evaluate the appropriateness of the model specification for Dare county. If the model is appropriate, the observed data should not be an outlier with respect to distributions of new data drawn from the posterior predictive distribution. To draw these comparisons, we first consider how to draw from the posterior predictive distribution. We produce a sample from the posterior predictive distribution by drawing new data conditional upon the parameters inferred in a single draw from the posterior distribution of the parameters, repeated over many draws from the posterior distribution.

Formally, we generate new data, which we denote $(\mathcal{I}_O^*, \mathcal{J}_O^*, \mathcal{T}_O^*, \mathcal{W}^*)$, conditional upon a set of inferred parameter values from

$$P(\mathcal{I}_O^*, \mathcal{J}_O^*, \mathcal{T}_O^*, \mathcal{W}^* | \mathcal{Z}, \mathcal{C}, \mathcal{B}, \delta, \mathbf{u}, \alpha, \mathbf{m}, \mu_{\mathbf{b}}, \Sigma_{\mathbf{b}}, \mu_{\delta}, \sigma_{\delta}^2), \quad (11)$$

and the pseudocode for generating new sets of data is given by Appendix E.

We use this posterior predictive checking to assess the extent to which our model is a “good fit” for the Dare County email network. For the test of goodness-of-fit in terms of network dynamics, we defined multiple network statistics that summarize meaningful aspects of Dare County data: indegree distribution for sender activities, outdegree distribution for receiver activities, document time-increment distributions, multicast (number of receivers) distribution, the edgewise shared partner distribution, and the geodesic distance distribution. For content-wise goodness-of-fit, we employed mutual information (MI) in Mimno and Blei (2011), which is often used to evaluate “bag of words” model assumptions. We then generated 1,000 synthetic networks and texts from the posterior predictive distribution implied by the IPTM and Dare County email data. We applied each discrepancy function to each synthetic network to yield the distributions over the values of the six network statistics and MI. If the IPTM is a suitable model for Dare County data, these distributions should be centered around the values of the corresponding discrepancy functions (except MI) when computed using the observed data.

Figure 9 illustrates the result of posterior predictive checks, showing IPTM’s goodness of fit for Dare County data. (Mutual information plots are in Appendix F)³ Four plots, indegree distribution, multicast distribution, geodesic distance distribution, and edgewise shared partner distribution, reveal that IPTM captures some important work features of the data, including spreadness and transitivity. However, outdegree distribution and time-increments from the simulated data are not quite similar to those of the observed data. The two aspects of data, sender and time-increments, are directly related from the time-generating process from Exponential distribution, thus we again learn that the current time-modeling framework of IPTM is not fully able to reflect or capture the time-generating process of real-world data. We plan on further dive into this issue and make adjustments in the future work.

³Number of samples currently used: 80 simulated corpus.

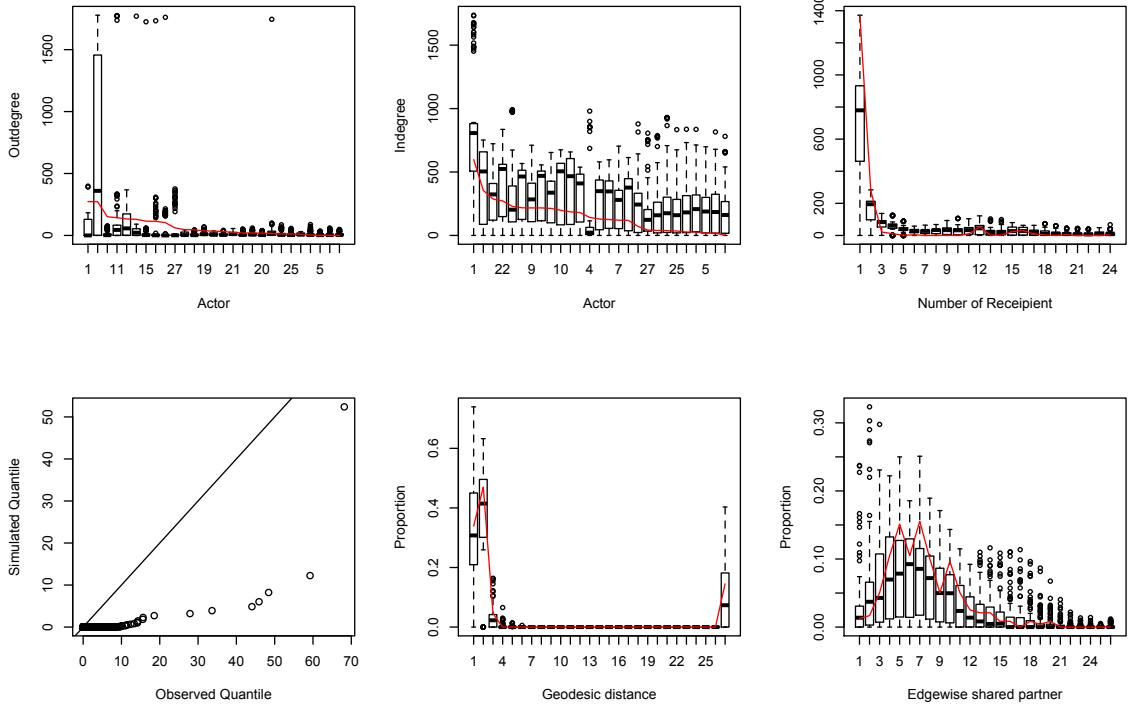


Figure 9: PPC results for Dare County data: (a) outdegree distribution, (b) indegree distribution, (c) multicast distribution, (d) QQplot of time-increments, (e) geodesic distance distribution, and (f) edgewise shared partner distribution. For (a) and (b), a box plot indicates the sampled degree of each actors in the synthetic networks, with actors sorted from highest to lowest observed degree and their observed degrees indicated by a line.

9 Conclusion

The IPTM is, to our knowledge, the first model to be capable of jointly modeling sender, receivers, time and contents in time stamped text valued networks. The IPTM incorporates innovative components, including the modeling of multicast tie formation and the conditioning of ERGM style network generative features on topic-based content. The application to North Carolina county government email data demonstrates, among other capabilities, the effectiveness at the IPTM in separating out both the content and relational structure underlying the normal day-to-day function of an organization and the management of a highly time-sensitive event—Hurricane Sandy. The IPTM is applicable to a variety of networks in which ties are attributed with textual documents. These include, for example, economic sanctions sent between countries and legislation attributed with sponsors and co-sponsors.

References

- Alemán, E. and Calvo, E. (2013). Explaining policy ties in presidential congresses: A network analysis of bill initiation data. *Political Studies*, 61(2):356–377.
- Altman, M., Gill, J., and McDonald, M. P. (2004). *Numerical issues in statistical computing for the social scientist*, volume 508. John Wiley & Sons.
- Barabási, A.-L. and Albert, R. (1999). Emergence of scaling in random networks. *science*, 286(5439):509–512.
- ben Aaron, J., Denny, M., Desmarais, B., and Wallach, H. (2017). Transparency by conformity: A field experiment evaluating openness in local governments. *Public Administration Review*, 77(1):68–77.
- Blei, D. M., Ng, A. Y., and Jordan, M. I. (2003). Latent dirichlet allocation. *J. Mach. Learn. Res.*, 3:993–1022.
- Bratton, K. A. and Rouse, S. M. (2011). Networks in the legislative arena: How group dynamics affect cosponsorship. *Legislative Studies Quarterly*, 36(3):423–460.
- Burda, Z., Jurkiewicz, J., and Krzywicki, A. (2004). Network transitivity and matrix models. *Physical Review E*, 69(2):026106.
- Burgess, A., Jackson, T., and Edwards, J. (2004). Email overload: Tolerance levels of employees within the workplace. In *Innovations Through Information Technology: 2004 Information Resources Management Association International Conference, New Orleans, Louisiana, USA, May 23-26, 2004*, volume 1, page 205. IGI Global.
- Butts, C. T. (2008). A relational event framework for social action. *Sociological Methodology*, 38(1):155–200.
- Camber Warren, T. (2010). The geometry of security: Modeling interstate alliances as evolving networks. *Journal of Peace Research*, 47(6):697–709.
- Chatterjee, S., Diaconis, P., et al. (2013). Estimating and understanding exponential random graph models. *The Annals of Statistics*, 41(5):2428–2461.
- Cranmer, S. J., Desmarais, B. A., and Kirkland, J. H. (2012a). Toward a network theory of alliance formation. *International Interactions*, 38(3):295–324.
- Cranmer, S. J., Desmarais, B. A., and Menninga, E. J. (2012b). Complex dependencies in the alliance network. *Conflict Management and Peace Science*, 29(3):279–313.
- Desmarais, B. A. and Cranmer, S. J. (2017). Statistical inference in political networks research. In Victor, J. N., Montgomery, A. H., and Lubell, M., editors, *The Oxford Handbook of Political Networks*. Oxford University Press.
- Fahmy, C. and Young, J. T. (2016). Gender inequality and knowledge production in criminology and criminal justice. *Journal of Criminal Justice Education*, pages 1–21.
- Fellows, I. and Handcock, M. (2017). Removing phase transitions from gibbs measures. In *Artificial Intelligence and Statistics*, pages 289–297.
- Geweke, J. (2004). Getting it right: Joint distribution tests of posterior simulators. *Journal of the American Statistical Association*, 99(467):799–804.
- Griffiths, T. (2002). Gibbs sampling in the generative model of latent dirichlet allocation.
- Griffiths, T. L. and Steyvers, M. (2004). Finding scientific topics. *Proceedings of the National academy of Sciences*, 101(suppl 1):5228–5235.
- Hammer, M. (1985). Implications of behavioral and cognitive reciprocity in social network data. *Social Networks*, 7(2):189–201.

- Hand, D. J. and Till, R. J. (2001). A simple generalisation of the area under the roc curve for multiple class classification problems. *Machine learning*, 45(2):171–186.
- Hunter, D. R., Handcock, M. S., Butts, C. T., Goodreau, S. M., and Morris, M. (2008). ergm: A package to fit, simulate and diagnose exponential-family models for networks. *Journal of statistical software*, 24(3):nihpa54860.
- Iribarren, J. L. and Moro, E. (2011). Branching dynamics of viral information spreading. *Physical Review E*, 84(4):046116.
- Jeong, H., Néda, Z., and Barabási, A.-L. (2003). Measuring preferential attachment in evolving networks. *EPL (Europhysics Letters)*, 61(4):567.
- Kanungo, S. and Jain, V. (2008). Modeling email use: a case of email system transition. *System Dynamics Review*, 24(3):299–319.
- Kinne, B. J. (2016). Agreeing to arm: Bilateral weapons agreements and the global arms trade. *Journal of Peace Research*, 53(3):359–377.
- Krafft, P., Moore, J., Desmarais, B., and Wallach, H. M. (2012). Topic-partitioned multinetword embeddings. In Pereira, F., Burges, C., Bottou, L., and Weinberger, K., editors, *Advances in Neural Information Processing Systems 25*, pages 2807–2815. Curran Associates, Inc.
- Kronegger, L., Mali, F., Ferligoj, A., and Doreian, P. (2011). Collaboration structures in slovenian scientific communities. *Scientometrics*, 90(2):631–647.
- Lai, C.-H., She, B., and Tao, C.-C. (2017). Connecting the dots: A longitudinal observation of relief organizations’ representational networks on social media. *Computers in Human Behavior*, 74:224–234.
- Liang, X. (2015). The changing impact of geographic distance: A preliminary analysis on the co-author networks in scientometrics (1983-2013). In *System Sciences (HICSS), 2015 48th Hawaii International Conference on*, pages 722–731. IEEE.
- Lim, K. W., Chen, C., and Buntine, W. (2013). Twitter-network topic model: A full bayesian treatment for social network and text modeling. In *NIPS2013 Topic Model workshop*, pages 1–5.
- Lin, Y., Desouza, K. C., and Roy, S. (2010). Measuring agility of networked organizational structures via network entropy and mutual information. *Applied Mathematics and Computation*, 216(10):2824–2836.
- Louch, H. (2000). Personal network integration: transitivity and homophily in strong-tie relations. *Social networks*, 22(1):45–64.
- McCallum, A., Corrada-Emmanuel, A., and Wang, X. (2005). The author-recipient-topic model for topic and role discovery in social networks, with application to enron and academic email. In *Workshop on Link Analysis, Counterterrorism and Security*, page 33.
- McCullough, B. D. (2009). The accuracy of econometric software. *Handbook of computational econometrics*, pages 55–79.
- Mimno, D. and Blei, D. (2011). Bayesian checking for topic models. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 227–237. Association for Computational Linguistics.
- Mimno, D., Wallach, H. M., Talley, E., Leenders, M., and McCallum, A. (2011). Optimizing semantic coherence in topic models. In *Proceedings of the conference on empirical methods in natural language processing*, pages 262–272. Association for Computational Linguistics.
- Peng, T.-Q., Liu, M., Wu, Y., and Liu, S. (2016). Follower-followee network, communication networks, and vote agreement of the us members of congress. *Communication Research*, 43(7):996–1024.
- Perry, P. O. and Wolfe, P. J. (2013). Point process modelling for directed interaction networks. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 75(5):821–849.

- Pew, R. C. (2016). Social media fact sheet. Accessed on 03/07/17.
- Rao, A. R. and Bandyopadhyay, S. (1987). Measures of reciprocity in a social network. *Sankhyā: The Indian Journal of Statistics, Series A*, pages 141–188.
- Robins, G., Pattison, P., Kalish, Y., and Lusher, D. (2007). An introduction to exponential random graph (p^*) models for social networks. *Social networks*, 29(2):173–191.
- Roca, C. P., Lozano, S., Arenas, A., and Sánchez, A. (2010). Topological traps control flow on real networks: The case of coordination failures. *PLoS One*, 5(12):e15210.
- Rubin, D. B. et al. (1984). Bayesianly justifiable and relevant frequency calculations for the applied statistician. *The Annals of Statistics*, 12(4):1151–1172.
- Szóstek, A. M. (2011). ?dealing with my emails?: Latent user needs in email management. *Computers in Human Behavior*, 27(2):723–729.
- Tadić, B. and Thurner, S. (2004). Information super-diffusion on structured networks. *Physica A: Statistical Mechanics and its Applications*, 332:566–584.
- Vázquez, A. (2003). Growing network with local rules: Preferential attachment, clustering hierarchy, and degree correlations. *Physical Review E*, 67(5):056104.
- Wallach, H. M. (2008). *Structured topic models for language*. PhD thesis, University of Cambridge.
- Yoon, H. Y. and Park, H. W. (2014). Strategies affecting twitter-based networking pattern of south korean politicians: social network analysis and exponential random graph model. *Quality & Quantity*, pages 1–15.

Appendix

A Normalizing constant of non-empty Gibbs measure

In Section 2.4, we define the non-empty Gibbs measure such that the probability of sender i selecting the binary receiver vector of length $(A - 1)$, $J_i^{(d)}$ is given by

$$P(J_i^{(d)}) = \frac{1}{Z(\delta, \log(\lambda_i^{(d)}))} \exp \left\{ \log(I(\|J_i^{(d)}\|_1 > 0)) + \sum_{j \in \mathcal{A}_{\setminus i}} (\delta + \log(\lambda_{ij}^{(d)})) J_{ij}^{(d)} \right\}.$$

To use this distribution efficiently, we derive a closed-form expression for $Z(\delta, \log(\lambda_i^{(d)}))$ that does not require brute-force summation over the support of $J_i^{(d)}$. We recognize that if $J_i^{(d)}$ were drawn via independent Bernoulli distributions in which $P(J_{ij}^{(d)}=1)$ was given by $\text{logit}(\delta + \lambda_{ij}^{(d)})$, then

$$P(J_i^{(d)}) \propto \exp \left\{ \sum_{j \in \mathcal{A}_{\setminus i}} (\delta + \log(\lambda_{ij}^{(d)})) J_{ij}^{(d)} \right\}.$$

This is straightforward to verify by looking at

$$\begin{aligned} P(J_{ij}^{(d)} = 1 | J_{i,-j}) &= \frac{\exp(\delta + \log(\lambda_{ij}^{(d)})) \exp \left\{ \sum_{h \neq i,j} (\delta + \log(\lambda_{ih}^{(d)})) J_{ih}^{(d)} \right\}}{\exp(\delta + \log(\lambda_{ij}^{(d)})) \exp \left\{ \sum_{h \neq i,j} (\delta + \log(\lambda_{ih}^{(d)})) J_{ih}^{(d)} \right\} + \exp(0) \exp \left\{ \sum_{h \neq i,j} (\delta + \log(\lambda_{ih}^{(d)})) J_{ih}^{(d)} \right\}}, \\ &= \frac{\exp(\delta + \log(\lambda_{ij}^{(d)}))}{\exp(\delta + \log(\lambda_{ij}^{(d)})) + 1}. \end{aligned}$$

We denote the logistic-Bernoulli normalizing constant as $Z^l(\delta, \lambda_i^{(d)})$, which is defined as

$$Z^l(\delta, \log(\lambda_i^{(d)})) = \sum_{J_i \in [0,1]^{(A-1)}} \exp \left\{ \sum_{j \neq i} (\delta + \log(\lambda_{ij}^{(d)})) J_{ij}^{(d)} \right\}.$$

Now, since

$$\exp \left\{ \log(I(\|J^{(d)}\|_1 > 0)) + \sum_{j \in \mathcal{A}_{\setminus i}} (\delta + \log(\lambda_{ij}^{(d)})) J_{ij}^{(d)} \right\} = \exp \left\{ \sum_{j \in \mathcal{A}_{\setminus i}} (\delta + \log(\lambda_{ij}^{(d)})) J_{ij}^{(d)} \right\},$$

except when $\|J^{(d)}\|_1 = 0$, in which case the left-hand side

$$\exp \left\{ \log(I(\|J^{(d)}\|_1 > 0)) + \sum_{j \in \mathcal{A}_{\setminus i}} (\delta + \log(\lambda_{ij}^{(d)})) J_{ij}^{(d)} \right\} = 0.$$

As such, we note that

$$\begin{aligned} Z(\delta, \log(\lambda_i^{(d)})) &= Z^l(\delta, \log(\lambda_i^{(d)})) - \exp \left\{ \sum_{j \in \mathcal{A}_{\setminus i}, J_{ij}^{(d)}=0} (\delta + \log(\lambda_{ij}^{(d)})) J_{ij}^{(d)} \right\} \\ &= Z^l(\delta, \log(\lambda_i^{(d)})) - 1. \end{aligned}$$

We can therefore derive a closed form expression for $Z(\delta, \log(\lambda_i^{(d)}))$ via a closed form expression for $Z^l(\delta, \log(\lambda_i^{(d)}))$. This can be done by looking at the probability of the zero vector under the logistic-Bernoulli model:

$$\begin{aligned} \frac{\exp \left\{ \sum_{j \neq i, J_{ij}^{(d)}=0} (\delta + \log(\lambda_{ij}^{(d)})) J_{ij}^{(d)} \right\}}{Z^l(\delta, \log(\lambda_i^{(d)}))} &= \prod_{j \in \mathcal{A}_{\setminus i}} \frac{\exp\{-(\delta + \log(\lambda_{ij}^{(d)}))\}}{\exp\{-(\delta + \log(\lambda_{ij}^{(d)}))\} + 1}, \\ Z^l(\delta, \log(\lambda_i^{(d)})) &= \frac{1}{\prod_{j \in \mathcal{A}_{\setminus i}} \frac{\exp\{-(\delta + \log(\lambda_{ij}^{(d)}))\}}{\exp\{-(\delta + \log(\lambda_{ij}^{(d)}))\} + 1}}. \end{aligned}$$

The closed form expression for the normalizing constant under the non-empty Gibbs measure is

$$Z(\delta, \lambda_i^{(d)}) = \left(\prod_{j \in \mathcal{A}_{\setminus i}} \left(\exp\{\delta + \log(\lambda_{ij}^{(d)})\} + 1 \right) \right) - 1.$$

B Sampling Equations

B.1 Joint distribution of latent and observed tie variables

As mentioned earlier in Section 2.4, we use data augmentation in the tie generating process. Since we should include both the observed and augmented data to make inferences on the related latent variables, the derivation steps for the contribution of tie data is not as simple as other variables. Therefore, here we provide the detailed derivation steps for the last term of joint posterior distribution in Equation (8):

$$\begin{aligned} & P(\mathcal{J}_a, \mathcal{T}_a, \mathcal{I}_o, \mathcal{J}_o, \mathcal{T}_o | \mathcal{Z}, \mathcal{C}, \mathcal{B}, \delta, \boldsymbol{\eta}, \sigma_T^2) \\ & \propto P(\mathcal{J}_a^{(1)}, \dots, \mathcal{J}_a^{(D)}, \mathcal{T}_a^{(1)}, \dots, \mathcal{T}_a^{(D)}, i_o^{(1)}, \dots, i_o^{(D)}, J_o^{(1)}, \dots, J_o^{(D)}, t_o^{(1)}, \dots, t_o^{(D)} | \mathcal{Z}, \mathcal{C}, \mathcal{B}, \delta, \boldsymbol{\eta}, \sigma_T^2) \quad (12) \\ & \propto P(\mathcal{J}_a^{(D)}, \mathcal{T}_a^{(D)}, i_o^{(D)}, J_o^{(D)}, t_o^{(D)} | \mathcal{I}_o^{(<D)}, \mathcal{J}_o^{(<D)}, \mathcal{T}_o^{(<D)}, \mathcal{Z}, \mathcal{C}, \mathcal{B}, \delta, \boldsymbol{\eta}, \sigma_T^2). \end{aligned}$$

Note that the joint likelihood of document $d = 1, \dots, D$ is proportional to the likelihood of D^{th} document, since the documents are not independent; all tie variables in D^{th} are conditioned on the earlier documents ($< D$) via network covariates \mathbf{x} or past interaction history.

Now we tackle the problem of evaluating the last line of Equation 12 by deriving the likelihood for d^{th} document, $P(\mathcal{J}_a^{(d)}, \mathcal{T}_a^{(d)}, i_o^{(d)}, J_o^{(d)}, t_o^{(d)} | \mathcal{I}_o^{(<d)}, \mathcal{J}_o^{(<d)}, \mathcal{T}_o^{(<d)}, \mathcal{Z}, \mathcal{C}, \mathcal{B}, \delta, \boldsymbol{\eta}, \sigma_T^2)$. There are three steps involved. First is the generation of the latent receivers J_i for each i ; second is the generation of time increments $\Delta T^{(d)} = t^{(d)} - t^{(d-1)}$ for each i ; and the last part is the simultaneous selection process of the observed sender, receivers, and timestamp, implying that the latent time increments generated from the latent sender-receiver pairs were greater than the observed time increment. Reflecting the three steps, the joint distribution is:

$$\begin{aligned} & P(\mathcal{J}_a^{(d)}, \mathcal{T}_a^{(d)}, i_o^{(d)}, J_o^{(d)}, t_o^{(d)} | \mathcal{I}_o^{(<d)}, \mathcal{J}_o^{(<d)}, \mathcal{T}_o^{(<d)}, \mathcal{Z}, \mathcal{C}, \mathcal{B}, \delta, \boldsymbol{\eta}, \sigma_T^2) \\ & = P(\text{latent receivers generation}) \times P(\text{observed time generation}) \times P(\text{observed time being the minimum}) \\ & = \prod_{i \in \mathcal{A}} \left(J_i^{(d)} \sim \text{Gibbs}(\{\lambda_{ij}^{(d)}\}_{j=1}^A, \delta) \right) \times \left(\Delta T_{i_o^{(d)} J_o^{(d)}}^{(d)} \sim \text{Lognormal}(\boldsymbol{\eta}^T \mathbf{y}_{i_o^{(d)} J_o^{(d)}}^{(d)}, \sigma_T^2) \right) \times P\left(\Delta T_{i_o^{(d)} J_o^{(d)}}^{(d)} < \min_{i \in \mathcal{A}_{\setminus i_o^{(d)}}} \Delta T_{i J_i}^{(d)}\right) \\ & = \prod_{i \in \mathcal{A}} \left(J_i^{(d)} \sim \text{Gibbs}(\{\lambda_{ij}^{(d)}\}_{j=1}^A, \delta) \right) \times \left(\Delta T_{i_o^{(d)} J_o^{(d)}}^{(d)} \sim \text{Lognormal}(\boldsymbol{\eta}^T \mathbf{y}_{i_o^{(d)} J_o^{(d)}}^{(d)}, \sigma_T^2) \right) \times \prod_{i \in \mathcal{A}_{\setminus i_o^{(d)}}} \left(1 - P(\Delta T_{i J_i}^{(d)} \leq \Delta T_{i_o^{(d)} J_o^{(d)}}^{(d)}) \right) \\ & = \left(\prod_{i \in \mathcal{A}} \frac{1}{Z(\delta, \log(\lambda_i^{(d)}))} \exp \left\{ \log(I(\sum_{j \in \mathcal{A}_{\setminus i}} J_{ij}^{(d)} > 0)) + \sum_{j \in \mathcal{A}_{\setminus i}} (\delta + \log(\lambda_{ij}^{(d)})) J_{ij}^{(d)} \right\} \right) \\ & \quad \times \left(\phi_l(\Delta T_{i_o^{(d)} J_o^{(d)}}^{(d)}; \boldsymbol{\eta}^T \mathbf{y}_{i_o^{(d)} J_o^{(d)}}^{(d)}, \sigma_T^2) \right) \times \left(\prod_{i \in \mathcal{A}_{\setminus i_o^{(d)}}} \left(1 - \Phi_l(\Delta T_{i_o^{(d)} J_o^{(d)}}^{(d)}; \boldsymbol{\eta}^T \mathbf{y}_{i J_i}^{(d)}, \sigma_T^2) \right) \right), \quad (13) \end{aligned}$$

where $\phi_l(\cdot; \mu, \sigma^2)$ is lognormal density function (p.d.f) with mean μ and variance σ^2 , and $\Phi_l(\cdot; \mu, \sigma^2)$ is lognormal cumulative density function (c.d.f), evaluating $P(x \leq \cdot)$ where x comes from lognormal distribution with mean μ and variance σ^2 . Plugging in the full Gibbs-measure equation,

$$\begin{aligned} & P(\mathcal{J}_a^{(d)}, \mathcal{T}_a^{(d)}, i_o^{(d)}, J_o^{(d)}, t_o^{(d)} | \mathcal{I}_o^{(<d)}, \mathcal{J}_o^{(<d)}, \mathcal{T}_o^{(<d)}, \mathcal{Z}, \mathcal{C}, \mathcal{B}, \delta, \boldsymbol{\eta}, \sigma_T^2) \\ & \propto \left(\prod_{i \in \mathcal{A}} \frac{1}{\left(\prod_{j \in \mathcal{A}_{\setminus i}} \left(\exp\{\delta + \log(\lambda_{ij}^{(d)})\} + 1 \right) \right) - 1} \exp \left\{ \sum_{j \in \mathcal{A}_{\setminus i}} (\delta + \log(\lambda_{ij}^{(d)})) J_{ij}^{(d)} \right\} \right) \\ & \quad \times \left(\phi_l(\Delta T_{i_o^{(d)} J_o^{(d)}}^{(d)}; \boldsymbol{\eta}^T \mathbf{y}_{i_o^{(d)} J_o^{(d)}}^{(d)}, \sigma_T^2) \right) \times \left(\prod_{i \in \mathcal{A}_{\setminus i_o^{(d)}}} \left(1 - \Phi_l(\Delta T_{i_o^{(d)} J_o^{(d)}}^{(d)}; \boldsymbol{\eta}^T \mathbf{y}_{i J_i}^{(d)}, \sigma_T^2) \right) \right), \quad (14) \end{aligned}$$

where this joint distribution can be interpreted as ‘probability of latent and observed edges from non-empty Gibbs measure \times probability of the observed time-increments from lognormal distribution \times probability of all latent time greater than the observed time.’ Finally for implementation, we need to compute these equations in log space to prevent numerical underflow.

B.2 Resampling \mathcal{J}_a

First of all, for each document d , we sample the latent sender-receiver(s) pairs as in pseudocode (Algorithm 6). That is, given the observed sender of the document $i_o^{(d)}$, we sample the latent receivers for each sender $i \in \mathcal{A}_{\setminus i_o^{(d)}}$. Here we illustrate how each sender-receiver pair in the document d is updated.

Define $\mathcal{J}_i^{(d)}$ be the $(A - 1)$ length random vector of indicators with its realization being $J_i^{(d)}$, representing the latent receivers corresponding to the sender i in the document d . For each latent sender i , we are going to resample $J_{ij}^{(d)}$, which is the j^{th} element of the receiver vector $J_i^{(d)}$, one at a time with random order. The full conditional probability of $J_{ij}^{(d)}$ is:

$$P(\mathcal{J}_{ij}^{(d)} = J_{ij}^{(d)} | \mathcal{J}_{i \setminus j}^{(d)}, \mathcal{Z}, \mathcal{C}, \mathcal{B}, \delta, \boldsymbol{\eta}, \sigma_T^2, \mathcal{W}, \mathcal{J}_{a,-i}, \mathcal{T}_{a,i}^{(d)}, \mathcal{I}_O, \mathcal{J}_O, \mathcal{T}_O, \beta, \mathbf{u}, \alpha, \mathbf{m}, \mu_b, \Sigma_b, \mu_\delta, \sigma_\delta^2, \mu_\eta, \Sigma_\eta, a_T, b_T), \quad (15)$$

which we can drop some independent terms and move to

$$\begin{aligned} P(\mathcal{J}_{ij}^{(d)} = J_{ij}^{(d)} | \mathcal{J}_{i \setminus j}^{(d)}, \mathcal{T}_{a,i}^{(d)}, i_o^{(d)}, J_O^{(d)}, t_O^{(d)}, \mathcal{I}_O^{(<d)}, \mathcal{J}_O^{(<d)}, \mathcal{T}_O^{(<d)}, \mathcal{Z}, \mathcal{C}, \mathcal{B}, \delta, \boldsymbol{\eta}, \sigma_T^2) \\ \propto P(\mathcal{J}_{ij}^{(d)} = J_{ij}^{(d)}, \mathcal{J}_{i \setminus j}^{(d)}, \mathcal{T}_{a,i}^{(d)}, i_o^{(d)}, J_O^{(d)}, t_O^{(d)} | \mathcal{I}_O^{(<d)}, \mathcal{J}_O^{(<d)}, \mathcal{T}_O^{(<d)}, \mathcal{Z}, \mathcal{C}, \mathcal{B}, \delta, \boldsymbol{\eta}, \sigma_T^2) \\ \propto \left(\frac{1}{\left(\prod_{j \in \mathcal{A}_{\setminus i}} \left(\exp\{\delta + \log(\lambda_{ij}^{(d)})\} + 1 \right) \right) - 1} \exp \left\{ \log(I(\sum_{j \in \mathcal{A}_{\setminus i}} J_{ij}^{(d)} > 0)) + \sum_{j \in \mathcal{A}_{\setminus i}} (\delta + \log(\lambda_{ij}^{(d)})) J_{ij}^{(d)} \right\} \right) \\ \times \left(\phi_l(\Delta T_{i_o^{(d)} J_o^{(d)}}^{(d)}; \boldsymbol{\eta}^T \mathbf{y}_{i_o^{(d)} J_o^{(d)}}^{(d)}, \sigma_T^2) \right) \times \left(1 - \Phi_l(\Delta T_{i_o^{(d)} J_o^{(d)}}^{(d)}; \boldsymbol{\eta}^T \mathbf{y}_{i J_i}^{(d)}, \sigma_T^2) \right) \\ \propto \left(\exp \left\{ \log(I(\|J_i^{(d)}\|_1 > 0)) + \sum_{j \in \mathcal{A}_{\setminus i}} (\delta + \log(\lambda_{ij}^{(d)})) J_{ij}^{(d)} \right\} \right) \times \left(1 - \Phi_l(\Delta T_{i_o^{(d)} J_o^{(d)}}^{(d)}; \boldsymbol{\eta}^T \mathbf{y}_{i J_i}^{(d)}, \sigma_T^2) \right), \end{aligned} \quad (16)$$

where we replace typical use of $(-d)$ to $(< d)$ on the right hand side, due to the fact that $d^{(th)}$ document only depends on the past documents. The last line of Equation (16) is obtained by dropping the terms that do not include $J_{ij}^{(d)}$, such as the normalizing constant of Gibbs measure.

To be more specific, since $J_{ij}^{(d)}$ could be either 1 or 0, we divide into two cases as below:

$$\begin{aligned} P(\mathcal{J}_{ij}^{(d)} = 1 | \mathcal{J}_{i \setminus j}^{(d)}, \mathcal{T}_{a,i}^{(d)}, i_o^{(d)}, J_O^{(d)}, t_O^{(d)}, \mathcal{I}_O^{(<d)}, \mathcal{J}_O^{(<d)}, \mathcal{T}_O^{(<d)}, \mathcal{Z}, \mathcal{C}, \mathcal{B}, \delta, \boldsymbol{\eta}, \sigma_T^2) \\ \propto \exp \left(\log(1) + \sum_{j \in \mathcal{A}_{\setminus i}} (\delta + \log(\lambda_{ij}^{(d)})) J_{ij}^{(d)} \right) \times \left(1 - \Phi_l(\Delta T_{i_o^{(d)} J_o^{(d)}}^{(d)}; \boldsymbol{\eta}^T \mathbf{y}_{i J_{i+[j]}}^{(d)}, \sigma_T^2) \right) \\ \propto \exp \left(\delta + \log(\lambda_{ij}^{(d)}) \right) \times \left(1 - \Phi_l(\Delta T_{i_o^{(d)} J_o^{(d)}}^{(d)}; \boldsymbol{\eta}^T \mathbf{y}_{i J_{i+[j]}}^{(d)}, \sigma_T^2) \right), \end{aligned} \quad (17)$$

where $J_{i+[j]}^{(d)}$ meaning that the j^{th} element of $J_i^{(d)}$ is fixed as 1 (thus making $\log(I(\|J_i^{(d)}\|_1 > 0)) = 0$ for sure). On the other hand,

$$\begin{aligned} P(\mathcal{J}_{ij}^{(d)} = 0 | \mathcal{J}_{i \setminus j}^{(d)}, \mathcal{T}_{a,i}^{(d)}, i_o^{(d)}, J_O^{(d)}, t_O^{(d)}, \mathcal{I}_O^{(<d)}, \mathcal{J}_O^{(<d)}, \mathcal{T}_O^{(<d)}, \mathcal{Z}, \mathcal{C}, \mathcal{B}, \delta, \boldsymbol{\eta}, \sigma_T^2) \\ \propto \exp \left(\log(I(\|J_{i[-j]}^{(d)}\|_1 > 0)) + \sum_{j \in \mathcal{A}_{\setminus i}} (\delta + \log(\lambda_{ij}^{(d)})) J_{ij}^{(d)} \right) \times \left(1 - \Phi_l(\Delta T_{i_o^{(d)} J_o^{(d)}}^{(d)}; \boldsymbol{\eta}^T \mathbf{y}_{i J_{i[-j]}}^{(d)}, \sigma_T^2) \right) \\ \propto \exp \left(\log(I(\|J_{i[-j]}^{(d)}\|_1 > 0)) \right) \times \left(1 - \Phi_l(\Delta T_{i_o^{(d)} J_o^{(d)}}^{(d)}; \boldsymbol{\eta}^T \mathbf{y}_{i J_{i[-j]}}^{(d)}, \sigma_T^2) \right), \end{aligned} \quad (18)$$

where $J_{i[-j]}^{(d)}$ meaning similarly that the j^{th} element of $J_i^{(d)}$ is fixed as 0. In this case, we cannot guarantee $I(\|J_{i[-j]}^{(d)}\|_1 > 0) = 1$, so we have to leave the term. When it is zero, $\exp\{\log(I(\|J_{i[-j]}^{(d)}\|_1 >$

$0)\} = 0$, thus we will sample 1 with probability 1. From this property of non-empty Gibbs measure, we prevent from the instances where the sender has no recipients to send the document. Now we can use multinomial sampling using the two probabilities, Equation (17) and Equation (18).

B.3 Resampling \mathcal{Z}

Second, we resample the topic assignments, one words in a document at a time. The new values of $z_n^{(d)}$ are sampled using the conditional posterior probability of being topic k , and we derive the sampling equation from the conditional distribution in Latent Dirichlet allocation (Blei et al., 2003):

$$\begin{aligned} & P(\mathbf{w}^{(d)}, \mathbf{z}^{(d)} | \mathcal{W}_{\setminus d}, \mathcal{Z}_{\setminus d}, \beta, \mathbf{u}, \alpha, \mathbf{m}) \\ & \propto \prod_{n=1}^{N^{(d)}} P(z_n^{(d)} = k, w_n^{(d)} = w | \mathcal{W}_{\setminus d, n}, \mathcal{Z}_{\setminus d, n}, \beta, \mathbf{u}, \alpha, \mathbf{m}). \end{aligned} \quad (19)$$

To obtain the Gibbs sampling equation, we need to obtain an expression for $P(z_n^{(d)} = k, w_n^{(d)} = w | \mathcal{W}_{\setminus d, n}, \mathcal{Z}_{\setminus d, n}, \beta, \mathbf{u}, \alpha, \mathbf{m})$. From Bayes' theorem and Gamma identity $\Gamma(k+1) = k\Gamma(k)$,

$$\begin{aligned} & P(z_n^{(d)} = k, w_n^{(d)} = w | \mathcal{W}_{\setminus d, n}, \mathcal{Z}_{\setminus d, n}, \beta, \mathbf{u}, \alpha, \mathbf{m}) \\ & \propto \frac{P(\mathcal{W}, \mathcal{Z} | \beta, \mathbf{u}, \alpha, \mathbf{m})}{P(\mathcal{W}_{\setminus d, n}, \mathcal{Z}_{\setminus d, n} | \beta, \mathbf{u}, \alpha, \mathbf{m})} \\ & \propto \frac{\prod_{k=1}^K \frac{\prod_{w=1}^W \Gamma(N_{wk}^{WK} + \beta u_w)}{\Gamma(\sum_{w=1}^W N_{wk}^{WK} + \beta)} \times \prod_{k=1}^K \frac{\Gamma(N_{k|d} + \alpha m_k)}{\Gamma(N_{\cdot|d} + \alpha)}}{\prod_{k=1}^K \frac{\prod_{w=1}^W \Gamma(N_{wk, \setminus d, n}^{WK} + \beta u_w)}{\Gamma(\sum_{w=1}^W N_{wk, \setminus d, n}^{WK} + \beta)} \times \prod_{k=1}^K \frac{\Gamma(N_{k|d, \setminus d, n} + \alpha m_k)}{\Gamma(N_{\cdot|d, \setminus d, n} + \alpha)}} \\ & \propto \frac{N_{wk, \setminus d, n}^{WK} + \frac{\beta}{W}}{\sum_{w=1}^W N_{wk, \setminus d, n}^{WK} + \beta} \times \frac{N_{k|d, \setminus d, n} + \alpha m_k}{N^{(d)} - 1 + \alpha}. \end{aligned} \quad (20)$$

Then, same as for LDA, we also know that the topic assignment $z_n^{(d)} = k$ is obtained by:

$$P(z_n^{(d)} = k | w_n^{(d)} = w, \mathcal{W}_{\setminus d, n}, \mathcal{Z}_{\setminus d, n}, \beta, \mathbf{u}, \alpha, \mathbf{m}) \propto \frac{N_{k|d, \setminus d, n} + \alpha m_k}{N^{(d)} - 1 + \alpha}$$

Now, considering the modeling framework of IPTM, we re-derive the sampling equation reflecting the network effects as well:

$$\begin{aligned} & P(z_n^{(d)} = k | \mathcal{Z}_{\setminus d, n}, \mathcal{C}, \mathcal{B}, \delta, \boldsymbol{\eta}, \sigma_T^2, \mathcal{W}, \mathcal{J}_a, \mathcal{T}_a, \mathcal{I}_a, \mathcal{J}_o, \mathcal{T}_o, \mathcal{I}_o, \beta, \mathbf{u}, \alpha, \mathbf{m}, \mu_b, \Sigma_b, \mu_\delta, \sigma_\delta^2, \mu_\eta, \Sigma_\eta, a_T, b_T) \\ & \propto P(z_n^{(d)} = k, w_n^{(d)} | \mathcal{J}_a^{(\geq d)}, \mathcal{T}_a^{(>d)}, i_o^{(\geq d)}, J_o^{(\geq d)}, t_o^{(\geq d)} | \mathcal{Z}_{\setminus d, n}, \mathcal{C}, \mathcal{B}, \delta, \boldsymbol{\eta}, \sigma_T^2, \mathcal{W}_{\setminus d, n}, \mathcal{I}_o^{(<d)}, \mathcal{J}_o^{(<d)}, \mathcal{T}_o^{(<d)}, \beta, \mathbf{u}, \alpha, \mathbf{m}) \\ & \propto P(z_n^{(d)} = k | \mathcal{Z}_{\setminus d, n}, \alpha, \mathbf{m}) P(w_n^{(d)} | z_n^{(d)} = k, \mathcal{W}_{\setminus d, n}, \mathcal{Z}_{\setminus d, n}, \beta, \mathbf{u}) \\ & \quad \times P(\mathcal{J}_a^{(d*)}, \mathcal{T}_a^{(d*)}, i_o^{(d*)}, J_o^{(d*)}, t_o^{(d*)} | z_n^{(d)} = k, \mathcal{Z}_{\setminus d, n}, \mathcal{C}, \mathcal{B}, \delta, \boldsymbol{\eta}, \sigma_T^2), \end{aligned} \quad (21)$$

where the subscript “ $\setminus d, n$ ” denotes the exclusion of position n in d^{th} document. Note that since selecting a topic for any token influences the histories acting on documents from the current one d to the future ones with the timepoints less than $t_o^{(d)} + 384$, we define $d^* = \text{argmax}_{d'} \{t_o^{(d')} \leq t_o^{(d)} + 384\}$ to correctly evaluate tie contribution part. From Equation (22), we know that:

$$P(z_n^{(d)} = k | \mathcal{Z}_{\setminus d, n}, \alpha, \mathbf{m}) = \frac{N_{\setminus d, n}^{(k|d)} + \alpha m_k}{N^{(d)} - 1 + \alpha} \quad (22)$$

which is the well-known form of collapsed Gibbs sampling equation for LDA. We also know that

$$P(w_n^{(d)} | z_n^{(d)} = k, \mathcal{W}_{\setminus d, n}, \mathcal{Z}_{\setminus d, n}, \beta, \mathbf{u}) = \frac{N_{\setminus d, n}^{(w_n^{(d)}|k)} + \frac{\beta}{W}}{N_{\setminus d, n}^{(k)} + \beta}, \quad (23)$$

where $N^{(w_n^{(d)}|k)}$ is the number of tokens assigned to topic k whose type is the same as that of $w_n^{(d)}$, excluding $w_n^{(d)}$ itself, and $N_{\setminus d,n}^{(k)} = \sum_{w=1}^W N_{\setminus d,n}^{(w^{(d)}|k)}$. We already have shown in Section B.1 that $P(\mathcal{J}_a^{(d)}, i_o^{(d)}, J_o^{(d)}, t_o^{(d)} | z_n^{(d)} = k, \mathcal{Z}_{\setminus d,n}, \mathcal{C}, \mathcal{B}, \delta, \boldsymbol{\eta}, \sigma_T^2)$ is Equation (14) where every part includes $p_c^{(d)}$ such that we cannot simplify any further. Therefore, assuming $N^{(d)} > 0$, the conditional probability of n^{th} word in document d being topic k is:

$$\begin{aligned} P(z_n^{(d)} = k | \mathcal{Z}_{\setminus d,n}, \mathcal{C}, \mathcal{B}, \delta, \boldsymbol{\eta}, \sigma_T^2, \mathcal{W}, \mathcal{J}_a, \mathcal{T}_a, \mathcal{I}_o, \mathcal{J}_o, \mathcal{T}_o, \beta, \mathbf{u}, \alpha, \mathbf{m}, \mu_b, \Sigma_b, \mu_\delta, \sigma_\delta^2, \mu_\eta, \Sigma_\eta, a_T, b_T) \\ \propto (N_{\setminus d,n}^{(k|d)} + \alpha m_k) \times \frac{N_{\setminus d,n}^{(w_n^{(d)}|k)} + \frac{\beta}{W}}{N_{\setminus d,n}^{(k)} + \beta} \\ \times \left(\prod_{i \in \mathcal{A}} \frac{1}{\left(\prod_{j \in \mathcal{A}_{\setminus i}} \left(\exp\{\delta + \log(\lambda_{ij}^{(d*)})\} + 1 \right) \right) - 1} \exp \left\{ \sum_{j \in \mathcal{A}_{\setminus i}} (\delta + \log(\lambda_{ij}^{(d*)})) J_{ij}^{(d*)} \right\} \right) \\ \times \left(\phi_l(\Delta T_{i_o^{(d*)} J_o^{(d*)}}^{(d*)}; \boldsymbol{\eta}^T \mathbf{y}_{i_o^{(d*)} J_o^{(d*)}}^{(d*)}, \sigma_T^2) \right) \times \left(\prod_{i \in \mathcal{A}_{\setminus i_o^{(d*)}}} (1 - \Phi_l(\Delta T_{i_o^{(d*)} J_o^{(d*)}}^{(d*)}; \boldsymbol{\eta}^T \mathbf{y}_{i J_i}^{(d*)}, \sigma_T^2)) \right). \end{aligned} \quad (24)$$

B.4 Resampling \mathcal{C}

The next variable to resample is the topic-interaction pattern assignments, one topic at a time. We derive the posterior conditional probability for the interaction pattern \mathcal{C} for k^{th} topic as below:

$$\begin{aligned} P(c_k = c | \mathcal{Z}, \mathcal{C}_{\setminus k}, \mathcal{B}, \delta, \boldsymbol{\eta}, \sigma_T^2, \mathcal{W}, \mathcal{J}_a, \mathcal{T}_a, \mathcal{I}_o, \mathcal{J}_o, \mathcal{T}_o, \beta, \mathbf{u}, \alpha, \mathbf{m}, \mu_b, \Sigma_b, \mu_\delta, \sigma_\delta^2, \mu_\eta, \Sigma_\eta, a_T, b_T) \\ \propto P(c_k = c, \mathcal{J}_a, \mathcal{T}_a, \mathcal{I}_o, \mathcal{J}_o, \mathcal{T}_o | \mathcal{Z}, \mathcal{C}_{\setminus k}, \mathcal{B}, \delta, \boldsymbol{\eta}, \sigma_T^2) \\ \propto P(c_k = c) P(\mathcal{J}_a, \mathcal{T}_a, \mathcal{I}_o, \mathcal{J}_o, \mathcal{T}_o | \mathcal{Z}, c_k = c, \mathcal{C}_{\setminus k}, \mathcal{B}, \delta, \boldsymbol{\eta}, \sigma_T^2) \end{aligned} \quad (25)$$

where $P(c_k = c) = \frac{1}{C}$ so this term disappears. Therefore, throughout $c_k = c$:

$$\begin{aligned} P(c_k = c | \mathcal{Z}, \mathcal{C}_{\setminus k}, \mathcal{B}, \delta, \boldsymbol{\eta}, \sigma_T^2, \mathcal{W}, \mathcal{J}_a, \mathcal{T}_a, \mathcal{I}_o, \mathcal{J}_o, \mathcal{T}_o, \beta, \mathbf{u}, \alpha, \mathbf{m}, \mu_b, \Sigma_b, \mu_\delta, \sigma_\delta^2, \mu_\eta, \Sigma_\eta, a_T, b_T) \\ \propto P(\mathcal{J}_a, \mathcal{T}_a, \mathcal{I}_o, \mathcal{J}_o, \mathcal{T}_o | \mathcal{Z}, c_k = c, \mathcal{C}_{\setminus k}, \mathcal{B}, \delta, \boldsymbol{\eta}, \sigma_T^2) \\ = \left(\prod_{i \in \mathcal{A}} \frac{1}{\left(\prod_{j \in \mathcal{A}_{\setminus i}} \left(\exp\{\delta + \log(\lambda_{ij}^{(D)})\} + 1 \right) \right) - 1} \exp \left\{ \sum_{j \in \mathcal{A}_{\setminus i}} (\delta + \log(\lambda_{ij}^{(D)})) J_{ij}^{(D)} \right\} \right) \\ \times \left(\phi_l(\Delta T_{i_o^{(D)} J_o^{(D)}}^{(D)}; \boldsymbol{\eta}^T \mathbf{y}_{i_o^{(D)} J_o^{(D)}}^{(D)}, \sigma_T^2) \right) \times \left(\prod_{i \in \mathcal{A}_{\setminus i_o^{(D)}}} (1 - \Phi_l(\Delta T_{i_o^{(D)} J_o^{(D)}}^{(D)}; \boldsymbol{\eta}^T \mathbf{y}_{i J_i}^{(D)}, \sigma_T^2)) \right). \end{aligned} \quad (26)$$

Note that for this one, \mathcal{C} , we need to consider all documents so we use the likelihood of the last D^{th} document.

B.5 Resampling \mathcal{B} and δ

Next, we jointly update $\mathcal{B} = \{\mathbf{b}^{(c)}\}_{c=1}^C$ and δ . We use the Metropolis-Hastings algorithm with a proposal density Q being the multivariate Gaussian distribution, with a diagonal covariance matrix multiplied by σ_Q^2 (proposal distribution variance parameters set by the user), centered on the current values of $\mathcal{B} = \{\mathbf{b}^{(c)}\}_{c=1}^C$ and δ . Under the symmetric proposals, we cancel out Q-ratio and then accept the new proposed value $\mathcal{B}' = \{\mathbf{b}'^{(c)}\}_{c=1}^C$ and δ' with probability equal to:

$$\text{Acceptance Probability} = \begin{cases} \frac{P(\mathcal{B}', \delta' | \mathcal{Z}, \mathcal{C}, \boldsymbol{\eta}, \sigma_T^2, \mathcal{W}, \mathcal{J}_a, \mathcal{T}_a, \mathcal{I}_o, \mathcal{J}_o, \mathcal{T}_o, \beta, \mathbf{u}, \alpha, \mathbf{m}, \mu_b, \Sigma_b, \mu_\delta, \sigma_\delta^2, \mu_\eta, \Sigma_\eta, a_T, b_T)}{P(\mathcal{B}, \delta | \mathcal{Z}, \mathcal{C}, \boldsymbol{\eta}, \sigma_T^2, \mathcal{W}, \mathcal{J}_a, \mathcal{T}_a, \mathcal{I}_o, \mathcal{J}_o, \mathcal{T}_o, \beta, \mathbf{u}, \alpha, \mathbf{m}, \mu_b, \Sigma_b, \mu_\delta, \sigma_\delta^2, \mu_\eta, \Sigma_\eta, a_T, b_T)} & \text{if } < 1 \\ 1 & \text{else} \end{cases} \quad (27)$$

After factorization, we get

$$\begin{aligned}
& \frac{P(\mathcal{B}', \delta' | \mathcal{Z}, \mathcal{C}, \boldsymbol{\eta}, \sigma_T^2, \mathcal{W}, \mathcal{J}_a, \mathcal{T}_a, \mathcal{I}_o, \mathcal{T}_o, \beta, \mathbf{u}, \alpha, \mathbf{m}, \mu_b, \Sigma_b, \mu_\delta, \sigma_\delta^2, \mu_{\boldsymbol{\eta}}, \Sigma_{\boldsymbol{\eta}}, a_T, b_T)}{P(\mathcal{B}, \delta | \mathcal{Z}, \mathcal{C}, \boldsymbol{\eta}, \sigma_T^2, \mathcal{W}, \mathcal{J}_a, \mathcal{T}_a, \mathcal{I}_o, \mathcal{T}_o, \beta, \mathbf{u}, \alpha, \mathbf{m}, \mu_b, \Sigma_b, \mu_\delta, \sigma_\delta^2, \mu_{\boldsymbol{\eta}}, \Sigma_{\boldsymbol{\eta}}, a_T, b_T)} \\
&= \frac{P(\mathcal{Z}, \mathcal{C}, \mathcal{B}', \delta, \boldsymbol{\eta}, \sigma_T^2, \mathcal{W}, \mathcal{J}_a, \mathcal{T}_a, \mathcal{I}_o, \mathcal{T}_o | \beta, \mathbf{u}, \alpha, \mathbf{m}, \mu_b, \Sigma_b, \mu_\delta, \sigma_\delta^2, \mu_{\boldsymbol{\eta}}, \Sigma_{\boldsymbol{\eta}}, a_T, b_T)}{P(\mathcal{Z}, \mathcal{C}, \mathcal{B}, \delta, \boldsymbol{\eta}, \sigma_T^2, \mathcal{W}, \mathcal{J}_a, \mathcal{T}_a, \mathcal{I}_o, \mathcal{T}_o | \beta, \mathbf{u}, \alpha, \mathbf{m}, \mu_b, \Sigma_b, \mu_\delta, \sigma_\delta^2, \mu_{\boldsymbol{\eta}}, \Sigma_{\boldsymbol{\eta}}, a_T, b_T)} \quad (28) \\
&= \frac{P(\mathcal{B}' | \mathcal{C}, \mu_b, \Sigma_b) P(\delta' | \mu_\delta, \sigma_\delta^2) P(\mathcal{J}_a, \mathcal{T}_a, \mathcal{I}_o, \mathcal{T}_o | \mathcal{Z}, \mathcal{C}, \mathcal{B}', \delta, \boldsymbol{\eta}, \sigma_T^2)}{P(\mathcal{B} | \mathcal{C}, \mu_b, \Sigma_b) P(\delta | \mu_\delta, \sigma_\delta^2) P(\mathcal{J}_a, \mathcal{T}_a, \mathcal{I}_o, \mathcal{T}_o | \mathcal{Z}, \mathcal{C}, \mathcal{B}, \delta, \boldsymbol{\eta}, \sigma_T^2)},
\end{aligned}$$

where $P(\mathcal{B} | \mathcal{C}, \mu_b, \Sigma_b)$ is calculated from the product of $\mathbf{b}^{(c)} \sim \text{Multivariate Normal}(\mu_b, \Sigma_b)$ over the interaction patterns $c \in \{1, \dots, C\}$, $P(\delta | \mu_\delta, \sigma_\delta^2)$ is also calculated from $\delta \sim \text{Normal}(\mu_\delta, \sigma_\delta^2)$, and $P(\mathcal{J}_a, \mathcal{T}_a, \mathcal{I}_o, \mathcal{T}_o | \mathcal{Z}, \mathcal{C}, \mathcal{B}, \delta, \boldsymbol{\eta}, \sigma_T^2)$ is the same as the log of Equation (26), without the time components that does not depend on either \mathcal{B} or δ . Again, we take the log and obtain:

$$\begin{aligned}
& \sum_{c=1}^C \log(\mathcal{N}(\mathbf{b}'^{(c)}; \mu_b, \Sigma_b)) - \sum_{c=1}^C \log(\mathcal{N}(\mathbf{b}^{(c)}; \mu_b, \Sigma_b)) + \log(\mathcal{N}(\delta'; \mu_\delta, \sigma_\delta^2)) - \log(\mathcal{N}(\delta; \mu_\delta, \sigma_\delta^2)) \\
&+ \left(\left(\sum_{i \in \mathcal{A}} \left(-\log \left(\prod_{j \in \mathcal{A}_{\setminus i}} (\exp\{\delta' + \log(\lambda_{ij}^{(D)})\} + 1) - 1 \right) + \sum_{j \in \mathcal{A}_{\setminus i}} (\delta' + \log(\lambda_{ij}^{(D)})) J_{ij}^{(D)} \right) \text{ given } \mathbf{b}' \right) \right. \\
&- \left. \left(\left(\sum_{i \in \mathcal{A}} \left(-\log \left(\prod_{j \in \mathcal{A}_{\setminus i}} (\exp\{\delta + \log(\lambda_{ij}^{(D)})\} + 1) - 1 \right) + \sum_{j \in \mathcal{A}_{\setminus i}} (\delta + \log(\lambda_{ij}^{(D)})) J_{ij}^{(D)} \right) \text{ given } \mathbf{b} \right) \right),
\end{aligned} \quad (29)$$

where \mathcal{N} is the multivariate normal density. If the log of a sample from Uniform(0,1) is less than the log-acceptance probability (29), we accept the proposal \mathbf{b}' and δ' . Otherwise, we reject.

B.6 Resampling $\boldsymbol{\eta}$ and σ_T^2

Lastly, we also jointly update $\boldsymbol{\eta}$ and σ_T^2 , time-increment related lognormal mean and variance parameters. Again, we use the Metropolis-Hastings algorithm with the multivariate Gaussian distribution proposals. We cancel out Q-ratio and then accept the new $\boldsymbol{\eta}'$ and $\sigma_T^{2'}$ with probability equal to:

$$\text{Acceptance Probability} = \begin{cases} \frac{P(\boldsymbol{\eta}' | \mu_{\boldsymbol{\eta}}, \Sigma_{\boldsymbol{\eta}}) P(\sigma_T^{2'} | a_T, b_T) P(\mathcal{J}_a, \mathcal{T}_a, \mathcal{I}_o, \mathcal{T}_o | \mathcal{Z}, \mathcal{C}, \mathcal{B}', \delta, \boldsymbol{\eta}, \sigma_T^2)}{P(\boldsymbol{\eta} | \mathcal{C}, \mu_{\boldsymbol{\eta}}, \Sigma_{\boldsymbol{\eta}}) P(\sigma_T^2 | a_T, b_T) P(\mathcal{J}_a, \mathcal{T}_a, \mathcal{I}_o, \mathcal{T}_o | \mathcal{Z}, \mathcal{C}, \mathcal{B}, \delta, \boldsymbol{\eta}, \sigma_T^2)} & \text{if } < 1 \\ 1 & \text{else} \end{cases} \quad (30)$$

where $P(\boldsymbol{\eta} | \mathcal{C}, \mu_{\boldsymbol{\eta}}, \Sigma_{\boldsymbol{\eta}})$ is calculated from $\boldsymbol{\eta} \sim \text{Multivariate Normal}(\mu_{\boldsymbol{\eta}}, \Sigma_{\boldsymbol{\eta}})$, $P(\sigma_T^2 | a_T, b_T)$ is calculated from $\sigma_T^2 \sim \text{Inverse-Gamma}(a_T, b_T)$, and $P(\mathcal{J}_a, \mathcal{T}_a, \mathcal{I}_o, \mathcal{T}_o | \mathcal{Z}, \mathcal{C}, \mathcal{B}, \delta, \boldsymbol{\eta}, \sigma_T^2)$ is the same as the log of Equation (26), without the latent receiver part that does not depend on either $\boldsymbol{\eta}$ or σ_T^2 . After taking the log, we obtain the log of acceptance ratio:

$$\begin{aligned}
& \log(\mathcal{N}(\boldsymbol{\eta}' | \mu_{\boldsymbol{\eta}}, \Sigma_{\boldsymbol{\eta}})) - \log(\mathcal{N}(\boldsymbol{\eta} | \mu_{\boldsymbol{\eta}}, \Sigma_{\boldsymbol{\eta}})) + \log(\mathcal{IG}(\sigma_T^{2'} | a_T, b_T)) - \log(\mathcal{IG}(\sigma_T^2 | a_T, b_T)) \\
&+ \log \left(\phi_l(\Delta T_{i_o^{(D)} J_o^{(D)}}^{(D)}; \boldsymbol{\eta}^{T'} \mathbf{y}_{i_o^{(D)} J_o^{(D)}}^{(D)}, \sigma_T^{2'}) \right) + \sum_{i \in \mathcal{A}_{\setminus i_o^{(D)}}} \log \left(1 - \Phi_l(\Delta T_{i_o^{(D)} J_o^{(D)}}^{(D)}; \boldsymbol{\eta}^{T'} \mathbf{y}_{i J_i}^{(D)}, \sigma_T^{2'}) \right) \quad (31) \\
&- \log \left(\phi_l(\Delta T_{i_o^{(D)} J_o^{(D)}}^{(D)}; \boldsymbol{\eta}^T \mathbf{y}_{i_o^{(D)} J_o^{(D)}}^{(D)}, \sigma_T^2) \right) - \sum_{i \in \mathcal{A}_{\setminus i_o^{(D)}}} \log \left(1 - \Phi_l(\Delta T_{i_o^{(D)} J_o^{(D)}}^{(D)}; \boldsymbol{\eta}^T \mathbf{y}_{i J_i}^{(D)}, \sigma_T^2) \right),
\end{aligned}$$

where \mathcal{IG} is the inverse-Gamma density. Then, if the log of a sample from Uniform(0,1) is less than the log-acceptance probability (31), we accept the proposal $\boldsymbol{\eta}'$ and $\sigma_T^{2'}$. Otherwise, we reject.

C Details on Getting It Right Test

C.1 Backward Generating Process

For backward sampling, we let NKV be a $V \times K$ dimensional matrix where each entry will record the count of the number of tokens of word-type v that are currently assigned to topic k . Also let NK be a K dimensional vector recording the total count of tokens currently assigned to topic k . Word-assignments are implemented via collapsed Gibbs sampling (Griffiths, 2002), while the generation of tie data directly follows the generating process in Section 2.4. This “backward” version of the generative process is detailed below in Algorithm below.

Algorithm 4 Generate data with backward sampling

Input:

- 1) token topic assignments $\{\{z_n^{(d)}\}_{n=1}^{N^{(d)}}\}_{d=1}^D$,
- 2) topic interaction pattern assignments, $\{C_k\}_{k=1}^K$,
- 3) interaction pattern parameters $\{\boldsymbol{b}^{(c)}\}_{c=1}^C$,
- 4) receiver size parameter δ ,
- 5) time-increment parameters $(\boldsymbol{\eta}, \sigma_T^2)$.

Set $NKV = 0$ and $NK = 0$

for $d=1$ to D **do**

```

for  $n=1$  to  $N^{(d)}$  do
  for  $v=1$  to  $V$  do
    token-word-type-distribution $_n^{(d)}[v] = \frac{NKV_{v,z_n^{(d)}} + \beta \mathbf{u}_v}{NK_{z_n^{(d)}} + \beta}$ 
  end
  draw  $w_n^{(d)} \sim (\text{token-word-type-distribution}_n^{(d)})$ 
   $NKV_{w_n^{(d)}, z_n^{(d)}} += 1$ 
   $NK_{z_n^{(d)}} += 1$ 
end
for  $i=1$  to  $A$  do
  for  $j=1$  to  $A$  do
    if  $j \neq i$  then
      calculate  $\boldsymbol{x}_{t_{+}^{(d-1)}}^{(c)}(i, j)$ , the network statistics evaluated at time  $t_{+}^{(d-1)}$ 
      set  $\lambda_{ij}^{(d)} = \sum_{c=1}^C p_c^{(d)} \cdot \exp\left\{\boldsymbol{b}^{(c)T} \boldsymbol{x}_{t_{+}^{(d-1)}}^{(c)}(i, j)\right\} \cdot 1\{j \in \mathcal{A} \setminus i\}$ 
    end
  end
  draw  $J_i^{(d)} \sim \text{Gibbs measure}(\{\lambda_{ij}^{(d)}\}_{j=1}^A, \delta)$ 
  calculate  $\boldsymbol{y}_{iJ_i}^{(d)}$  (including  $\boldsymbol{x}_{iJ_i}^{(d)}$ )
  draw  $\Delta T_{iJ_i}^{(d)} \sim \text{Lognormal}(\boldsymbol{\eta}^T \boldsymbol{y}_{iJ_i}^{(d)}, \sigma_T^2)$ 
end
  set  $i_o^{(d)} = i_{\min(\Delta T_{iJ_i}^{(d)})}$ ,  $J_o^{(d)} = J_{i^{(d)}}$ , and  $t_o^{(d)} = t_o^{(d-1)} + \min(\Delta T_{iJ_i}^{(d)})$ 
end

```

C.2 Initialization of History $\boldsymbol{x}_t^{(c)}$

Considering that our network statistics $\boldsymbol{x}_t^{(c)}$ are generated as a function of the network history, it is necessary to use the same initial value of $\boldsymbol{x}_t^{(c)}$ across the forward and backward samples. If not, when we generate fixed number of documents, we cannot guarantee the same number of documents

used for the inference, since only the documents with its timestamp greater than 384 hours are used in the inference. In the extreme cases, we may end up with two types of failure:

1. Zero document generated after 384 hours (i.e. $t^{(10)} < 384$), making no documents to be used for inference,
2. Zero document generated before 384 hours (i.e. $t^{(1)} > 384$), making the estimate of \mathcal{B} totally biased since $\forall \mathbf{x}_t^{(c)}(i, j) = 0$.

Therefore, we fix the initial state of $\mathbf{x}_t^{(c)}$ over the entire GiR process. Specifically, we fix some baseline documents where the timestamps are all smaller than 384 and use as an input for forward sampling, backward sampling, and the inference. Then, in the forward and backward generative process, we set the starting point of the timestamp as $t^{(0)} = 384$ and generate fixed number of documents given the initial $\mathbf{x}_{t^{(0)}=384}^{(c)}$ so that we can achieve consistency in the generated number of documents with $t^{(d)} > 384$.

C.3 GiR Implementation Details

While we tried a number of different parameter combinations in the course of testing, we outline our standard setup. We selected the following parameter values:

- D (number of documents) = 5
- $N^{(d)}$ (tokens per document) = 4
- A (number of actors) = 4
- W (unique word types) = 5
- C (number of interaction patterns) = 2
- K (number of topics) = 4
- α (Dirichlet concentration prior) = 2
- \mathbf{m} (Dirichlet base prior) = \mathbf{u}
- β (Dirichlet concentration prior) = 2
- \mathbf{n} (Dirichlet base prior) = \mathbf{u}
- netstat = “intercept” and “dyadic”
- prior for $\mathbf{b}^{(c)}$: $\mu_{\mathbf{b}^{(c)}} = (-3, \mathbf{0}_6)$, $\Sigma_{\mathbf{b}^{(c)}} = 0.05 \times I_7$
- prior for δ : $\mu_\delta = 2.5$, $\sigma_\delta^2 = 0.0001$
- prior for η : $\mu_\eta = 2.5$, $\sigma_\eta^2 = 0.0001$
- I (outer iteration) = 5
- n_1 (hyperparameter optimization) = 0
- n_2 (M-H sampling iteration of \mathcal{B}) = 5500
- burn (M-H sampling burn-in of \mathcal{B}) = 500
- thin (M-H sampling thinning of \mathcal{B}) = 5
- σ_{Q1}^2 (proposal variance for \mathcal{B}) = 0.1
- n_3 (M-H sampling iteration of δ) = 500
- σ_{Q2}^2 (proposal variance for δ) = 0.0002

D Details on Document Prediction Experiments

D.1 Prediction using IPTM

Here, one important point is that we do not have $\mathcal{Z}^{(D)}$ and $\mathcal{J}_a^{(D)}$ that came directly from inference on the observed data, as they represent document-specific (i.e., ‘local’) latent variables. We need to somehow initialize them conditional on the inferred $\mathcal{Z}^{1:(D-1)}$ and $\mathcal{J}_a^{1:(D-1)}$, and update them multiple times such that they are independent of the initialized ones. This can be done via MCMC by iteratively drawing the new data, then the new latent recipients conditional upon the new data and the new token topic assignments conditional upon the new data and the new latent recipients. We present detailed pseudocode in Algorithm 9 below.

Algorithm 5 Predicting tie data for document D

Inputs:

1. D , the document to make a prediction
2. O , number of outer iterations of inference from which to generate predictions
3. R , the number of iterations to sample predicted data within each outer iteration
4. $\mathbf{w}^{(D)}$, the observed words in the document to predict (D)
5. b , the number of burnin iterations of sampling predicted data within outer iteration
6. q , the thinning interval at which to keep predicted data within outer iteration

Run burn-in iterations for inference on documents 1 through $(D - 1)$

for $o=1$ to O **do**

run an outer iteration of inference on documents 1 through $(D - 1)$
obtain global latent variables $(\mathcal{B}, \mathcal{C}, \delta, \eta)$ and local variables $(\mathcal{Z}^{(1:D-1)}, \mathcal{J}_a^{(1:D-1)})$
initialize $N_{v|k}$ and N_k counts from the inference on $\mathcal{Z}^{(1:D-1)}$ and set $N_{k|D} = 0$

for $n = 1$ to $N^{(D)}$ **do**

for $k = 1$ to K **do**

set $P(z_n^{(D)} = k | z_{1:n-1}^{(D)}, w_{1:n}^{(D)}, \mathcal{Z}^{(1:D-1)}, \mathbf{w}^{(1:D-1)}) = \frac{(N_{k|D} + \alpha m_k)}{(n-1+\alpha)} \times \frac{(N_{w_n^{(D)}|k} + \beta)}{(N_k + \beta)}$

end

draw $z_n^{(D)} \sim P(z_n^{(D)} = k | z_{1:n-1}^{(D)}, w_{1:n}^{(D)}, \mathcal{Z}^{(1:D-1)}, \mathbf{w}^{(1:D-1)})$

increment $N_{w_n^{(D)}|z_n^{(D)}}$ and $N_{z_n^{(D)}}$ and $N_{z_n^{(D)}|D}$

end

initialize $\mathcal{J}_a^{(D)}$ from the inference on \mathcal{J}_a for documents 1 through $(D - 1)$ as below:

for $i = 1$ to A **do**

sample one document $d^* \sim \text{Discrete-uniform}(1, D - 1)$

set $\mathcal{J}_{a,i}^{(D)} = \mathcal{J}_{a,i}^{(d^*)}$

(NOTE: can't directly sample from the generative process due to computational complexity)

end

for $r=1$ to $b+R$ **do**

sample $i_o^{(D)}$, $J_o^{(D)}$, and $t_o^{(D)}$ following generative process

for $n = 1$ to $N^{(D)}$ **do**

Sample $z_n^{(D)}$ using Equation (27):

$$P(z_n^{(D)} = k | \mathcal{Z}_{\setminus D,n}, \mathcal{C}, \mathcal{B}, \mathbf{w}^{(D)}, \mathcal{J}_a, \mathcal{I}_O^{(D)}, \mathcal{J}_O^{(D)}, \mathcal{T}_O^{(D)}, \mathcal{Z}_{1:(D-1)}, \mathbf{w}_{1:(D-1)}, \beta, \mathbf{u}, \alpha, \mathbf{m}).$$

end

for $i \neq i_o^{*(d)}$ **do**

for $j \neq i$ **do**

Sample $\mathcal{J}_{ij}^{(d)}$ using Equations (19) and (20):

$$P(\mathcal{J}_{ij}^{(D)} = J_{ij}^{(D)} | \mathcal{J}_{i \setminus j}^{(D)}, \mathcal{Z}, \mathcal{C}, \mathcal{B}, \delta, \eta, \mathbf{w}^{(D)}, \mathcal{J}_{a,-i}^{(D)}, \mathcal{I}_O^{(D)}, \mathcal{J}_O^{(D)}, \mathcal{T}_O^{(D)}).$$

end

end

if $r > b$ & $r = qn$ (where n is an integer),

store $(i_o^{(D)}, J_o^{(D)}, t_o^{(D)})$

end

end

D.2 Prediction using Regressions

Algorithm 6 Predicting tie data for document D

Input

1. $\{i_o^{(d)}\}_{d=1}^{D-1}$, Observed senders for documents 1 through $D - 1$
2. $\{J_o^{(d)}\}_{d=1}^{D-1}$, Observed recipients for documents 1 through $D - 1$
3. $\{t_o^{(d)}\}_{d=1}^{D-1}$, Observed timing for documents 1 through $D - 1$
4. R , Number of predictions to make regarding document D

Train models defined by Equations (9) and (10) by Bayesian estimation using MCMC given the input data.

for $r = 1$ to R **do**

 Draw b_0 , \mathbf{b} , and η from their sampling distributions conditional on the input data.

for $i=1$ to A **do**

 set $J_i^{(D)}$ to a vector of zeros

while $\sum_{j \neq i} J_{ij}^{(D)} = 0$ **do**

for $j \neq i$ **do**

 | Draw $J_{ij}^{(D)}$ using Equation (9)

 | **end**

 | **end**

 | Draw $t_i^{(D)}$ using Equation (10)

 | **end**

 Set $s = \min_{i \in 1, 2, \dots, A} t_i^{(D)}$

 Store $i_o^{(D)} = s$

 Store $t_o^{(D)} = t_s^{*(D)}$

 Store $J_o^{(D)} = J_s^{*(D)}$

end

E Details on Posterior Predictive Checks

Algorithm 7 Generate new data from posterior predictive distribution

Input data :

- 1) \mathcal{Z} , estimates of token topic assignments from document 1 through D
- 2) \mathcal{C} , estimates of topic interaction pattern assignments for topic $k = 1, \dots, K$
- 3) \mathcal{B} , estimates of interaction pattern parameters for intercation pattern $c = 1, \dots, C$
- 4) δ , receiver size parameter
- 5) η , time parameter

Sample $(\mathcal{I}_o^*, \mathcal{J}_o^*, \mathcal{T}_o^*, \mathcal{W}^*)$ from $P(\mathcal{I}_o^*, \mathcal{J}_o^*, \mathcal{T}_o^*, \mathcal{W}^* | \mathcal{Z}, \mathcal{C}, \mathcal{B}, \delta, \eta, \mathbf{u}, \alpha, \mathbf{m})$ using Algorithm 8.

One difference from Algorithm 8 is that we initialize \mathcal{W}^* as below:

Initialize NKV and NK counts from the inference on \mathcal{Z} and observed \mathcal{W} (instead of zeros)

Return $(\mathcal{I}_o^*, \mathcal{J}_o^*, \mathcal{T}_o^*, \mathcal{W}^*)$ as a draw from the posterior distribution of the data

F Mutual Information in Posterior Predictive Checks

G Plots for further updates

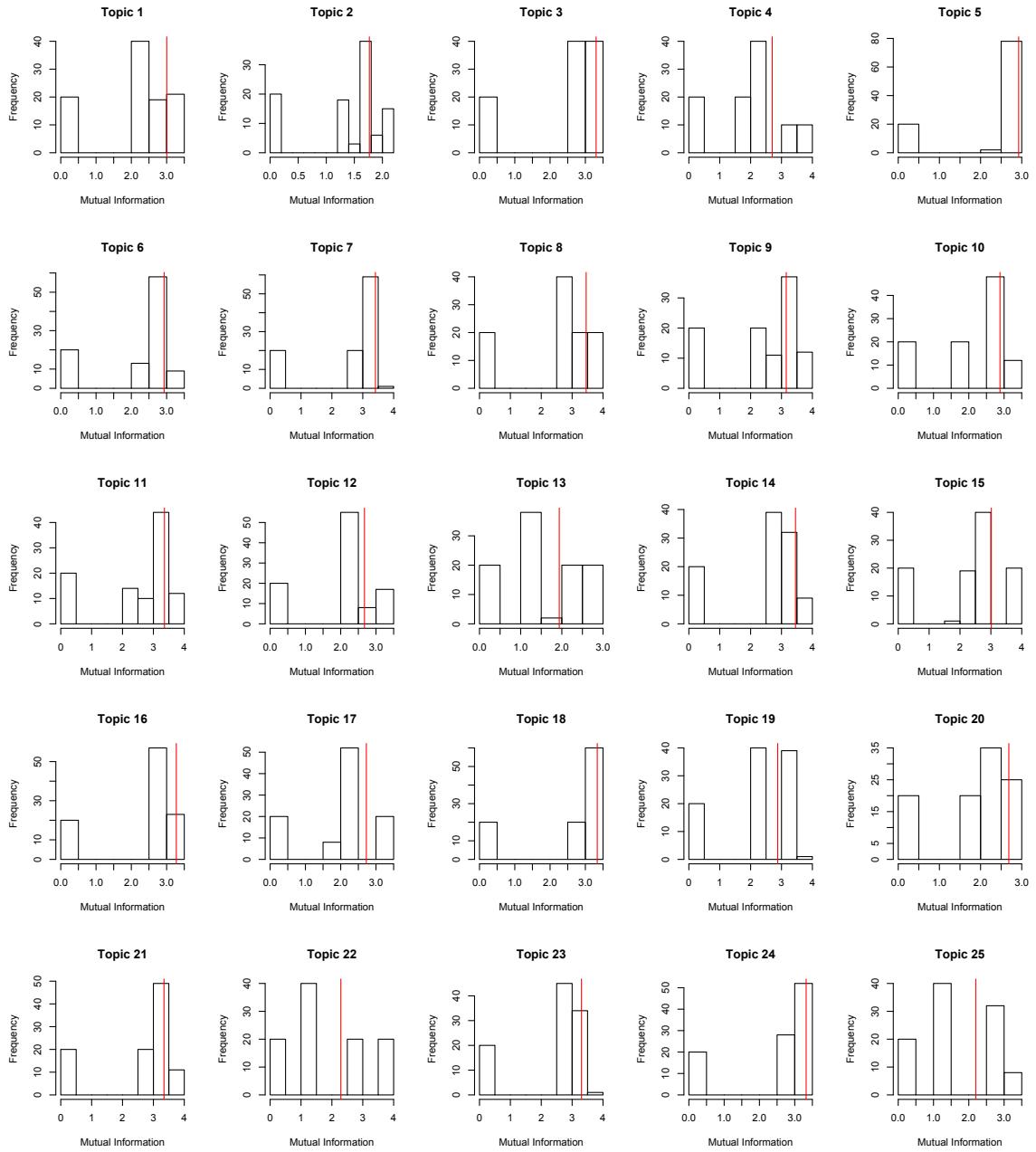
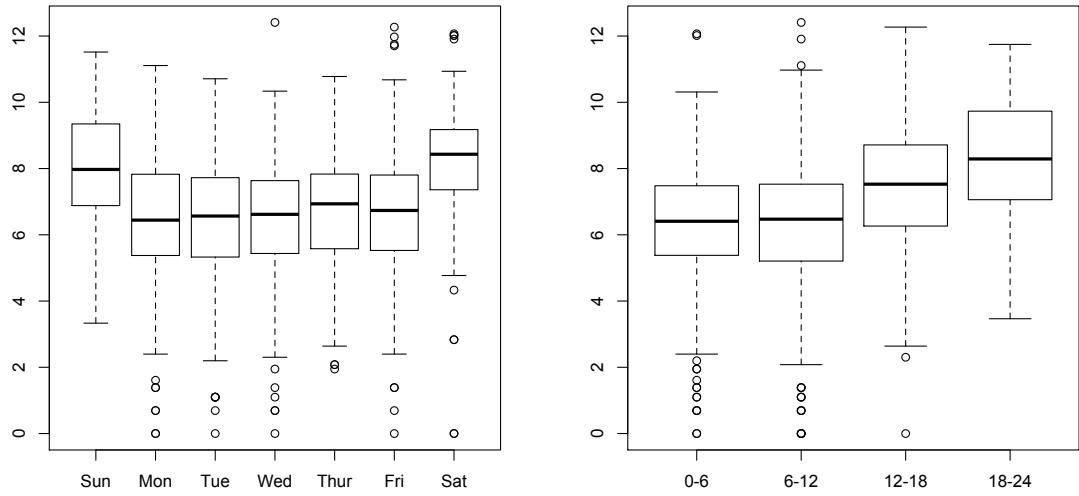
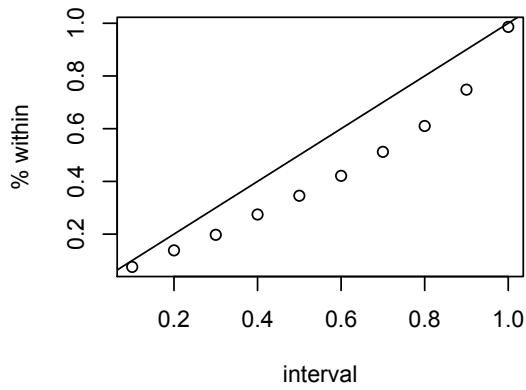


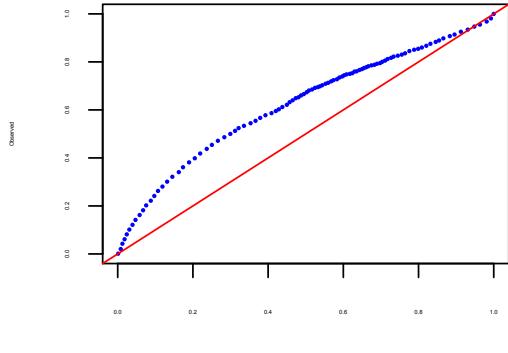
Figure 10: Mutual information (MI) score in posterior predictive check for Dare County (using $C = 2$ and $K = 25$)



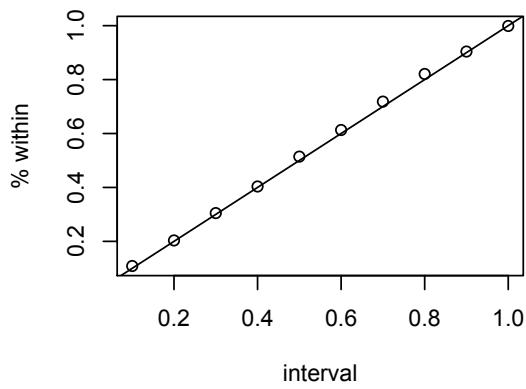
Exponential



sender+X+W+D (Exponential)



Lognormal



sender+X+W+D (Lognormal)

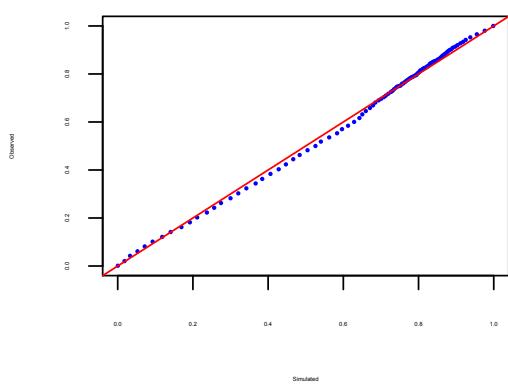


Figure 11: Boxplots for time-increments based on the week of the day (left) and time (right) of the document being sent