### 

#### 

#### 

# 

# 028 029

#### 

# 

#### 

#### 

# Suplementary Materials for "A Network Model for Dynamic Textual Communications with Application to Government Email Corpora"

#### Anonymous Authors<sup>1</sup>

#### 1. Normalizing constant of Gibbs measure

The non-empty Gibbs measure defines the probability of author i selecting the binary recipient vector  $u_{id}$  as

$$P(\boldsymbol{u}_{id}|\delta, \boldsymbol{\lambda}_{id})$$

$$= \frac{\exp\left\{\log(\mathbf{I}(\|\boldsymbol{u}_{id}\|_{1} > 0)) + \sum_{j \neq i} (\delta + \lambda_{idj}) u_{idj}\right\}}{Z(\delta, \boldsymbol{\lambda}_{id})}.$$

To use this distribution efficiently, we derive a closed-form expression for  $Z(\delta, \lambda_{id})$  that does not require brute-force summation over the support of  $u_{id}$  (i.e.  $\forall u_{id} \in [0,1]^A$ ). We recognize that if  $u_{id}$  were drawn via independent Bernoulli distributions in which  $P(u_{idj} = 1 | \delta, \lambda_{id})$  was given by  $\log \operatorname{it}(\delta + \lambda_{idj})$ , then

$$P(\boldsymbol{u}_{id}|\delta,\boldsymbol{\lambda}_{id}) \propto \exp\Big\{\sum_{j\neq i}(\delta+\lambda_{idj})u_{idj}\Big\}.$$

This is straightforward to verify by looking at

$$P(u_{idj} = 1 | \boldsymbol{u}_{id[-j]}, \delta, \boldsymbol{\lambda}_{id}) = \frac{\exp(\delta + \lambda_{idj})}{\exp(\delta + \lambda_{idj}) + 1}.$$

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

$$Z^{l}(\delta, \boldsymbol{\lambda}_{id}) = \sum_{\boldsymbol{u}_{id} \in [0,1]^{A}} \exp \Big\{ \sum_{j \neq i} (\delta + \lambda_{idj}) u_{idj} \Big\}.$$

Now, since

$$\exp\left\{\log\left(\mathbf{I}(\|\boldsymbol{u}_{id}\|_{1}>0)\right) + \sum_{j\neq i}(\delta + \lambda_{idj})u_{idj}\right\}$$
$$= \exp\left\{\sum_{j\neq i}(\delta + \lambda_{idj})u_{idj}\right\},$$

except when  $\|\boldsymbol{u}_{id}\|_1 = 0$ , we note that

$$Z(\delta, \lambda_{id}) = Z^{l}(\delta, \lambda_{id}) - \exp\left\{ \sum_{\forall u_{idj} = 0} (\delta + \lambda_{idj}) u_{idj} \right\}$$
$$= Z^{l}(\delta, \lambda_{id}) - 1.$$

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute. We can therefore derive a closed form expression for  $Z(\delta, \lambda_{id})$  via a closed form expression for  $Z^l(\delta, \lambda_{id})$ . This can be done by looking at the probability of the zero vector under the logistic-Bernoulli model:

$$\frac{\exp\left\{\sum_{\forall u_{idj}=0} (\delta + \lambda_{idj}) u_{idj}\right\}}{Z^{l}(\delta, \lambda_{id})} = \prod_{j \neq i} \left(1 - \frac{\exp\left(\delta + \lambda_{idj}\right)}{\exp\left(\delta + \lambda_{idj}\right) + 1}\right).$$

Then, we have

$$\frac{1}{Z^{l}(\delta, \boldsymbol{\lambda}_{id})} = \prod_{i \neq i} \frac{1}{\exp(\delta + \lambda_{idj}) + 1}.$$

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

$$Z(\delta, \lambda_{id}) = \prod_{j \neq i} (\exp\{\delta + \lambda_{idj}\} + 1) - 1.$$

## 2. Specification of Network Features

We provide the details on the specification of  $\boldsymbol{x}_{idjc} = (x_{idjc}^1, x_{idjc}^2, x_{idjc}^3)$  in Section 4.2. We first partitioned the interval  $[t_d - 16d, t_d)$  into L = 3 sub-intervals with equal length in the log-scale, by setting the difference  $\Delta_l = (6 \text{ hours}) \times 4^l$  for l = 1, 2, 3. In other words, we define the intervals  $I_d^l$  by

$$\begin{aligned} &[t_d - 384h, t_d) \\ &= [t_d - 384h, t_d - 96h) \cup [t_d - 96h, t_d - 24h) \cup [t_d - 24h, t_d) \\ &= I_d^3 \cup I_d^2 \cup I_d^1, \end{aligned}$$

where  $I_d^l$  is the half-open interval  $[t_d - \Delta_l, t_d - \Delta_{l-1})$ .

Then, for each time interval l=1,2,3, the degree and dyadic statistics are defined as:

1. outdegree 
$$_{id \cdot c}^l = \sum\limits_{d \in I^l} \pi_{dc} I\{i \rightarrow \forall j\}$$

2. indegree 
$$_{.djc}^{l}=\sum\limits_{d\in I_{s}^{l}}\pi_{dc}I\{\forall i
ightarrow j\}$$

3. 
$$\mathbf{send}_{idjc}^{l} = \sum_{d \in I_{d}^{l}} \pi_{dc} I\{i \to j\}$$

<sup>&</sup>lt;sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

4. 
$$\mathbf{receive}_{idjc}^l = \sum\limits_{d \in I_d^l} \pi_{dc} I\{j \to i\}$$

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 earlier works (Perry & 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  $3\times 3$  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. **2-send**<sup>$$l$$</sup> <sub>$idjc = 
$$\sum_{\substack{\max(l_1, l_2) = l \ h \neq i \\ h \neq j}} (\sum_{d \in I_d^{l_1}} \pi_{dc} I\{i \to h\}) (\sum_{d \in I_d^{l_2}} \pi_{dc} I\{h \to j\})$$$</sub> 

6. **2-receive**<sup>$$l$$</sup> <sub>$idjc = 
$$\sum_{\substack{\max(l_1,l_2)=l}} \sum_{\substack{h\neq i\\h\neq j}} (\sum_{d\in I_d^{l_1}} \pi_{dc} I\{h\rightarrow i\}) (\sum_{d\in I_d^{l_2}} \pi_{dc} I\{j\rightarrow h\})$$$</sub> 

7. 
$$\begin{aligned} \text{sibling}_{idjc}^{l} &= \\ &\sum_{\max(l_1,l_2)=l} \sum_{\substack{h\neq i\\h\neq j}} (\sum_{d\in I_d^{l_1}} \pi_{dc} I\{h\rightarrow i\}) (\sum_{d\in I_d^{l_2}} \pi_{dc} I\{h\rightarrow j\}) \end{aligned}$$

8. 
$$\begin{aligned} \text{cosibling}_{idjc}^l &= \\ &\sum_{\max(l_1,l_2)=l} \sum_{\substack{h\neq i\\h\neq j}} (\sum_{d\in I_d^{l_1}} \pi_{dc} I\{i\rightarrow h\}) (\sum_{d\in I_d^{l_2}} \pi_{dc} I\{j\rightarrow h\}) \end{aligned}$$

where  $l_1 = 1, 2, 3$  and  $l_2 = 1, 2, 3$ .

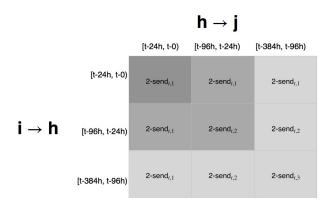


Figure 1. 2-send defined for each interval l=1,2,3. Cells with same shades sum up to one statistic.