

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

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

¹ The Pennsylvania State University

² Microsoft Research NYC

³ University of Massachusetts Amherst

June 8, 2017

Interaction-Partitioned Topic Model (IPTM)

- ▶ Probabilistic model for time-stamped textual communications (e.g. emails, cosponsorship of bills, international sanctions)
- ▶ Integration of two generative models:
 - Latent Dirichlet allocation (LDA) for topic-based contents
 - Dynamic exponential random graph model (ERGM) for ties
- ▶ IPTM assigns each topic to an “interaction pattern,” which is governed by a set of dynamic network features

“who communicates with whom about what, and when?”

Content Generating Process: LDA (Blei et al., 2003)

- ▶ For each topic $k = 1, \dots, K$:
 1. Topic-word distribution $\phi^{(k)} \sim \text{Dirichlet}(\beta, \mathbf{u})$
 - A topic k is characterized by a discrete distribution over V word types with probability vector $\phi^{(k)}$.
 2. Topic-IP distribution $c_k \sim \text{Uniform}(1, C)$
 - Each topic is associated with a single interaction pattern.

- ▶ For each document $d = 1, \dots, D$:
 - 3-1. Document-topic distribution $\theta^{(d)} \sim \text{Dirichlet}(\alpha, \mathbf{m})$
 - A document d is characterized by a discrete distribution over K topics with probability vector $\theta^{(d)}$.
 - 3-2. For each word in a document $n = 1$ to $N^{(d)}$:
 - (a) Choose a topic $z_n^{(d)} \sim \text{Multinomial}(\theta^{(d)})$
 - (b) Choose a word $w_n^{(d)} \sim \text{Multinomial}(\phi^{(z_n^{(d)})})$

Dynamic Network Features (Perry and Wolfe, 2012)

send

$i \longrightarrow j$

2-send

$i \longrightarrow h \longrightarrow j$

receive

$i \longleftarrow j$

2-receive

$i \longleftarrow h \longleftarrow j$

sibling



cosibling

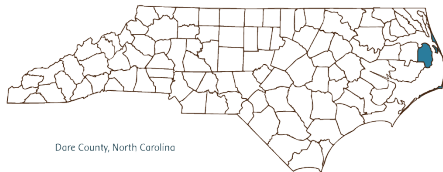


Tie Generating Process

Inference - Pseudocode

Data: North Carolina Dare county email data

- ▶ $D = 1456$ emails between $A = 27$ county government managers, covering 2 month periods (October 1 - November 30) in 2013



Dare County, North Carolina

Effect of Hurricane Sandy

IPTM Result

Conclusion

- ▶ Joint modeling of ties (sender, receiver, time) and contents
- ▶ Allowance of multicast – multiple senders and/or receivers
- ▶ Possible application to