

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

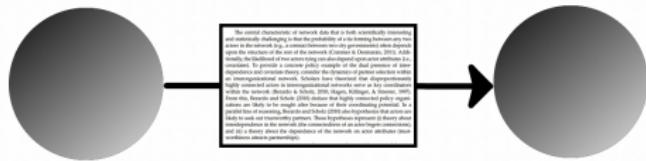
November 6, 2017

Work supported by NSF grants SES-1558661, SES-1619644, SES-1637089, and CISE-1320219)



Motivation

- ▶ Many networks have ties that are attributed with text:
 - ▶ International treaties, legislative cosponsorship, online discussions



- ▶ Network models can't account for text
- ▶ Models for text either can't account for ties or only account for simplistic network structure

Interaction-Partitioned Topic Model (IPTM)

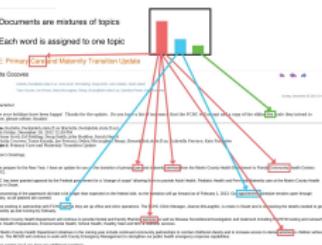
- ▶ Model for networks with timestamped, text-valued ties
- ▶ Draws on ideas from two existing models:
 - ▶ Latent Dirichlet allocation (LDA) for topic-based document content
 - ▶ Dynamic exponential random graph model (ERGM) for ties
- ▶ *Who communicates with whom about what, and when?*

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

- ▶ For each topic $k = 1, \dots, K$:
 - ▶ Draw a distribution over vocabulary:
 $\phi_k \sim \text{Dir}(\beta, (\frac{1}{V}, \dots, \frac{1}{V}))$

- ▶ For each email $d = 1, \dots, D$:
 - ▶ Draw a distribution over topics:
 $\theta_d \sim \text{Dir}(\alpha, (m_1, \dots, m_K))$
 - ▶ For each token $n = 1$ to N_d :
 - ▶ Draw a topic: $z_{dn} \sim \theta_d$
 - ▶ Draw a word type: $w_{dn} \sim \phi_{z_{dn}}$

	$k = 1$	$k = 2$	$k = 3$
support			
position			
fill			
staff			
desk			
service			
customer			
begin			
duties			
vacancy			
⋮			
IP = 1			
services			
care			
child			
information			
system			
community			
nurse			
completed			
provided			
pregnancy			
⋮			
IP = 2			
budget			
funds			
money			
budgeted			
including			
cost			
salary			
amount			
revenues			
debt			
⋮			
IP = 1			



Key Idea: Interaction Patterns

- ▶ Different topics associated with different interaction patterns
- ▶ For each topic $k = 1, \dots, K$:
 - ▶ Assign topic k to an interaction pattern: $l_k \sim \text{Unif}(1, C)$
- ▶ For each email $d = 1, \dots, D$:
 - ▶ Calculate distribution over interaction patterns:

$$\pi_{dc} = \frac{\sum_{k:l_k=c} N_{dk}}{N_d}$$

Generating Ties and Timestamps

- ▶ Continuous-time process
- ▶ Ties predicted using dynamic network features:
 - ▶ Popularity, activity
 - ▶ Repetition, reciprocity
 - ▶ Transitivity
- ▶ Author determines recipients and timestamp
- ▶ Innovative approach to modeling multiple recipients

Generating Hypothetical Authors and Recipients

1. For each author $i \in \{1, \dots, A\}$ and recipient $j \in \{1, \dots, A\}$ ($i \neq j$), calculate the stochastic intensity between i and j :

$$\nu_{idj} = \exp(\mathbf{b}_c^\top \mathbf{x}_{idj})$$

$$\lambda_{idj} = \sum_{c=1}^C \pi_{dc} \nu_{idj}$$

which combines information about content and network structure

2. For each author $i \in \{1, \dots, A\}$, draw a binary vector \mathbf{u}_{id} of length A from a non-empty Gibbs measure (Fellows and Handcock, 2017)

$$\mathbf{u}_{id} = (u_{id1}, \dots, u_{idA}) \sim \text{Gibbs}(\delta, (\lambda_{id1}, \dots, \lambda_{idA}))$$

where δ is a real-valued value that controls the number of recipients

Generating Hypothetical Timestamps and Actual Data

3. For each author $i \in \{1, \dots, A\}$, generate a timestamp:

$$\mu_{id} = \sum_{c=1}^C \pi_{dc} \text{GeomMean}(\{\nu_{idjc}\}_{j:u_{idj}=1})$$

$$\tau_{id} \sim \text{Exp}(\eta \mu_{id})$$

$$t_{id} = t_{d-1} + \tau_{id}$$

where η is a positive real-valued value

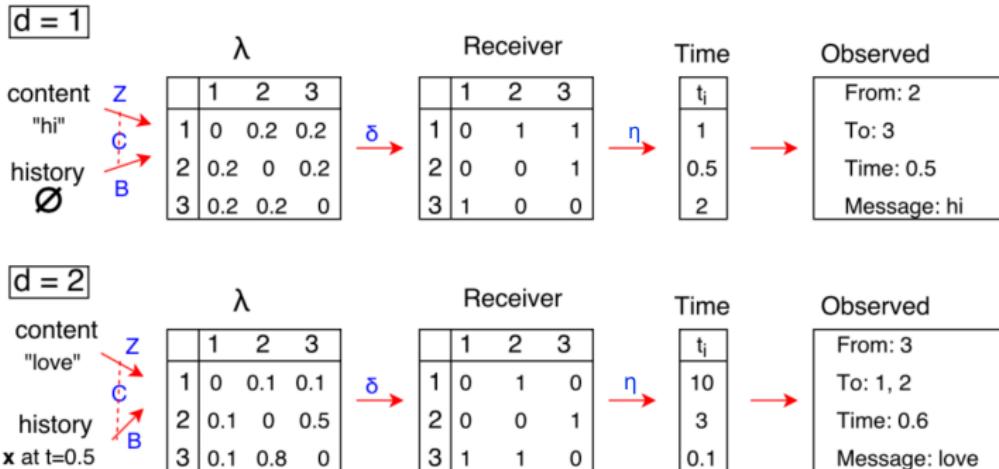
4. Select the actual author, recipients, and timestamp:

$$a_d = \operatorname{argmin}_i(t_{id})$$

$$\mathbf{y}_d = \mathbf{u}_{a_d d}$$

$$t_d = t_{a_d d}$$

Generating Ties and Timestamps



Network Features (Perry and Wolfe, 2012)

- ▶ Features capture different types of interaction: popularity, activity, repetition, reciprocity, transitivity

outdegree ($i \rightarrow \forall j$) **send** ($i \rightarrow j$)

indegree ($i \leftarrow \forall j$) **receive** ($i \leftarrow j$)

2-send $\sum_h (i \rightarrow h \rightarrow j)$ **sibling** $\sum_h (h \xleftrightarrow{i} j)$

2-receive $\sum_h (i \leftarrow h \leftarrow j)$ **cosibling** $\sum_h (h \leftarrow i \leftarrow j)$

Dynamic Network Features

- ▶ To form \mathbf{x}_{idjc} , consider 3 time intervals prior to t_{d-1}^+ :
 - ▶ 3–16 days, 1–3 days, 0–1 days
- ▶ For each interval, compute each network feature focusing on emails in that interval: 24 dynamic network features in \mathbf{x}_{idjc}

		$\mathbf{h} \rightarrow \mathbf{j}$		
		[t-24h, t-0)	[t-96h, t-24h)	[t-384h, t-96h)
[t-24h, t-0)		2-send _{i,j}	2-send _{i,j}	2-send _{i,j}
$\mathbf{i} \rightarrow \mathbf{h}$	[t-96h, t-24h)	2-send _{i,j}	2-send _{i,j}	2-send _{i,j}
	[t-384h, t-96h)	2-send _{i,j}	2-send _{i,j}	2-send _{i,j}

Inference

Algorithm 1 MCMC: Metropolis-Hastings within Gibbs

Initialize latent variables

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

 Sample token-topic assignments z_{dn}

 Sample topic-interaction pattern assignments l_k

 Sample hypothetical author-recipient set pairs u_{ad}

 Sample coefficients b_c using Metropolis-Hastings

 Sample recipient parameter δ using Metropolis-Hastings

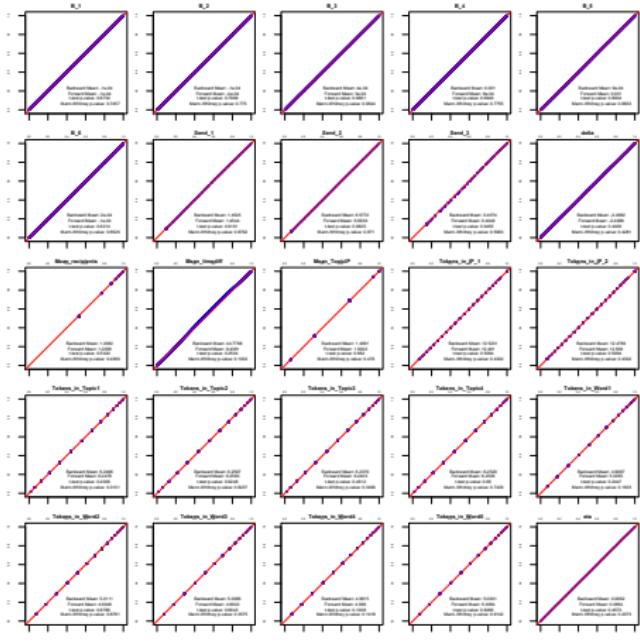
 Sample timestamp parameter η using Metropolis-Hastings

end

Getting it Right: Testing Math and Code (Geweke, 2004)

- ▶ *Generate forward samples:*
 1. Draw latent variables from priors
 2. Draw synthetic data conditioned on latent variables
 3. Repeat 1 and 2 many times
- ▶ *Backward samples:*
 1. Start with a forward sample
 2. Draw latent variables from posterior using inference code
 3. Draw synthetic data conditioned on latent variables
 4. Repeat 2 and 3 many times
- ▶ Forward samples and backward samples should match

Getting it Right Results



NC Email Corpus: Dare County



- ▶ $D = 2,210$ emails
- ▶ $A = 27$ department managers
- ▶ $V = 2,907$ words in vocabulary
- ▶ Covers 3 months (September 1 to November 30) in 2012:
 - ▶ Hurricane Sandy: October 26 to October 30

Example Email

- ▶ From: Health
- ▶ To: County Manager
- ▶ At: 8 April 2012 8:46:00
- ▶ Content: ... A few months ago Jennifer was on call and Friday night she was out till 2 am and started the calls before 5pm, there were 2 bites at one time and then a goat in the highway, she took the goat to her barn because we did not have anywhere to keep it...

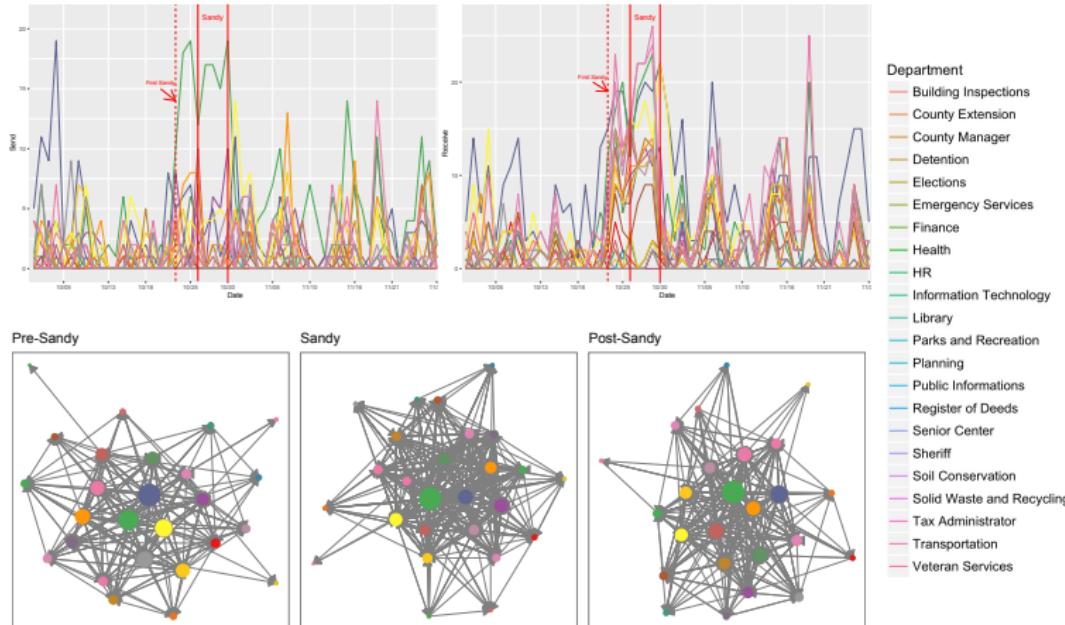


(Not the actual goat.)

Hypotheses

- ▶ Personal/social topics exhibit reciprocity and transitivity
- ▶ Dissemination topics don't exhibit reciprocity
- ▶ Topics about Hurricane Sandy look different

Data Exploration



Interaction Pattern 1: Top 5 Topics

23	18	11	20	8
change	will	sandy	time	services
order	winds	munis	hours	public
manager	location	hurricane	monday	white
storm	beach	position	leave	director
emergency	hydrant	monday	employees	fyi
coastal	water	point	timesheets	tim
statute	relocation	power	storm	update
evacuation	mirlo	update	employee	status
track	road	storm	tomorrow	board
couple	high	hey	work	approval
changes	moving	release	regular	wanted
well	gas	weeks	period	rec
concerns	forecast	weekend	comp	adult
things	saturday	working	sheets	older
consistent	project	month	vacation	today
boat	outer	problems	administrative	storm
misdemeanor	map	strong	operation	reminder
program	airport	impacts	timesheet	called
system	called	three	personnel	sep
powerpoint	banks	ncdot	will	charlotte

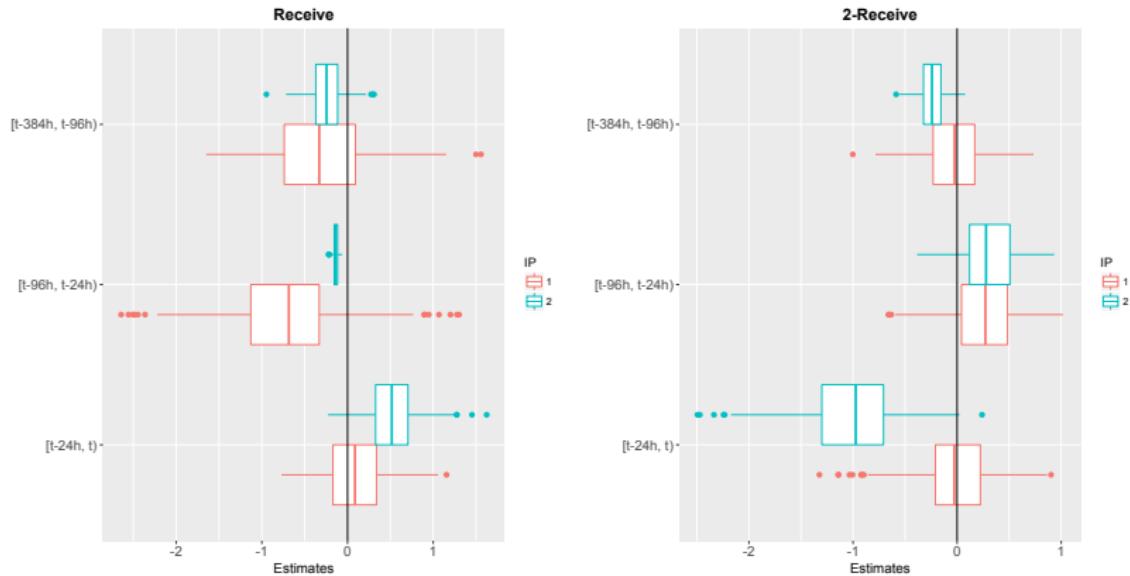
Example Email

- ▶ From: Public Information
- ▶ To: Emergency Services, Detention
- ▶ At: 25 Oct 2012 19:14:47
- ▶ Content: will will storm good sure sandy morning morning morning
morning send change plan hurricane wanted copy base afternoon
keep description description saturday saturday saturday saturday
touch touch release signature plans prior ready duration submitted
eoc running exactly activation activation released mind joint jis
pasted activate

Interaction Pattern 2: Top 5 Topics

21	12	14	4	10
library	marshall	will	board	survey
best	collins	center	meeting	request
start	drive	day	property	call
web	manteo	office	planning	sure
place	phone	great	will	people
visit	box	manager	notice	emergency
albermarle	fax	assistance	january	thought
east	resources	problem	review	seafood
regional	director	lot	commissioners	mail
system	human	call	list	grant
voice	phr	currently	thoughts	ncacc
librarian	policy	september	amendment	community
learning	time	early	members	check
discovery	october	parking	dot	rodanthe
manteo	hire	received	budget	district
expressed	will	monday	building	complete
views	offices	baum	inspections	option
box	approved	questions	night	best
conversion	hour	year	text	residents
director	told	today	copy	talk

Coefficients



Predicting Ties and Timestamps

Algorithm 2 Predict ties and timestamp for next email

Run burn-in iterations of MCMC on emails 1 through $d - 1$

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

 Run single iteration of MCMC on emails 1 through $d - 1$

 Initialize a_d , y_d , t_d , and z_d

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

 Sample a_d , y_d , t_d conditioned on z_d using generative process

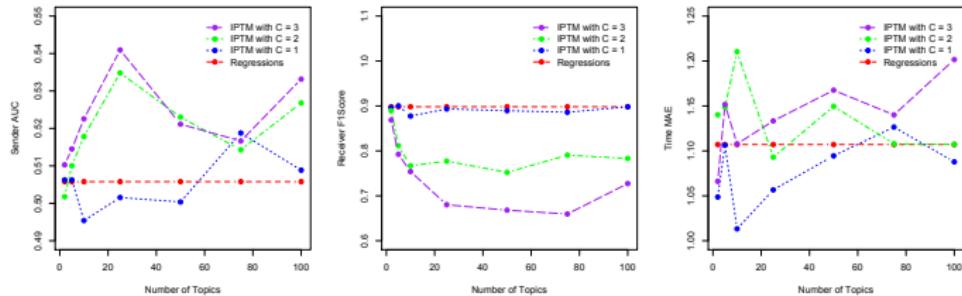
 Sample z_d from conditional posterior

end

 Save a_d , y_d , t_d , and z_d

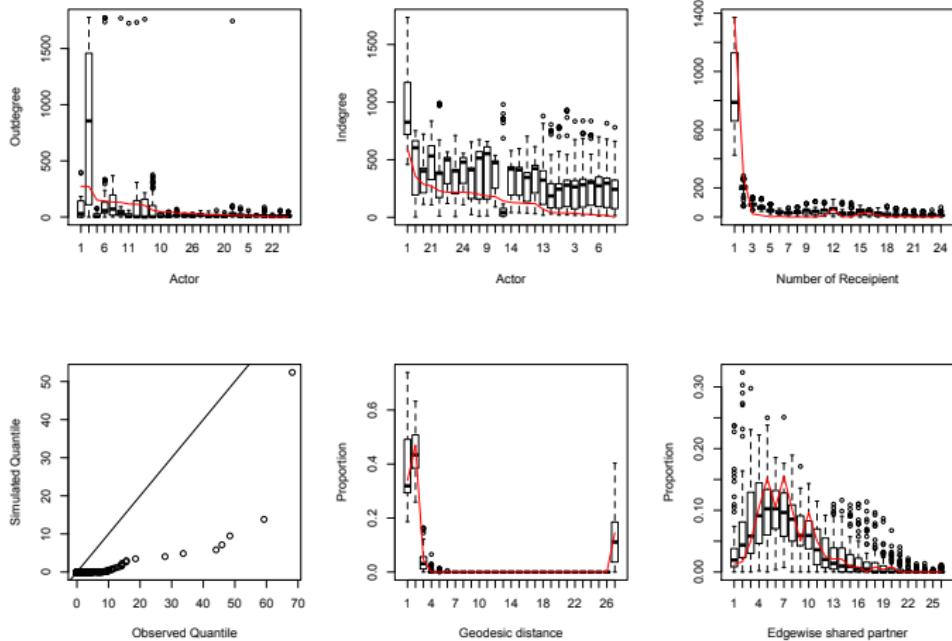
end

Example Results



- ▶ Model does poorly at predicting ties and timestamps

Posterior Predictive Checks



To Conclude...

- ▶ Model for networks with timestamped, text-valued ties
- ▶ Many potential applications in political science
- ▶ Development of R package ‘IPTM’

I otter know the answer...



I don't have the necessary koalaifications...



If they'd gibbon me more time...



If I'd implemented the IPTM with my bear hands...



I'd be lion if I said yes...



Let's table bat discussion...

