

# The Hyperedge Event Model

**Bomin Kim**

Department of Statistics  
Pennsylvania State University

July 29, 2018

Joint Statistical Meetings 2018

# Collaborators



- ▶ **Aaron Schein**  
College of Information and Computer Sciences, UMass Amherst
- ▶ **Bruce Desmarais**  
Department of Political Science, Pennsylvania State University
- ▶ **Hanna Wallach**  
Microsoft Research NYC

# Motivations: Hyperedge Event Model

- ▶ **Hyperedge**: directed edges from one sender to multiple receivers or from multiple senders to one receiver
- ▶ **Event**: timestamped events in the continuous-time scale
- ▶ **Model**: statistical framework to jointly understand

“who interacts with whom, and when?”

## Generative Process: “Who Interacts with Whom”

For event  $e = 1, \dots, E$ , between  $i \in \{1, \dots, A\}$  &  $j \in \{1, \dots, A\}$ ,

- Receiver intensity for every sender-receiver pair  $(i, j)_{i \neq j}$

$$\lambda_{iej} = \mathbf{b}^T \mathbf{x}_{iej},$$

where  $\mathbf{x}_{iej}$  is a set of receiver selection features or covariates and  $\mathbf{b}$  is the corresponding  $P$ -dimensional coefficient.

- Every sender  $i$  selects candidate receivers from non-empty multivariate Bernoulli distribution<sup>1</sup>  $\mathbf{u}_{ie} \sim \text{MB}_G(\lambda_{ie1}, \dots, \lambda_{ieA})$

$$P(\mathbf{u}_{ie} | \mathbf{b}, \mathbf{x}_{iej}) \propto \exp \left( \log(I(\|\mathbf{u}_{ie}\|_1 > 0)) + \sum_{j \neq i} \lambda_{iej} u_{iej} \right)$$

---

<sup>1</sup>Fellows and Handcock (2017); Dai et al. (2013)

# Generative Process: “and When”

- ▶ **Timing rate** for each sender  $i$

$$\mu_{ie} = g^{-1}(\boldsymbol{\eta}^T \mathbf{y}_{ie}),$$

where  $\mathbf{y}_{ie}$  is a set of timing features or covariates and  $\boldsymbol{\eta}$  is the corresponding  $Q$ -dimensional coefficient.

- ▶ **Generalized linear model** (GLM) for time increment  $\tau_{ie}$  so that

$$E(\tau_{ie}) = \mu_{ie} \text{ and } V(\tau_{ie}) = V(\mu_{ie}),$$

with a choice of distribution from exponential family.

- ▶ Select the sender-receiver-set with **the smallest time increment**<sup>2</sup>

$$s_e = \operatorname{argmin}_i(\tau_{ie}),$$

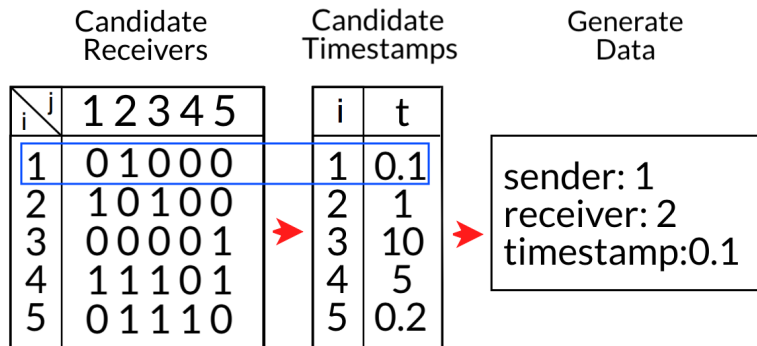
$$\mathbf{r}_e = \mathbf{u}_{s_e e},$$

$$t_e = t_{e-1} + \tau_{s_e e}.$$

---

<sup>2</sup>Snijders (1996)

# Generative Process: Sender, Receivers, and Timestamps<sup>3</sup>



- **Bayesian inference:** invert the generative process to obtain the posterior distribution over the latent variables—i.e.,  $(\mathbf{u}, \mathbf{b}, \eta)$ .

<sup>3</sup>Assuming  $t_{e-1} = 0$  for simplicity.

# Application: Montgomery County Government Email Data<sup>4</sup>

- ▶ Email corpora covering inboxes and outboxes of **Montgomery county government managers** in North Carolina
- ▶ Contains  $E = 680$  emails, sent and received by  $A = 18$  department managers over 3 months (March–May) in 2012.

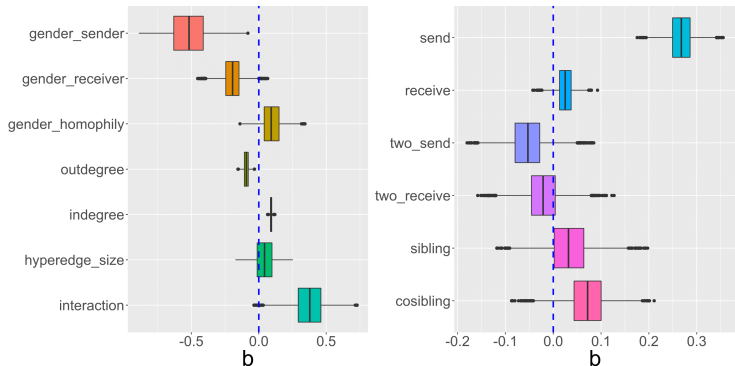
“To what extent are **nodal, dyadic or triadic network effects** relevant to predicting future emails?”

---

<sup>4</sup>ben Aaron et al. (2017)

## Results: Exploratory Analysis on $b$

$$\text{logit}(\lambda_{iej}) = \log\left(\frac{\lambda_{iej}}{1 - \lambda_{iej}}\right) = b_1 + b_2 x_{iej2} \dots + b_{14} x_{iej14},$$

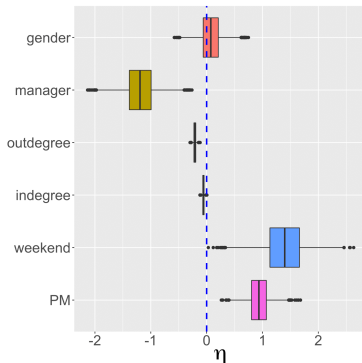


- ▶ Log odds is two times less if the sender is a woman.
- ▶ If  $i$  sent  $n$  number of emails to  $j$  last week, then  $i$  is  $e^{0.27n} \approx (1.32)^n$  times more likely to send an email to  $j$ .



## Results: Exploratory Analysis on $\eta$

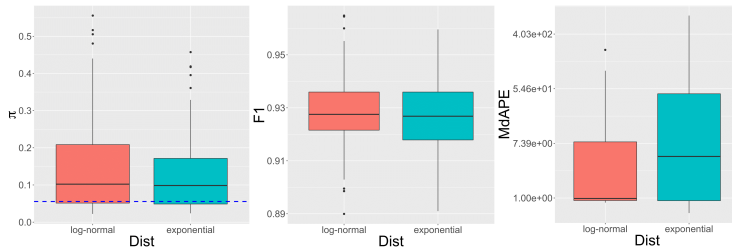
$$\log(\tau_{ie}) \sim N(\mu_{ie}, \sigma_{\tau}^2), \text{ with } \mu_{ie} = \eta_1 + \eta_2 y_{ie2} \dots + \eta_7 y_{ie7}.$$



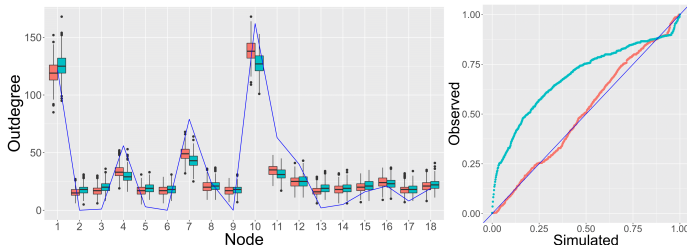
- ▶ If an email was sent during weekend or PM, then time to next email takes  $e^{1.55} \approx 4.72$  and  $e^{0.98} \approx 2.67$  hours longer.
- ▶ manager, outdegree, and indegree shorten time to next email.

# Comparison: Lognormal vs. Exponential

- Out-of-sample predictions: sender, receiver, and timestamp



- Posterior predictive checks (PPC)



# Conclusion and Discussion

- ▶ Account for **hyperedges** without duplications
- ▶ Flexible choice of **continuous-time distribution** via GLM
- ▶ Reverse the process for **multiple senders to one receiver** (e.g., international sanctions and co-sponsorship of bills)
- ▶ Sources: <http://arxiv.org/abs/1807.08225>  
<https://github.com/desmarais-lab/MulticastNetwork>