
Simultaneous text & Network modeling in MMSB for Learning Structured Interactions on Online Forums

Anonymous Author(s)

Affiliation

Address

email

1 Tentative Ideas

We discuss here the ideas and questions that are important for modeling structure in online forums. Following are the questions that we want answered:

1. When two persons interact in a thread or a post which topic/community they each belong to
2. When a user U is in community C what type of text does he use to communicate
3. Multi-user-Interaction: In a thread a user can post by addressing to a specific user but he is also talking to other users in the thread simultaneously. Can we model this phenomenon
4. There is an inherent bias towards the thread starter or in turn topic of the thread; can such an information be utilised in some form of a prior value/input
5. Multi-layer-Interaction: On the network side of things there are multiple signals which cannot be simply added to make a single signal e.g. different types of edges in the graphs (user calling by username and nick-name). Can the model take this into account. We are not doing this at present.
6. User posts aggregation; there are multiple ways to aggregate
 - (a) Network Layer aggregation: We call all types of edges as a single edge type and use this combined signal.
 - (b) aggregating user posts across multiple threads in the forum.
 - (c) Aggregating user post only in the same thread
 - (d) Aggregating user post only for same user-user pair interaction; i.e. a user might have posted multiple replies to another user and we aggregate all such replies into one for this user pair interaction.
 - (e) No aggregation at all.

2 Graphical Model & Generative Story

Based on the discussions above we came up with the following final model shown in figure 1. In this model, figure 1 below, we aggregate the posts of a given user in a given thread into one document called R_p .

The generative process for the figure is as follows:

Assuming that there are total N_t users in the thread t .

- For each Thread t
 - For each user $p \in \mathcal{N}_t$

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

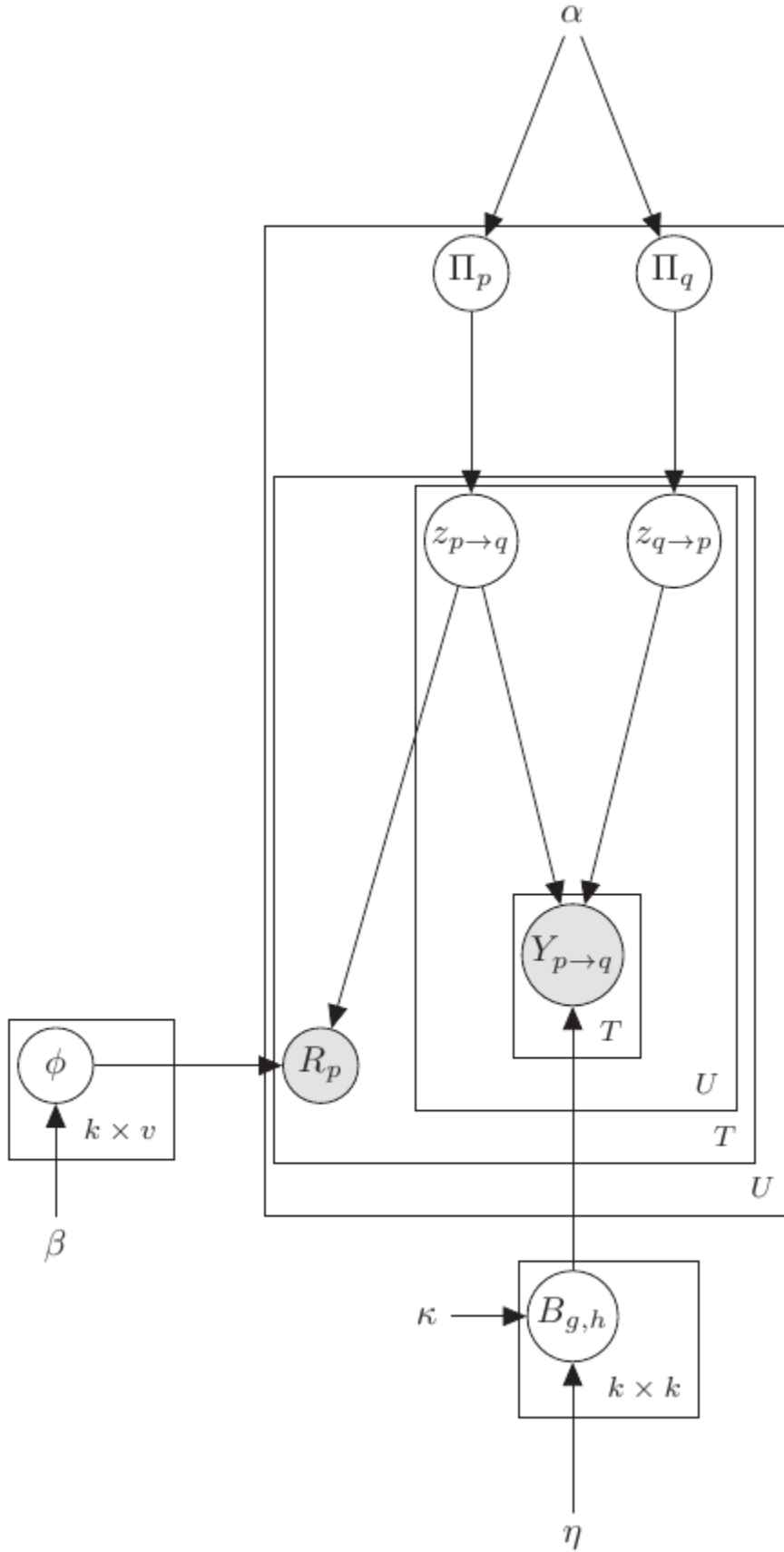


Figure 1: This graphical model takes into account multi-way interaction among users in a thread simultaneously

- * Draw a K dimensional mixed membership vector $\vec{\pi}_p \sim \text{Dirichlet}(\alpha)$
- * Draw $B(g, h) \sim \text{Gamma}(\kappa, \eta)$; where κ, η are parameters of the gamma distribution.
- For each pair of users $(p, q) \in \mathcal{N}_t \times \mathcal{N}_t$:
 - * Draw membership indicator for the indicator, $\vec{z}_{(p \rightarrow q, t)} \sim \text{Multinomial}(\pi_p)$.
 - * Draw membership indicator for the receiver, $\vec{z}_{(q \rightarrow p, t)} \sim \text{Multinomial}(\pi_q)$.
 - * Sample the value of their interaction, $Y(p, q, t) \sim \text{Poisson}(\vec{z}_{(p \rightarrow q, t)}^\top B \vec{z}_{(p \leftarrow q, t)})$.
- For each user $p \in \mathcal{N}_t$
 - * Draw ϕ_k from $\text{Dirichlet}(\beta)$.
 - * Form the set $Q_{p, t}$ that contains all the users that p interacts to on thread t
 - For each word $w \in R_{p, t}$
 - Draw $w \sim \phi(w|z_{(p \rightarrow q, t)}, \forall q \in Q_{p, t})$

The data likelihood for the model in figure 1

$$\begin{aligned}
 P(Y, R_p | \alpha, \beta, \kappa, \eta) &= \int_{\Phi} \int_{\Pi} \sum_z P(Y, R_p, z_{p \rightarrow q}, z_{p \leftarrow q}, \Phi, \Pi | \alpha, \beta, \kappa, \eta) \\
 &= \int_{\Phi} \int_{\Pi} \sum_z \left[\prod_{p, q} \prod_t P(Y_{pq}^t | z_{p \rightarrow q}^t, z_{p \leftarrow q}^t, B) \cdot P(z_{p \rightarrow q}^t | \Pi_p) \cdot P(z_{p \leftarrow q}^t | \Pi_q) \right. \\
 &\quad \cdot \left(\prod_p P(\Pi_p | \alpha) \prod_t \prod_p P(R_p^t | z_{p \rightarrow q}^t, \Phi) \cdot \prod_k P(\Phi_k | \beta) \right) \cdot \left. \prod_{g, h} P(B_{gh} | \eta, \kappa) \right] \quad (1)
 \end{aligned}$$

The complete log likelihood of the model is:

$$\begin{aligned}
 \log P(Y, W, z_{\rightarrow}, z_{\leftarrow}, \Phi, \Pi, B | \kappa, \eta, \beta, \alpha) &= \sum_t \sum_{p, q} \log P(Y_{pq}^t | z_{p \rightarrow q}^t, z_{p \leftarrow q}^t, B) + \\
 &\quad \sum_t \sum_{p, q} (\log P(z_{p \rightarrow q}^t | \Pi_p) + \log P(z_{p \leftarrow q}^t | \Pi_p)) + \sum_p \log P(\Pi_p | \alpha) + \\
 &\quad \sum_t \sum_p \sum_{w \in R_p^t} \log P(w | z_{p \rightarrow}, \Phi) + \sum_k \log P(\Phi_k | \beta) + \sum_{gh} \log P(B_{gh} | \eta, \kappa) \quad (2)
 \end{aligned}$$

The mean field variational approximation for the posterior is

$$\begin{aligned}
 q(z, \Phi, \Pi, B | \Delta_{z_{\rightarrow}}, \Delta_{\Phi}, \Delta_B, \Delta_{z_{\leftarrow}}, \Delta_{B_{\kappa}}) &= \prod_{t, p, q} \left(q_1(z_{p \rightarrow q}^t | \Delta_{z_{p \rightarrow q}}) + q_1(z_{p \leftarrow q}^t | \Delta_{z_{p \leftarrow q}}) \right) \\
 &\quad \cdot \prod_p q_4(\Pi_p | \Delta_{\Pi_p}) \prod_k q_3(\Phi_k | \Delta_{\Phi_k}) \prod_{g, h} q(B_{g, h} | \Delta_{B_{\eta}}, \Delta_{B_{\kappa}}) \quad (3)
 \end{aligned}$$

The lower bound for the data log-likelihood from jensen's inequality is:

$$L_{\Delta} = E_q \left[\log P(Y, W, z_{\rightarrow}, z_{\leftarrow}, \Phi, \Pi, B | \kappa, \eta, \beta, \alpha) - \log q \right] \quad (4)$$

$$\begin{aligned}
L_{\Delta} = E_q & \left[\sum_t \sum_{p,q} \log \left(B_{g,h}^{Y_{p,q}^t} \frac{e^{-B_{g,h}}}{Y_{pq}^t!} \right) + \sum_t \sum_{pq} \log \left(\prod_k (\pi_{p,k}^{z_{p \rightarrow q} = k}) \right) + \sum_t \sum_{p,q} \log \left(\prod_k (\pi_{q,k}^{z_{p \leftarrow q} = k}) \right) + \right. \\
& \sum_p \log \left[\prod_k (\Pi_{p,k})^{\alpha_k - 1} \cdot \frac{\Gamma(\sum \alpha_k)}{\prod_k \Gamma(\alpha_k)} \right] + \sum_t \sum_p \sum_{w \in R_p^t} \log \left(\prod_{u \in V} (\bar{z}^T \phi_u)^{w=u} \right) + \\
& \sum_k \log \left[\prod_{u \in V} (\phi_{k,u})^{\beta_k - 1} \cdot \frac{\Gamma(\sum \beta_k)}{\prod_k \Gamma(\beta_k)} \right] + \sum_{g,h} \log \left(B_{g,h}^{\kappa-1} / \eta^{\kappa} \Gamma(\kappa) \cdot \exp(-B_{g,h}/\eta) \right) \\
& - E_q \left[\sum_t \sum_{p,q} \log \left(\prod_k (\Delta_{z_{p \rightarrow q}, k})^{z_{p \rightarrow q} = k} \right) + \sum_t \sum_{p,q} \log \left(\prod_k (\Delta_{z_{p \leftarrow q}, k})^{z_{p \leftarrow q} = k} \right) + \right. \\
& \sum_k \log \left[\prod_k (\Pi_{p,k})^{\Delta_{\Pi_{p,k}} - 1} \frac{\Gamma(\Delta_{\Pi_p})}{\prod_{k=1} \Gamma(\Delta_{\Pi_{p,k}})} \right] + \sum_k \log \left[\prod_{u \in v} (\Phi_{k,u})^{\Delta_{\Phi_{k,u}} - 1} \frac{\Gamma(\Delta_{\Phi_k})}{\prod_{u \in v} \Gamma(\Delta_{\Phi_{k,u}})} \right] + \\
& \left. \sum_{g,h} \log \left[\frac{B_{g,h}^{\Delta_{\kappa}=1}}{\Delta_{\eta}^{\Delta_{\kappa}} \Gamma(\Delta_{\kappa})} \exp(-B_{g,h}/\Delta_{\eta}) \right] \right] \tag{5}
\end{aligned}$$

Equation 5 is the variational lower bound of the log likelihood function which is to be maximized. There are terms like $E_q \left[\sum_{g,h} \log \left(B_{g,h}^{\kappa-1} / \eta^{\kappa} \Gamma(\kappa) \cdot \exp(-B_{g,h}/\eta) \right) \right]$ which can be obtained by taking derivation of the partition function of the exponential family form of gamma distribution. I still have to figure out an effective way to evaluate $E_q \left[\sum_t \sum_p \sum_{w \in R_p^t} \log \left(\prod_{u \in V} (\bar{z}^T \phi_u)^{w=u} \right) \right]$ which Chong suggested (and Eric too in the last meeting) to evaluate by intriducing an additional latent variable \bar{z}_p which is a realization of the average $\frac{\sum_{q \in Q} z_{p \rightarrow q}}{|Q|}$.

References