
Modeling Structured Interaction in Large Online User Forums

Anonymous Author(s)

Affiliation

Address

email

Abstract

We present here an approach to model online social forums that respects the structure of the discussion and thus inturn provides researcher with unique insights into the myriads of user forums in the todays online communties. We bring together the structure of the forum network as well as texts posted in a model that respects the structure of the conversation in the forum. The work also focuses on large scal interaction and provides an efficient approximate estimation technique that is scalable **i am thinking of using stochastic variational for scalability if time permits** . Analysing Wikipidia edit and cancer patient online user forums using this technique provides us interesting insights into these online user communities.

1 Introduction

There have been flood of online forums in recent decade and consequently so have been the focus of academic research and industry on online social networks. Analysing online social networks and user forums have been approached using various perspective such as graph/network based [11, 10] , probabilistic graphical model [1], combined network & text mining based [7, 8] approaches. But very few of the approaches have taken into the structural framework in which the conversation in online forums happen. This is important to correctly model the interaction as well as the content posted by the users while interacting with the user community. E.g. in an onlne forum there are topic-threads and users post their responses on this thread after possibly reading trhough the responses of other users in this thread. And the users possibly posts multiple times on the thread in froms of replies to other posts in the thread. For analysing the user interaction, besides taking into account the user interaction network and the text posted, it becomes imperative that the structure of the conversation must also be taken into account. This will enable us to gain deeper insights into user behavior in the online community. Very few research work has tried to bring the forum structure in the analysis of online communities.Our work here besides bringing network modeling and text mining together puts in the forum structure in the model to provide a more robust analysis of the user interactions.

The current work also focuses on analysing large scale user interaction in big online social forums. We provide a variational approximation based estimation technique that is scalable to big forums with thousands of users. **Also write about stochastic approximation if given we have time and we can get it working**

2 Related Work

Our project lies at the intersection of graph clustering and social emergence. Wikipedia's talk pages are an instance of a large social community where we can observe social emergence. Sociologists

divide emergents into various levels [9] e.g. Individual-specific interactions, Ephemeral emergent interactions, Stable emergent interaction etc.

For role-identification and clustering users based on roles in online communities, White et al.[13] proposed a mixed-membership model that obtained membership probabilities for discussion-forum users for each statistic (in- and out-degrees, initiation rate and reciprocity) in various profiles and clustered the users into “extreme profiles”. Ho et al. [7] presented TopicBlock that combines text and network data for building a taxonomy for a corpus.

Griffiths et al. [5] described a generative model for Blei et al.’s LDA model [2] using an MCMC algorithm for bayesian inference queries on the model.

The LDA model and MMSB models were combined by Nallapati et al. [8] using the Pairwise-Link-LDA and Link-PLSA-LDA models where documents are assigned membership probabilities into bins obtained by topic-models.

For simultaneously modeling topics in bilingual-corpora, Smet et al. [12] proposed the Bi-LDA model which generates topics from the target languages for paired documents in these very languages. The end-goal of their approach is to classify any document to one of the obtained set of topics.

For modeling the behavioral aspects of entities and discovering communities in social networks, several game-theoretic approaches have been proposed (Chen et al. [3], Yadati and Narayanam [14]).

Our work is unique in this context as it tries to bridge the gap between community discovery and social emergence. A very popular approach for network clustering is to use a generative framework to infer the underlying structure of communities in a graph. Airolodi et al.[1] proposed a mixed membership stochastic block model to infer community structure in a network. Here they use a generative scheme which captures membership of a user to multiple communities using latent variables. There are other variants of mixed membership stochastic models, e.g. Ho et al.[6] that study evolving clusters over time varying networks while Fu et al.[4] explore dynamic mixed membership models.

3 Approach

Online forums have a generally a specific structure that provides a lot of context to all the interactions among the users. Ignoring this in the analysis makes us lose a lot of precious information as we will see in later sections. Here we describe a typical forum and the answers that we plan to obtain.

3.1 structure in online forums

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.

3.2 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 $p \in \mathcal{N}_t$
 - * 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 $w \sim \phi(w|z_{(p \rightarrow q, t)}, \forall q \in Q_{p, t})$

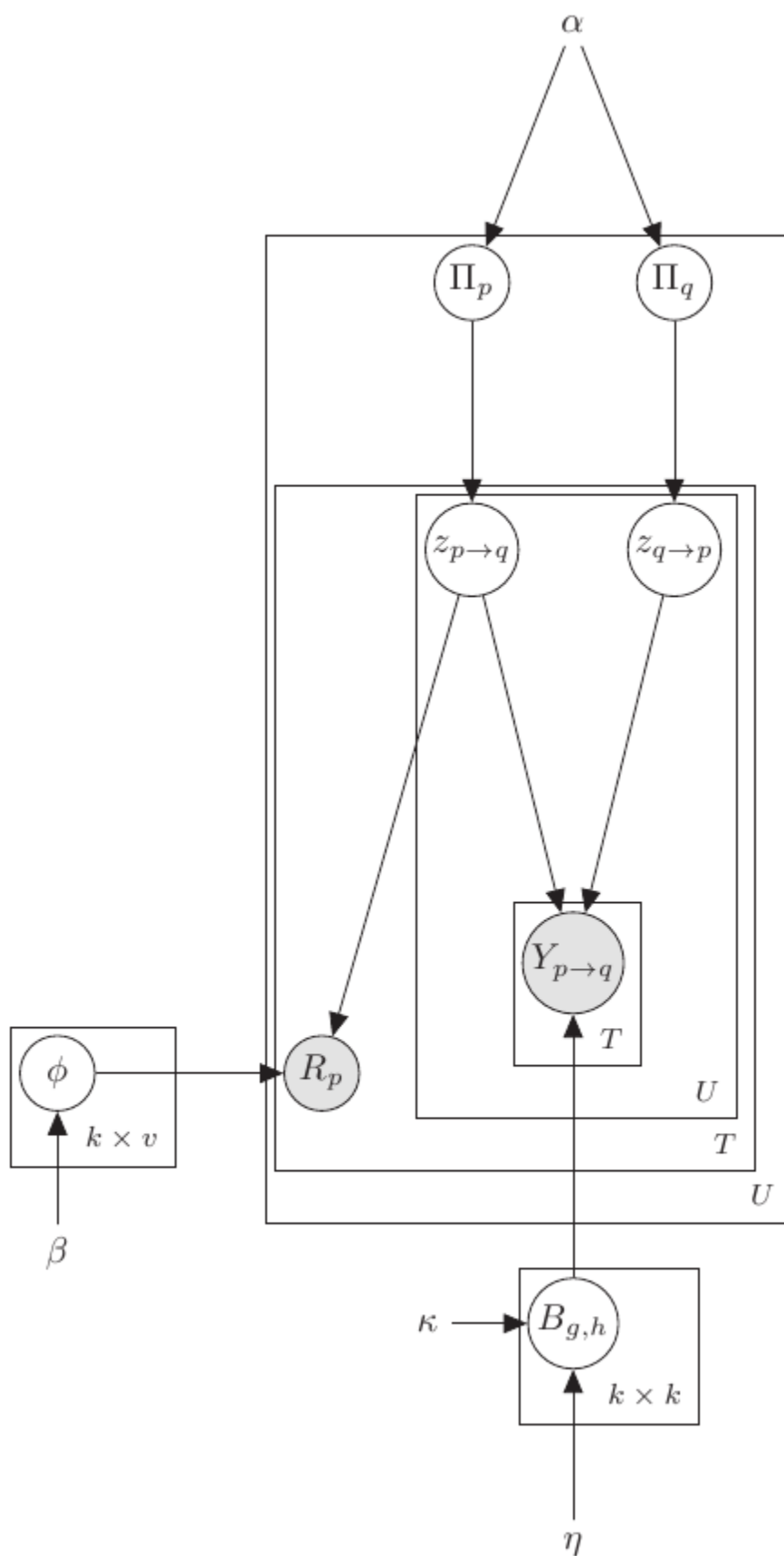
The data likelihood for the model in figure 1

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

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



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

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

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

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

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 $\frac{\sum_{q \in Q} z_{p \rightarrow q}}{|Q|}$.

4 Datasets

5 Experiments & Results & Evaluations

6 Conclusion & Future Work

References

- [1] Edoardo M. Airolidi, David M. Blei, Stephen E. Fienberg, and Eric P. Xing. Mixed membership stochastic blockmodels. *J. Mach. Learn. Res.*, 9:1981–2014, June 2008.
- [2] D.M. Blei, A.Y. Ng, and M.I. Jordan. Latent dirichlet allocation. *the Journal of machine Learning research*, 3:993–1022, 2003.

- 270 [3] Wei Chen, Zhenming Liu, Xiaorui Sun, and Yajun Wang. A game-theoretic framework to
271 identify overlapping communities in social networks. *Data Min. Knowl. Discov.*, 21(2):224–
272 240, September 2010.
- 273 [4] Wenjie Fu, Le Song, and Eric P. Xing. Dynamic mixed membership blockmodel for evolving
274 networks. In *Proceedings of the 26th Annual International Conference on Machine Learning*,
275 ICML ’09, pages 329–336, New York, NY, USA, 2009. ACM.
- 276 [5] T.L. Griffiths and M. Steyvers. Finding scientific topics. *Proceedings of the National Academy*
277 *of Sciences of the United States of America*, 101(Suppl 1):5228–5235, 2004.
- 278 [6] Q. Ho, L. Song, and E. Xing. Evolving cluster mixed-membership blockmodel for time-
279 evolving networks. In *International Conference on Artificial Intelligence and Statistics (AIS-*
280 *TATS)*, 2011.
- 281 [7] Qirong Ho, Jacob Eisenstein, and Eric P. Xing. Document hierarchies from text and links.
282 In *Proceedings of the 21st international conference on World Wide Web*, WWW ’12, pages
283 739–748, New York, NY, USA, 2012. ACM.
- 284 [8] Ramesh M. Nallapati, Amr Ahmed, Eric P. Xing, and William W. Cohen. Joint latent topic
285 models for text and citations. In *Proceedings of the 14th ACM SIGKDD international con-*
286 *ference on Knowledge discovery and data mining*, KDD ’08, pages 542–550, New York, NY,
287 USA, 2008. ACM.
- 288 [9] R. Keith (Robert Keith) Sawyer. *Social emergence : societies as complex systems / R. Keith*
289 *Sawyer*. Cambridge ; New York : Cambridge University Press, 2005. Formerly CIP.
- 290 [10] Jianbo Shi. Learning segmentation by random walks. In *In Advances in Neural Information*
291 *Processing*, pages 470–477. MIT Press, 2000.
- 292 [11] Jianbo Shi and Jitendra Malik. Normalized cuts and image segmentation. *IEEE Trans. Pattern*
293 *Anal. Mach. Intell.*, 22(8):888–905, August 2000.
- 294 [12] Wim De Smet, Jie Tang, and Marie-Francine Moens. Knowledge transfer across multilingual
295 corpora via latent topics. In *Proceedings of the 15th Pacific-Asia conference on Advances in*
296 *knowledge discovery and data mining - Volume Part I*, PAKDD’11, pages 549–560, Berlin,
297 Heidelberg, 2011. Springer-Verlag.
- 298 [13] Arthur White, Jeffrey Chan, Conor Hayes, and Brendan Murphy. Mixed membership models
299 for exploring user roles in online fora, 2012.
- 300 [14] Narahari Yadati and Ramasuri Narayanam. Game theoretic models for social network analysis.
301 In *Proceedings of the 20th international conference companion on World wide web*, WWW
302 ’11, pages 291–292, New York, NY, USA, 2011. ACM.
- 303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323