
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 it provides researchers with unique insights into the myriads of user forums in the today's online world. We bring together the structure of the forum as well as texts posted there in a model that respects the structure of the conversation in the forum. The model incorporates the strength of interaction among users as well, as compared to binary case that cares just for presence or absence of interaction. The work has a scalability focus and provides an efficient approximate estimation technique that is scalable to very large forums (*I am thinking of using stochastic variational for scalability if time permits.*). Analysing Wikipedia edits and cancer patients online user forums using this technique provides us interesting insights into these communities. We also perform a set of prediction tasks to validate our claim further (*We need to decide this soon, though it hinges on once the final model is coded and we have analysed the data using it*)

1 Introduction

There have been a 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 [7, 6], probabilistic graphical model [1], combined network & text mining based [3, 5] based approaches. But very few of these have taken into account the structural framework in which the conversation in online forums happen. This is important to correctly model the interaction as well as the contents posted by the users during their conversation with the user community. E.g. in an online forum there are topic-threads and users post their responses on this thread after possibly reading through the responses of other users in this thread. And the users possibly posts multiple times on the thread in the form of replies to other posts in the thread. For analysing such a user interaction it becomes imperative that the structure of the conversation must also be taken into account besides taking into account the user interaction network and the text posted. This enables us to gain deeper insights into user behavior in the online community that was not possible earlier. Very few research works have tried to bring the forum structure in the analysis of online communities. This is what set our work apart from the past works, our approach here besides bringing network modeling and text mining together adds in the forum structure in the model to provide a more robust analysis of the user interactions. The model also incorporates strength of interaction among the users by incorporating interaction counts as compared to MMSB model which just looks presence or absence of link [1]. In the process we discover interesting online communities and social phenomena.

The current work also focuses on analysing large scale user interactions 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 community discovery, structure modelling and text mining. Wikipedia's talk pages are an instance of a large social community where we can observe users networking with each other as well as posting content in a structured way. Similar phenomena are observed across almost all social networking websites and online forums.

For role-identification and clustering users based on roles in online communities, White et al.[9] 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". In a similar work, Ho et al. [3] presented TopicBlock that combines text and network data for building a taxonomy for a corpus. The LDA model and MMSB models were combined by Nallapati et al. [5] 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. [8] proposed the Bi-LDA model that generates topics from the target languages for paired documents in these very languages. The end-goal of their approach is to classify any document into 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. [2], Yadati and Narayanam [11]).

Ho et al. [4] provide a unique triangulated sampling schemes for scaling mixed membership stochastic block models [1] to hundreds of thousands users based communities.

None of the works above address the structure of the information flow in an online community. Our work is unique in this context as it tries to bridge the gap between community discovery and structured interaction. We propose a novel modelling scheme that takes into account the network information, user contents as well as structure of the interaction and is scalable to big online forums.

3 Approach

Online forums generally have a specific structure that provides a lot of context to all the interactions occurring among the users. Ignoring this in the analysis makes researchers lose a lot of precious information as we will see in later sections. Here we describe a typical forum & related structure and the answers that we are looking for.

3.1 Structure in online forums

In an online forum when two users interact in a thread or through a post they probably bring from their own individual point of views or come from possibly different communities. It is valuable to know which topic/community they each belong to in that interaction. When a user U is representing community C out of all communities that he is part of, he tailors his post content accordingly to suit the explicit or implicit community norms. Knowing the style of community specific text content provides a lot of information about the community in general. It also provides information about what role user U plays when he is in community C . In online forums multi-user interactions happen a lot i.e. in a thread a user can post by addressing to another specific user but he is also addressing other users in the thread explicitly or implicitly. Modeling this phenomenon would bring our model closer to reality. This knowledge can be modelled by aggregating users posts across a thread, though not across the whole of the forum. We will elaborate on this more in the generative story section. Another interesting property of such structured conversations is that there is an inherent bias towards the thread starter or in turn topic of the thread. It would be interesting to see what insights this knowledge provides given that a model can make use of such an information (right now our graphical model doesn't support this but it would be interesting to see how this can be brought in. It might not be too difficult to do this).

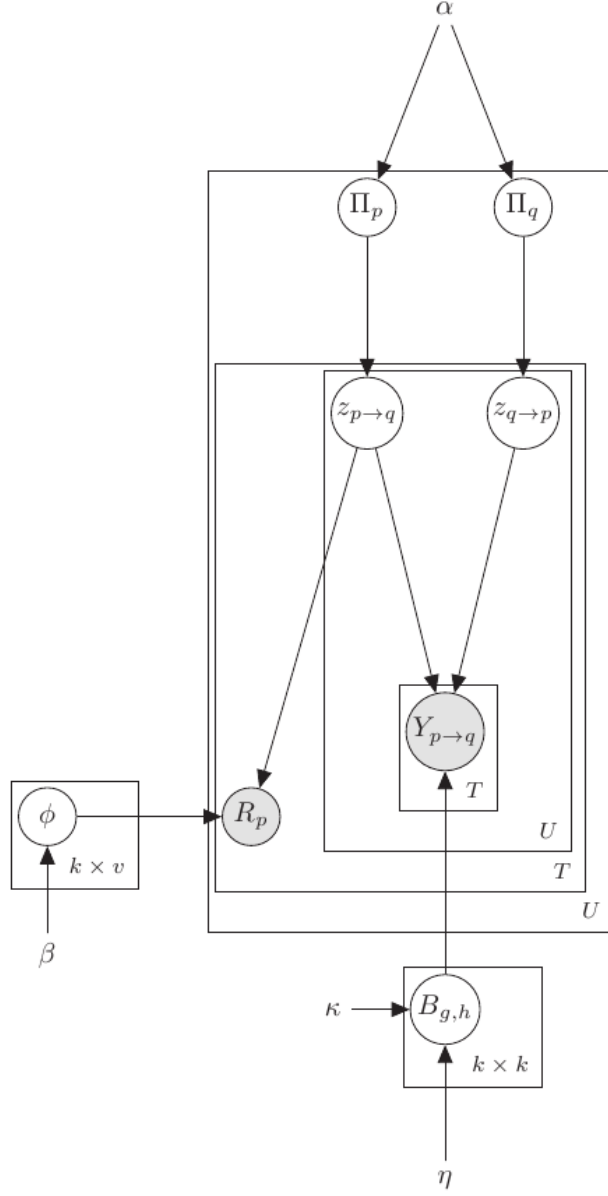


Figure 1: This graphical model takes into account multi-way interaction among users in a thread simultaneously

3.2 Graphical model & generative story

Based on the discussions above our graphical model is designed as 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 . This helps us incorporate the knowledge that a user's post is influenced by all the posts of other users present on the thread.

The generative process for the model 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

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

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

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

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

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. An effective way to evaluate $E_q \left[\sum_t \sum_p \sum_{w \in R_p^t} \log \left(\prod_{u \in V} (\bar{z}^T \phi_u)^{w=u} \right) \right]$ is by introducing an additional latent variable \bar{z}_p which is a realization of the average $\frac{\sum_{q \in Q} z_{p \rightarrow q}}{|Q|}$. So figure 1 will be modified slightly in future where R_p is drawn from \bar{z} and \bar{Z} is drawn from $z_{p \rightarrow q}$. This was suggested by Chong as well as Eric but I finally figured out the equation of the variational approximation for this. I will update the equations once I have coded it and verified..

4 Datasets

We analyse two different forums: 1) Wikipedia talk pages, and 2) Cancer forum. The two datasets mentioned above represent two different sets of online gatherings which helps us generalize our claims.

4.1 Wikipedia talk pages

Wikipedia currently hosts more than four million articles on a wide range of topics. Quality control on Wikipedia occurs through discussions on the Wikipedia talk pages. Every article on Wikipedia has a corresponding talk-page. Contributors to Wikipedia discuss edits by other users, topics that can be used to extend the article, the veracity of the article's contents etc. Talk-pages provide functionality for threaded discussions that are used as dialog among users. This rich structured discussion manifests itself as a social network that can be mined and studied. A standard Wikipedia talk page consists of topics which hold discussion threads. For building our dataset, we used a snapshot of Wikipedia on the 1st of October 2012 [10]. We built a parser and extracted the thread structure in the talk-pages to build the matrices. There are 20,000 users in our datasets that span across 30,000 talk pages (These figures will change depending on whether we want to incorporate more or less users in future). The talk pages become the threads in context of our graphical model.

4.2 Cancer Forum

The cancer forum is an online forum where users discuss about their cancer treatment and anything else under the sun. Here again the conversation happens in a structured way where users post their responses on a thread by thread basis. Users also call each other by their names (or nick-names)

while posting in many cases. This forum has around 3000 users and 10,000 threads, and a user on average posts around 120 words in a post.

Need to expand more in the dataset section with more numbers and stats. Though the stats would be much clearer as and when we perform the experiments

5 Experiments & Results & Evaluations

Following are the experiments that we are planning to do

1. Analyse and interpret the communities discovered on Wikipedia talk pages and Cancer forum providing interesting insights and telling how modelling forum structure helps.
2. Validate the above discovered communities either through held out likelihood or perplexity.
3. Introduce a prediction task; there were several suggestions:
 - a) predict the topic of the posts by a user in a thread where for training we have hand-labeled thread posts (suggested by Chong)
 - b) Hold out some users on some threads and predict whether the user is going to post on the held-out threads

6 Conclusion & Future Work

1) Currently our model just picks up one signal for the network component (MMSB part) i.e. in other words we are just modelling one type of interaction or just one graph, but as we did for the SEI model (tensor model) there are multiple types of networks/graphs/interactions. Future work can incorporate this.

2) Adding temporal dimension to this model would be a very interesting idea. E.g. how threads evolve over time, or how user behavior changes, or how new communities emerge in the forum etc.

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] Wei Chen, Zhenming Liu, Xiaorui Sun, and Yajun Wang. A game-theoretic framework to identify overlapping communities in social networks. *Data Min. Knowl. Discov.*, 21(2):224–240, September 2010.
- [3] Qirong Ho, Jacob Eisenstein, and Eric P. Xing. Document hierarchies from text and links. In *Proceedings of the 21st international conference on World Wide Web, WWW '12*, pages 739–748, New York, NY, USA, 2012. ACM.
- [4] Qirong Ho, Junming Yin, and Eric P. Xing. On triangular versus edge representations — towards scalable modeling of networks. In Peter L. Bartlett, Fernando C. N. Pereira, Christopher J. C. Burges, Lon Bottou, and Kilian Q. Weinberger, editors, *NIPS*, pages 2141–2149, 2012.
- [5] Ramesh M. Nallapati, Amr Ahmed, Eric P. Xing, and William W. Cohen. Joint latent topic models for text and citations. In *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining, KDD '08*, pages 542–550, New York, NY, USA, 2008. ACM.
- [6] Jianbo Shi. Learning segmentation by random walks. In *In Advances in Neural Information Processing*, pages 470–477. MIT Press, 2000.
- [7] Jianbo Shi and Jitendra Malik. Normalized cuts and image segmentation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 22(8):888–905, August 2000.
- [8] Wim De Smet, Jie Tang, and Marie-Francine Moens. Knowledge transfer across multilingual corpora via latent topics. In *Proceedings of the 15th Pacific-Asia conference on Advances in knowledge discovery and data mining - Volume Part I, PAKDD'11*, pages 549–560, Berlin, Heidelberg, 2011. Springer-Verlag.
- [9] Arthur White, Jeffrey Chan, Conor Hayes, and Brendan Murphy. Mixed membership models for exploring user roles in online fora, 2012.

324 [10] Wikipedia. Wikipedia data dump october-01-2012. [http://dumps.wikimedia.org/](http://dumps.wikimedia.org/enwiki/20121001/)
325 [enwiki/20121001/](http://dumps.wikimedia.org/enwiki/20121001/), October 2012.
326
327 [11] Narahari Yadati and Ramasuri Narayanam. Game theoretic models for social network analysis.
328 In *Proceedings of the 20th international conference companion on World wide web*, WWW
329 '11, pages 291–292, New York, NY, USA, 2011. ACM.
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377