

Stochastic Block Models with Multiple Continuous Attributes

Natalie Stanley, Thomas Bonacci, Roland Kwitt, Marc Niethammer, Peter J. Mucha

ABSTRACT

Stochastic block models (SBMs) are probabilistic models for community structure in networks in which nodes within a community are assumed to be connected to nodes within and between communities in a uniform, characteristic way. Typically, only the adjacency matrix is used to perform SBM parameter inference. In this paper, we consider circumstances in which nodes have an associated vector of continuous attributes. The model assumes that the attributes associated with nodes in a network's community can be described by a particular multivariate Gaussian model. Moreover, in this augmented SBM, the objective is to learn the SBM and multivariate gaussian parameters describing each community. While there are recent examples in the literature that combine connectivity and attribute information, we are to our knowledge the first to consider the effect of multiple continuous attributes. We highlight the usefulness of our model for two network prediction tasks: collaborative filtering and link prediction. As a result of fitting this attributed stochastic block model, one can predict the attribute vector or connectivity patterns for a new node in the event of the complementary source of information (connectivity or attributes, respectively). While our approach is the first stochastic block model (to our knowledge) to include multiple continuous attributes, we highlight the ability of our approach in two tasks: link prediction in collaborative filtering. In this work, we demonstrate the usefulness of our model in two different types of biological networks.

First, we consider a network between a set of subjects encoding the similarity in the bacterial species counts in their microbiomes. For each subject, we consider attributes, such as, BMI, nationality, and age. Second, We consider a protein interaction network encoding regulatory relationships, augmented with experimental data about function and abundance.

Data for tasks 1 and 2: <https://www.nature.com/articles/ncomms5344> <http://pubs.acs.org/doi/full/10.1021/pr401258d>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

WOODSTOCK '97 El Paso, Texas USA

© 2017 ACM. ISBN 123-4567-24-567/08/06...\$15.00

DOI: 10.475/123-4

Keywords

Stochastic Block Model, Networks, Community Detection, Attributes, Image Analysis

1. INTRODUCTION

1.1 Network data and node attributes

Uncovering patterns in network data is a common pursuit across a range of fields, such as in biology [10], medicine [8] and computational social science [7]. A powerful way to analyze mesoscale structural organization within a network is with community structure [14, 11]. In this pursuit, the objective is to identify cohesive groups of nodes with a high density of within-group connections and few between-group connections. Numerous approaches exist to accomplish this task [14, 11], but typically only the adjacency matrix conveying the wiring of the graph is taken into account. In certain applications, each node in a network is equipped with additional information (or particular attributes) that was not implicitly taken into account in the construction of the network. For example, one could consider a collection of attributes measured for individuals in a social network (i.e., age, income, level of education). In this work, we seek to incorporate these attributes in the community detection task, such that communities are able to effectively combine both sources of information to understand the organization of components in a system.

Further motivation for the incorporation of node attributes in community detection is also inspired by detectability of communities, as well as possible impact in the image analysis community. The detectability limit in networks is a theoretical notion of the relationship between the number of within-community connections to the number of between-community connections that allows for algorithms to accurately detect community structure [12, 5]. In circumstances where the number of within-community connections is sparse, we propose that the incorporation of attribute information will make communities easier to detect. Further, graph-based image segmentation has shown promise [15, 2], however the probabilistic notion of communities in this context is not well explored.

1.2 Related work

The incorporation of node attributes in network analysis tasks is not well explored, particularly in the context of probabilistic network models. In a modularity based approach, the authors of [3] modified the classical modularity

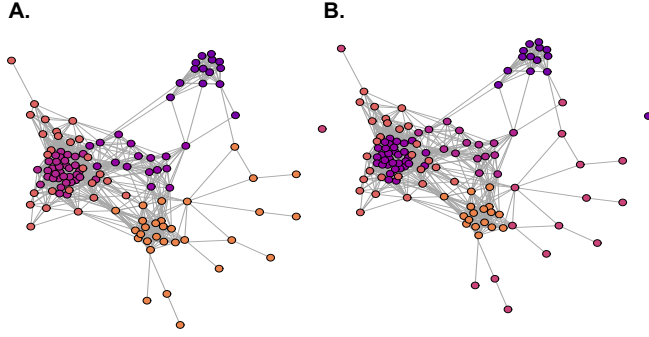


Figure 3: **Visualization of fitting attributed model to microbiome subject similarity network:** Nodes are colored by one of 5 community assignments as predicted with the attributed SBM.

Pre-Processing We extracted a subset of the subjects being from Eastern Europe, Southern Europe, Scandinavia, and the United States and constructed a between-subject similarity network between only individuals who had a BMI measurement. This resulted in a network of 121 nodes, where each edge is the pearson correlation between their microbial compositions. We then removed all edges in the network with weight (correlation) < 0.7 .

Constructing Node Attributes We clustered samples (nodes) in this network using the Louvain community detection method into one of six clusters. We then built a random forest model to predict each sample’s cluster assignment according to their metagenomic profiles. The attribute for each node was then the probability distribution of belonging to each of the 6 clusters of samples. This example highlights a novel way to combine two sources of information, between subject similarity and extra knowledge about the group memberships of the subjects. In figure 3A-B, we plot this subject similarity network and color the nodes by their community assignments using the classic and attributes stochastic block models, respectively. We also computed the normalized mutual information (NMI) between the node-to-community assignments identified from the partition under the louvain algorithm with the results obtained fitting the class and attributed SBM. Since the attributes of the nodes were the probability of belonging to each of the clusters identified under the Louvain algorithm partition, it was expected to observe a higher similarity between the Louvain partition with the attributed SBM partition than with the SBM partition. This is exactly what we observed, with 0.78 being the NMI between the louvain algorithm partition and the classic SBM. Using the attributed SBM, we observed an NMI of 0.94 with the Louvain algorithm partition.

Microbiome Link Prediction In figure 4a-b we show the results for the link prediction in the sample similarity network. To perform link prediction, we selected 25 pairs of nodes (i.e. edge stubs) from our 121-node subject network. In our prediction task, we sought to model the ability to use only the attribute information and our fitted model to predict whether there was a link. Repeating this experiment 25 times, we obtained a distribution of area under the curve (AUC) values in figure 4A. and plotted the ROC curves

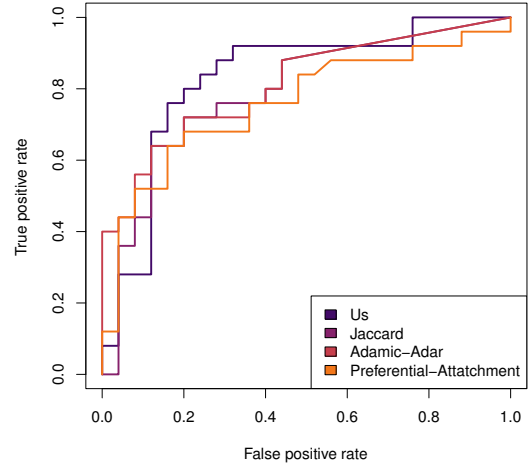


Figure 4: **Link Prediction on the microbiome subject similarity matrix:**

corresponding to the experiment closest to the median for all link prediction methods.

Microbiome Collaborative Filtering Figure X shows the results for the collaborative filtering task on the subject microbiome network

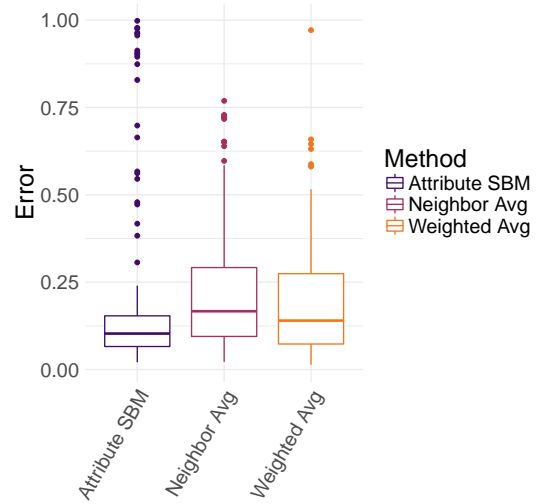


Figure 5: **Collaborative Filtering Accuracy in Microbiome Data:** For each of the 121 nodes, we fit a model to the remaining 120 node network and given the node’s closest network neighbors try to compute its attribute value. Error is the 2-norm between true and predicted attributes.

Here, we sought to predict the 6 dimensional attribute vector for each node.

sahfkahfkjh insert more text

5.2 Protein Interaction Network Results

We also apply our attributed SBM approach to the protein interaction network presented in [1]. This network represents interactions between proteins, predicted from the liter-

ature. Associated with each node (protein), is a vector of experimentally observed modifications resulting from the exposure of cancer cells to a chemotherapeutic drug. Proteins were able to undergo six possible modifications. While communities in this network should reflect functional relatedness among proteins, we also expect that members of a community should share similarities in the observed modification type.

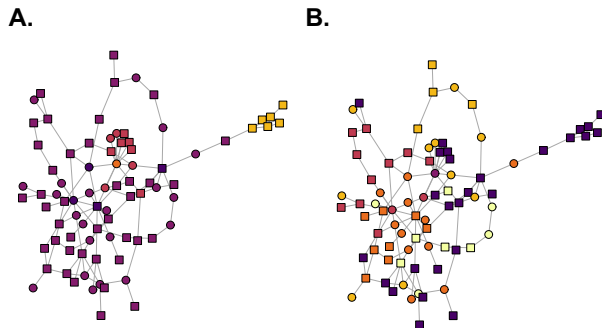


Figure 6: **Protein interaction network**

Data Pre-Processing: We downloaded the network data and the modification information from the supplement of [1].

Constructing Node Attributes: For each node, we constructed its attributes vector as a vector of length 6, where each entry is a binary indicator for which of the 6 modifications was experimentally observed.

Figure 6A-B show the results of fitting a classic SBM and attributed SBM, respectively. The 6 possible modifications exist for 3 biological processes that can either increase or decrease. The node shape reflects whether the modification for a node was an increase (square) or decrease (circle). Nodes are colored by their assignment into 1 of x communities. (mention something about the intuition behind the color groupings).

Next, we studied the entropies of these binary node classifications in each of the communities according to the regular and attributed SBM partitions. The hypothesis was that using the attributed SBM, we should have lower entropy of these two labels within communities because the attribute component of the model should assist in creating communities that are not only spatially relevant but also agree in attributes.

Link Prediction in the Protein interaction network

6. DISCUSSION

Detectability problems

7. DISCUSSION

To be filled in.

8. REFERENCES

- [1] T. Bonacci, S. Audebert, Stephane, L. Camoin, E. Baudet, and G. Bidaut. Identification of new mechanisms of cellular responses to chemotherapy by

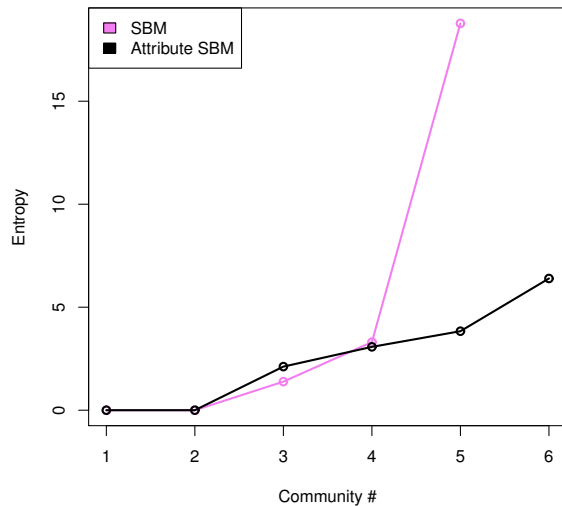


Figure 7: **Community Entropies**

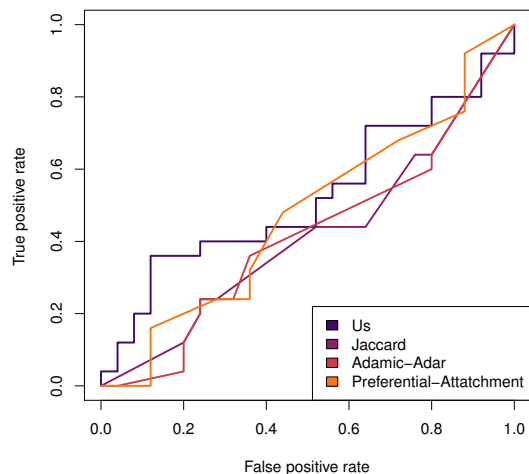


Figure 8: **Community Entropies**

tracking changes in post-translational modifications by ubiquitin and ubiquitin-like proteins.

- [2] A. Browet, P.-A. Absil, and P. Van Dooren. Community detection for hierarchical image segmentation. In *Combinatorial Image Analysis*, pages 358–371. Springer, 2011.
- [3] D. Combe, C. Largeron, M. Géry, and E. Eged-Zsigmond. I-louvain: An attributed graph clustering method. In *Advances in Intelligent Data Analysis XIV*, pages 181–192. Springer, 2015.
- [4] J.-J. Daudin, F. Picard, and S. Robin. A mixture model for random graphs. *Statistics and computing*, 18(2):173–183, 2008.
- [5] A. Decelle, F. Krzakala, C. Moore, and L. Zdeborová. Asymptotic analysis of the stochastic block model for

modular networks and its algorithmic applications. *Physical Review E*, 84(6):066106, 2011.

- [6] A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the royal statistical society. Series B (methodological)*, pages 1–38, 1977.
- [7] D. Greene and P. Cunningham. Producing a unified graph representation from multiple social network views. In *Proceedings of the 5th Annual ACM Web Science Conference*, pages 118–121. ACM, 2013.
- [8] J. Guinney, R. Dienstmann, X. Wang, A. de Reyniès, A. Schlicker, C. Soneson, L. Marisa, P. Roepman, G. Nyamundanda, P. Angelino, et al. The consensus molecular subtypes of colorectal cancer. *Nature medicine*, 2015.
- [9] P. W. Holland, K. B. Laskey, and S. Leinhardt. Stochastic blockmodels: First steps. *Social networks*, 5(2):109–137, 1983.
- [10] D. B. Larremore, A. Clauset, and C. O. Buckee. A network approach to analyzing highly recombinant malaria parasite genes. *PLoS Comput Biol*, 9(10):e1003268, 2013.
- [11] J. Leskovec, K. J. Lang, and M. Mahoney. Empirical comparison of algorithms for network community detection. In *Proceedings of the 19th international conference on World wide web*, pages 631–640. ACM, 2010.
- [12] R. R. Nadakuditi and M. E. Newman. Graph spectra and the detectability of community structure in networks. *Physical review letters*, 108(18):188701, 2012.
- [13] M. Newman and A. Clauset. Structure and inference in annotated networks. *arXiv preprint arXiv:1507.04001*, 2015.
- [14] M. A. Porter, J.-P. Onnela, and P. J. Mucha. Communities in networks. *Notices of the AMS*, 56(9):1082–1097, 2009.
- [15] J. Shi and J. Malik. Normalized cuts and image segmentation. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 22(8):888–905, 2000.
- [16] J. Yang, J. McAuley, and J. Leskovec. Community detection in networks with node attributes. In *Data mining (ICDM), 2013 IEEE 13th international conference on*, pages 1151–1156. IEEE, 2013.