

DISCUSSION OF “COAUTHORSHIP AND CITATION NETWORKS FOR STATISTICIANS”

BY PEDRO REGUEIRO, ABEL RODRÍGUEZ AND JUAN SOSA

University of California

1. General comments. We would like to start by congratulating the authors for taking the initiative in gathering this very interesting and novel data set. Given the substantial amount of work involved in scrapping the data, their focus on four periodicals (Journal of the American Statistical Association, Journal of the Royal Statistical Society Series B, *Biometrika* and *Annals of Statistics*) is understandable. However, this relatively narrow choice raises some concerns. The most obvious one relates to the robustness of the results to the choice of periodicals, particularly for authors/papers concentrating on areas for which specialized high-quality alternative publications exist. Two examples are biomedical applications and Bayesian methods. Furthermore, although the four journals selected are mainly methodological, the inclusion of the Applications and Case Studies section of JASA was unfortunate. Manuscripts published there can be expected to have more in common with papers published in the *Annals of Applied Statistics* or the Journal of the Royal Statistical Society, Series C than with manuscripts in the Theory and Methods section of JASA itself.

The analysis in the paper feels a little bit like a “fishing expedition.” The paper lacks a clear question that motivates and shapes the data collection. The use of multiple alternative methods (both for constructing the networks and for analyzing them) yielding different results also detracts from a sense of purpose. This is a pity because there are a number of interesting questions that could be explored if the data collection exercise had been slightly expanded with a clear objective in mind. Some examples include the following:

1. What are the main drivers of collaboration in statistics?
2. How have the collaboration networks evolved over time?
3. How likely are researchers to publish with their Ph.D. mentors as time goes by?
4. Are there regional biases in citation and/or publication patterns?
5. How prevalent is “self-referencing” (both at the author and journal level)?

The feeling of lack of focus is reinforced by the fact that the clusters generated by the community identification methods in the paper are puzzling. For example, the fact that only three clusters are identified in the connected component of the author citation network is quite surprising. This small number could be driven by

TABLE 1
*Top-ten authors based on eigenvalue centrality for the
 Coauthorship (A) and Coauthorship (B) networks*

Coauthorship (A)	Coauthorship (B)
Peter Hall	Joseph G Ibrahim
Raymond J Carroll	Hongtu Zhu
Yanyuan Ma	Weili Lin
Aurore Delaigle	Yimei Li
Hans-Georg Müller	Xiaoyan Shi
Enno Mammen	Bradley S Peterson
Hua Liang	Daniel B Rowe
Alexander Meister	Hongyu An
Fang Yao	Wei Gao
Naisyin Wang	Yashen Chen

the fact that the two-mode relational data has been projected into one-mode networks, by the fact that the resulting one-mode network is converted into a binary network instead of treated as weighted, or by the lack of formality in the choice of the number of networks. This observation also suggests that coauthorship and citation data separately are not enough to create a taxonomy of the statistical literature; multiple sources of information are necessary to produce a more fine-grained partition that better reflects most people's understanding of the community structure. The remainder of the discussion explores some of these issues.

2. Eigenvalue centrality. We complemented the centrality measures presented in the paper with the eigenvalue centrality (see Table 1). Interestingly, note that while the top-ten list contains two of the three highly central authors identified in Section 3 of the manuscript (Peter Hall and Raymond Carroll), it does not contain the third (Jinquiang Fan). This suggests that, even though all these three authors are highly collaborative themselves, the coauthors of Jinquiang Fan tend to be less collaborative than those of Peter Hall and Raymond Carroll. Furthermore, note that whereas the top-ten lists generated by other centrality measures substantially overlap for the two networks, the lists for eigenvalue centrality are completely different, suggesting that this metric is much more sensitive to the procedure used to dichotomize the weighted network.

3. How many communities? Reanalyses using stochastic block models. In this section we explore using a stochastic block model to fit the Coauthorship (A) and Coauthorship (B) networks, and compare the results to those presented in Section 4 of the manuscript. The model assumes that the edges in the network $(y_{i,i'})$ are conditionally independent given the set of interaction probabilities Θ , and that

the probability of observing an edge between two vertices i and i' depends exclusively on the community membership of i and i' ,

$$(1) \quad y_{i,i'} \stackrel{\text{ind}}{\sim} \text{Ber}(\theta_{\xi_i, \xi_{i'}}),$$

where ξ is a vector of *community indicators* taking values in $\{1, 2, \dots, K\}$ and K is the maximum number of communities. In a Bayesian setting the model is completed with priors for the parameters Θ and ξ . For simplicity, the interaction probabilities are assigned independent uniform priors, $\theta_{k,l} \stackrel{\text{ind}}{\sim} \text{Uni}[0, 1]$ and the prior on the community indicators are constructed by assuming the entries of ξ are exchangeable and follow a categorical distribution in $\{1, 2, \dots, K\}$,

$$(2) \quad \Pr(\xi_i = k | w_k) = w_k; \quad i = 1, 2, \dots, I,$$

with weights vector $\mathbf{w} \sim \text{Dir}(\frac{\alpha}{K}, \frac{\alpha}{K}, \dots, \frac{\alpha}{K})$, such that the marginal likelihood from this model converges to the marginal likelihood of a mixture model with a Chinese restaurant process prior. The parameter α , which controls the effective number of components $K^* \leq K$, is assigned a Gamma prior.

Notice that the communities in the stochastic block model have an interpretation that is slightly different from the communities obtained from the algorithms considered by the authors (NSC, BCPL, APL and SCORE). Specifically, rather than groups of vertices with a relatively large number of edges *within* and small number of edges *across*, communities in the stochastic block model are formed by vertices that interact similarly across the network and, thus, these clusters can be thought of as functional structures in the network. In the setting of coauthorship networks this distinction turns out to be relevant as, a priori, one would expect to observe disassortative communities that arise from multiple students collaborating almost exclusively with their advisors and, at the same time, assortative communities that represent close-knit research groups with few outside collaborators. Therefore, the stochastic block model seems a natural modeling choice, as it is capable of simultaneously recovering assortative and disassortative mixing in a network.

3.1. Coauthorship network (A). In this section we examine *Coauthorship network (A)*. Following the manuscript, we focus on the largest connected component of this network. A first difference that can be appreciated in Figure 1 is the fact that the stochastic block model supports the existence of three—rather than two—communities. As seen in this plot, Peter Hall, Raymond Carroll, Jianqing Fan and Tony Cai are clustered into a single community that can be interpreted as being composed by the network’s “hubs.” Although a direction for further investigation would be the use of degree-corrected block models [Karrer and Newman (2011)], the fact that Joseph Ibrahim is not included in this community, despite having the fourth largest number of publications in the network, is evidence that the partition obtained by the stochastic block model is not exclusively driven by vertex degree.

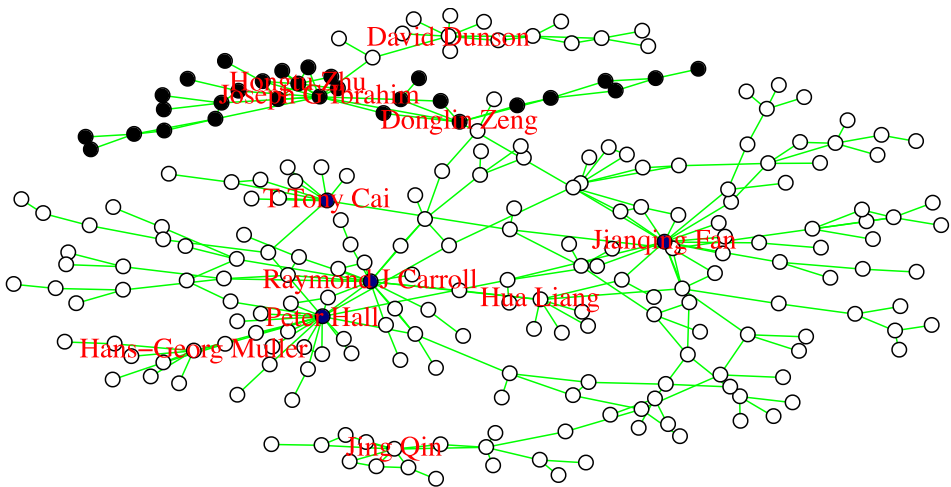


FIG. 1. Communities resulting from fitting a stochastic block model to Coauthorship Network (A).

Table 2 compares the partition from the stochastic block model to those obtained from NSC, BCPL, APL and SCORE using the Adjusted Random Index (ARI) and the Variation of Information (VI). Here, it can be seen that the communities from the stochastic block model are closest to those from APL. In particular, Community 2 in the SBM corresponds almost perfectly to the Carroll–Hall community identified in the main paper, and Community 3 corresponds to the North Carolina community (Community 1 in the SBM is made of the four high-degree authors identified above, which APL assigns to the Carroll–Hall community).

3.2. *Coauthorship network (B)*. We also examined *Coauthorship network (B)* where two researchers are connected with an edge if they share one or more publications, focusing again on the largest connected component. In this case the stochastic block model suggests six communities in the data, although two of them contain only a very small fraction of the vertices in the network (6 and 2 observations, respectively).

To compare the partitions obtained from the stochastic block model with those derived from NSC, BCPL, APL and SCORE, Table 3 presents ARI and VI mea-

TABLE 2
Adjusted Random Index and Variation of Information comparing the communities from the stochastic block model to the communities obtained by the different methods presented in Ji and Jin (2016) using the giant component of Coauthorship (A)

ARI/VI	SCORE	NSC	BCPL	APL
SBM	0.64/0.43	−0.05/0.99	0.08/0.95	0.90/0.13

TABLE 3

Adjusted Random Index and Variation of Information comparing the communities from the stochastic block model to the communities obtained by the different methods presented in Ji and Jin (2016) using the giant component of Coauthorship (B)

ARI/VI	SCORE	NSC	BCPL	APL
SBM	0.04/1.57	0.03/1.38	0.00/2.09	0.04/1.11

tures. These indexes suggest that, unlike the case of Coauthorship (A), the communities identified by the stochastic block model have little overlap with any of those identified by other metrics. An inspection of the estimated interaction probabilities Θ suggests that these differences might be driven by the fact that the stochastic block model identifies a couple of disassortive communities.

To investigate this relationship further, Table 4 shows the intersection of the communities from the stochastic block models with those from APL. The stochastic block model suggests that the “HDDA” community can be further partitioned into smaller blocks.

4. Embeddings and combining information from citation and coauthorship networks. An alternative approach to community identification involves first embedding the probabilities in a Euclidean latent “social” space, and then clustering the nodes according to their position in the latent space [e.g., see [Handcock, Raftery and Tantrum \(2007\)](#)]. For example, for an undirected network we could proceed with a two-step approach, where

$$y_{i,i'} \stackrel{\text{ind}}{\sim} \text{Ber}(\Phi(\beta + \mathbf{u}_i^T \mathbf{u}_{i'})), \quad \mathbf{u}_i \stackrel{\text{ind}}{\sim} N_L(\mathbf{0}, \sigma^2 \mathbf{I})$$

with further hyperpriors for β and σ^2 . The dimension L of the latent space is selected using the *Deviance Information Criterion* (DIC) [[Gelman, Hwang and Vehtari \(2014\)](#), Chapter 6]. Once point estimates $\tilde{\mathbf{u}}_1, \dots, \tilde{\mathbf{u}}_L$ are obtained (e.g., the posterior means after the enforcement of an appropriate identifiability constraint),

TABLE 4

Comparison of communities obtained for the stochastic block model and the APL algorithm

		APL		
		Bayes	Biostat	HDDA
SBM	Community 1	14	12	211
	Community 2	2	1	284
	Community 3	2	5	202
	Community 4	8	9	1505
	Community 5	0	0	6
	Community 6	0	0	2

TABLE 5
Comparison of communities obtained for the latent space modeling (LS) and the APL algorithm

		APL	
		North Carolina	Carroll–Hall
LS	1	23	159
	2	8	46

communities can be determined using a finite mixture model for clustering, such as that implemented in the R package `mclust` [Fraley and Raftery (2002), Fraley et al. (2012)].

We use this procedure on the giant component of the Coauthorship (A) network. DIC selects a three-dimensional social space (i.e., $L = 3$), and `mclust` identifies $K = 2$ communities. Table 5 compares the communities obtained using this procedure with those identified by APL; note that the results vary substantially.

The approach we just described can be extended to two or more adjacency matrices $\mathbf{Y}_1, \dots, \mathbf{Y}_J$ defined over a common set of I actors by letting

$$y_{i,i',j} | \beta_j, \mathbf{u}_{i,j}, \mathbf{u}_{i',j} \stackrel{\text{ind}}{\sim} \text{Ber}(\Phi(\beta_j + \mathbf{u}_{i,j}^T \mathbf{u}_{i',j})),$$

with

$$\beta_j \stackrel{\text{ind}}{\sim} N(\mu, \tau^2), \quad \mathbf{u}_{i,j} | \eta_i, \sigma^2 \stackrel{\text{ind}}{\sim} N(\eta_i, \sigma^2 \mathbf{I})$$

and $\mu \sim N(0, b_\mu^2)$, $\eta_i \stackrel{\text{ind}}{\sim} N(\mathbf{0}, b_\eta^2 \mathbf{I})$, $\tau^2 \sim \text{IG}(a_\tau, b_\tau)$ and $\sigma^2 \sim \text{IG}(a_\sigma, b_\sigma)$. Community identification proceeds then by clustering the “average” random position η_1, \dots, η_I .

We used the extended model to obtain a set of communities of authors that combines coauthorship and citation information. To facilitate comparisons, we focus again only on those authors included in the giant component of the Coauthorship (A) network. The joint model identifies $K = 5$ communities, again with $L = 3$. Table 6 compares these 5 communities to those identified by the model

TABLE 6
Comparison of communities obtained for the latent space modeling based only on coauthorship data (Co A) and both coauthorship and citation information (Joint)

		Joint				
		1	2	3	4	5
Co A	1	48	37	40	34	23
	2	0	1	2	4	47

based only on coauthorship data. Note that while the second original community remains largely unaffected by the inclusion of citation information (roughly corresponding to our new community 5), the first one is split into four subgroups.

REFERENCES

- FRALEY, C. and RAFTERY, A. E. (2002). Model-based clustering, discriminant analysis, and density estimation. *J. Amer. Statist. Assoc.* **97** 611–631. [MR1951635](#)
- FRALEY, C., RAFTERY, A. E., MURPHY, T. B. and SCRUCICA, L. (2012). mclust Version 4 for R: Normal mixture modeling for model-based clustering, classification, and density estimation.
- GELMAN, A., HWANG, J. and VEHTARI, A. (2014). Understanding predictive information criteria for Bayesian models. *Stat. Comput.* **24** 997–1016. [MR3253850](#)
- HANDCOCK, M. S., RAFTERY, A. E. and TANTRUM, J. M. (2007). Model-based clustering for social networks. *J. Roy. Statist. Soc. Ser. A* **170** 301–354. [MR2364300](#)
- Ji, P. and JIN, J. (2016). Coauthorship and citation networks for statisticians. *Ann. Appl. Stat.* To appear.
- KARRER, B. and NEWMAN, M. E. J. (2011). Stochastic blockmodels and community structure in networks. *Phys. Rev. E* (3) **83** 016107, 10. [MR2788206](#)

DEPARTMENT OF APPLIED MATHEMATICS AND STATISTICS
BASKIN SCHOOL OF ENGINEERING
UNIVERSITY OF CALIFORNIA
1156 HIGH STREET
SANTA CRUZ, CALIFORNIA 95064
USA
E-MAIL: pregueir@soe.ucsc.edu
abel@soe.ucsc.edu
jsosamar@soe.ucsc.edu