

structure analysis, where the positions of atoms on each protein can be determined by X-ray crystallography and where it is suspected that (subregions of) proteins of similar shape have similar biological function, but where the labellings are unknown. There are some similarities to the problem of comparing clusterings determined by the $\{\mu_g\}$ in different Markov chain Monte Carlo simulations.

Tony Lawrance (*University of Warwick, Coventry*)

It is a pleasure to contribute to the discussion after the experts have spoken. My experience of this area and paper is the 2-hour train journey from Warwick to London, a distance of nearly zero according to the metric which this paper induced. My first point is to enquire how the analysis can address the measurement of friendliness or interactions of the actors that are involved, rather than their groupings. Secondly, I noticed that the modelling is predicated on a conditional independence assumption, which it must be tough to validate and is probably a matter of faith. I could not immediately see any attention in the paper given to assessing the fit of the model, and the choice of prior forms seemed clever, but I wonder how much can they influence the final groupings? Wider empirical validation, replacing monks and monasteries by lecturers and departments, would satisfy me more generally. My final observation concerns the microscopic pie charts, noting that they take the development of invisible graphics to a new level, at least judging by my monochrome preprint. Overall, I thought that this paper was a nice blend of methodology and application.

The following contributions were received in writing after the meeting.

Edoardo M. Airolidi (*Carnegie Mellon University, Pittsburgh*)

The authors' work with the *latent space clustering* methodology provides an impressive demonstration of the use of hierarchical models for identifying groups of nodes from observed connectivity patterns. Modelling choices based on sociological principles, i.e. transitivity and homophily, increase its appeal as an exploratory tool for the analysis of social networks. The methodology proposed goes only part way, however, towards addressing fundamental issues that arise in the statistical analysis of social networks.

The *stochastic blockmodel of mixed membership* in Airolidi (2006) and Airolidi *et al.* (2007a) offers an alternative approach with different insights on latent aspects underlying network structure. Models in this family also posit the existence of an unknown number of clusters; however, they replace latent positions with mixed memberships $\pi_{1:N}$, which map nodes to (one or more) clusters, and add a *latent blockmodel* B that specifies cluster-to-cluster hierarchical relations. These parameters are directly interpretable in terms of notions and concepts that are relevant to social scientists, and better suited to assist them in extracting substantive knowledge from noisy data, ultimately to inform or support the development of new hypotheses and theories. Therefore, inference about $\pi_{1:N}$ and B is crucial for the analysis of data.

Applying this to Sampson's data demonstrates both linkages and differences. Our version of the Bayes information criterion also suggests the existence of three factions among the 18 monks, but our groupings are different. In Fig. 12, Romul and Victor (two of Sampson's Waverers) stand out; and so do Greg and John who were expelled first from the monastery. The mixed membership map is specified by using node-specific latent vectors $\pi_{1:18}$, independent and identically distributed samples from a three-dimensional symmetric Dirichlet(α) distribution. The map of hierarchical relationships among factions is specified by a 3×3 matrix of Bernoulli hyperparameters B , where $B(i, j)$ is the probability that monks in the i th faction relate to those in the j th faction. Other features that are relevant to data analysis are the marginal probability of a relation ($\pi'_n B \pi_m$) and the relation between the number of clusters and dimensionality of the latent simplex.

Our models allow a focus on issues such as membership of monks in factions, and this could lead to the formation of a social theory of failure in isolated communities, which is capable of testing with longitudinal data. In Airolidi *et al.* (2007), we provide full details on specification, estimation and interpretation for both the Sampson and the adolescent friendship network examples.

Julian Besag (*University of Washington, Seattle*)

I would like to comment on the authors' choice of examples. After all, social networks have been around for a long time and there is an abundance of data, so we should be expecting more than purely illustrative analyses by now.

In their first example, the authors deem an edge from i to j to exist if i cites j at any of his three interviews. In general, if clusters change over time, such a rule could lead to spurious results. Moreover, as regards social science, I would assume that the temporal development of clusters, including their creation, coalescence, fragmentation and destruction, is of more interest than their static properties. Although three

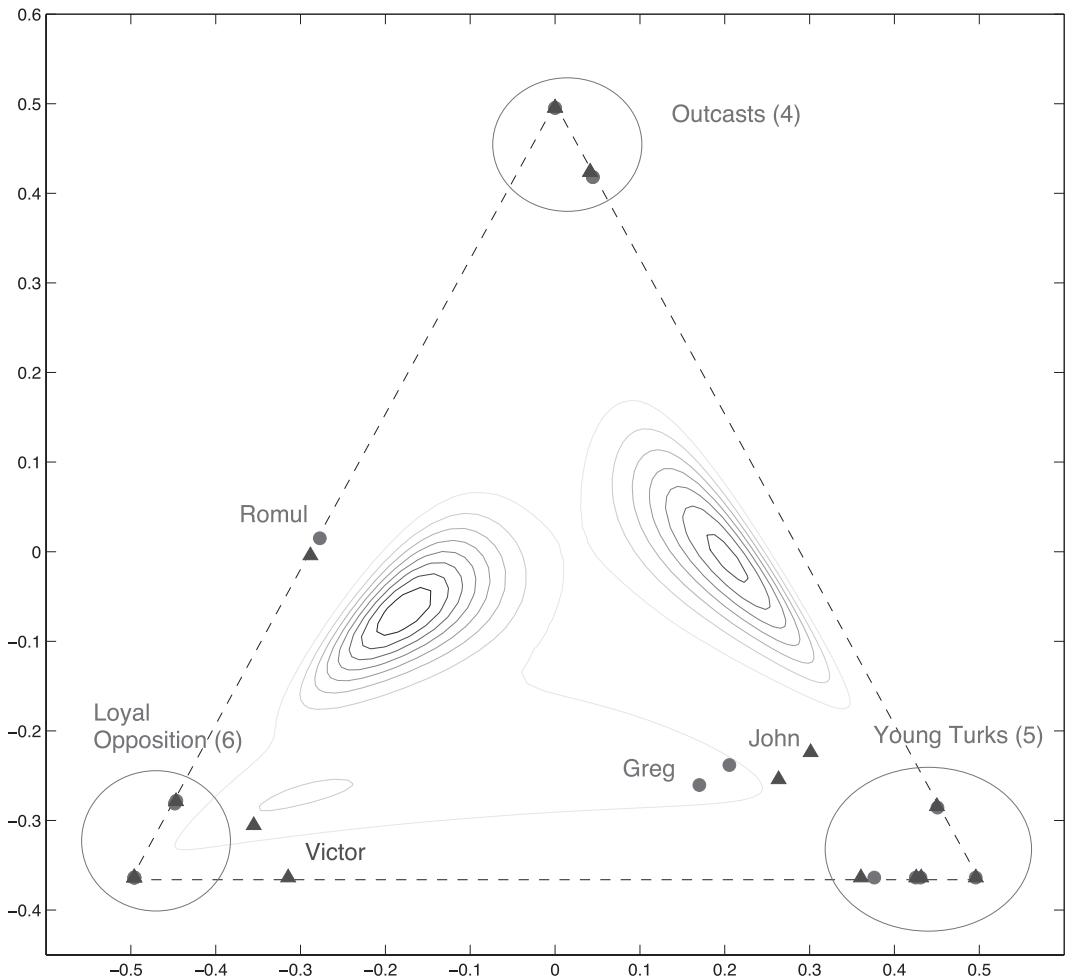


Fig. 12. In the reference simplex, circles and triangles correspond to mixed membership vectors of individual monks, $\pi_{1:18}$ (circles were obtained with $B = I_3$ and $\alpha = 0.01$, whereas triangles were obtained with $B = I_3$ and $\alpha = 0.58$ —estimated via an empirical Bayes method): an arbitrary one-to-one projection situates the Gaussian mixture of Table 1 in the simplex

time points are probably too few for meaningful analysis, more extensive space–time networks could have been chosen. Such analysis is particularly important for communicable diseases. Note that, in setting up space–time models, multiple changes in edge configurations can occur (almost) instantaneously, though this is sometimes overlooked.

As regards their second example, do the authors have a justification for focusing on one particular school out of 132? It seems to me that they should at least have analysed a small sample of schools. And why were no covariates included, particularly the grade of student? To claim success in extracting grade as an important clustering attribute suggests to me that the authors are too easily satisfied. Their secondary conclusions are plausible and could have been checked in other schools. The general point here is the effect of including cluster attributes as covariates, which is allowed in their original formulation but apparently not in their examples. How does this affect cluster identification?

Lastly, do the authors have anything to add about the relevance of their approach to the huge networks that for example AT&T and Microsoft researchers must deal with and for which quite different methods are used? Is this merely a computational issue or is it that exploratory techniques are more appropriate?