

# Clearing the FOG : Fuzzy, Overlapping Groups for Social Networks

George B. Davis\*, Kathleen M. Carley

*CASOS, ISRI, SCS, Carnegie Mellon University, 5000 Forbes Ave., Pittsburgh, PA, 15223, USA*

---

**Abstract:** Humans are well known to belong to many associative groups simultaneously, with various levels of affiliation. However, most group detection algorithms for social networks impose a strict partitioning on nodes, forcing entities to belong to a single group. Link analysis research has produced several methods which detect multiple memberships but force equal membership. This paper extends these approaches by introducing the FOG framework, a stochastic model and group detection algorithm for fuzzy, overlapping groups. We apply our algorithm to both link data and network data, where we use a random walk approach to generate rich links from networks. The results demonstrate that not only can fuzzy groups be located, but also that the strength of membership in a group and the fraction of individuals with exclusive membership are highly informative of emerging group dynamics.

---

---

This work was supported in part by the National Science Foundation under the IGERT program, 9972762, for training and research in CASOS and the Office of Naval Research under Dynamic Network Analysis program (N00014-02-1-0973). Additional support was provided by CASOS - the Center for Computational Analysis of Social and Organizational Systems at Carnegie Mellon University. The views and conclusions contained in this document are those of the author and should not be interpreted as representing the official policies, either expressed or implied of the National Science Foundation, the Office of Naval Research, or the U.S. government.

- Corresponding author; Tel.: +01 412 580 3790; E-Mail: [gbd@cs.cmu.edu](mailto:gbd@cs.cmu.edu) (George B. Davis)

## Introduction

Since the earliest days of social network research, accurate detection of cohesive group entities has been an attractive and elusive goal. Group structure can be used for high-level descriptions of complex networks, to support or contest theories about underlying processes influencing social interactions, and to detect strengths or vulnerabilities of social structures and individual positions in a variety of contexts. These goals have important applications in a wide range of fields, including anthropology, sociology, organization science, economics, management, and security and intelligence programs.

Typically, group detection has consisted of dividing nodes into discrete partitions indicating mutual association. However, common sense and empirical analysis (Freeman, 1992) support the view that humans are capable of simultaneously filling many roles in many contexts, such that a strict partitioning may prevent detection of the true group entities in a graph. To better understand modular structure in networks, we must develop models which allow for multiple memberships and varied levels of membership.

In this paper we build off several link analytic group detection methods, due to Kubica, Schneider and Moore (2003b) and Battacharya and Getoor (2004), which allow for relaxed partitioning by permitting individuals under certain conditions to participate in multiple groups. We refine the representation of group structure by permitting varying strengths of association from members to group entities, and present an algorithm that generates such groupings from link data using a stochastic model of link emission from group entities and a maximum-likelihood clustering method. To analyze the utility of the fuzzy overlapping group model, we make comparison to groupings by anthropological observations and prior algorithms. Our results suggest that this approach is capable of identifying groups that are confirmed by existing quantitative methods as well as expert ethnographic analysis, while providing additional information about overlap between groups and individuals who play multiple roles. This additional information facilitates understanding emergent behavior in the groups.

The remainder of the paper is organized as follows. In section 2, we provide a brief background on existing group detection methods. We also discuss the generation of link data from networks, so that we can apply our link analytic method to network datasets. In section 3 we describe our approach (termed *FOG*

for “Fuzzy Overlapping Grouper”) in two subsections: one proposing a stochastic model of the way groups generate link data, and another introducing a corresponding maximum likelihood method for inferring groupings based on evidence. In section 4, we present performance results on the FOG algorithm and use FOG to analyze two well-studied real world datasets: Sampson’s monastery survey data (1968) and Davis, Gardner and Gardner’s southern women (1941), comparing our results to previous groupings on the same. In the conclusion, we discuss FOG’s potential contribution to group analysis based on our results, and identify additional work necessary.

## **Background and Related Work**

### **Defining “Group”**

Theoretically, we consider a group a set of entities which experience the same membership relation with respect to the same external entity, real or abstract. In the social sphere, this can take many forms. For example: a formal organization like a board of directors, an implicit organization like a circle of friends, a demographic quality such as hair color, or even the set of individuals uniquely affected by an external force, such as the victims of a flu epidemic.

This paper is concerned specifically with *cohesive groups*, divisions which exhibit more associations within groups than between them. This operational definition may at first seem to restrict the varieties of group we can detect, yet we can imagine interaction data in which each of the above categories of group would leave such a trace. For example, members of a board of directors might occur together on the recipient list for formal memos and meeting announcements. Individuals afflicted by the same communicative disease might tend to be clustered in space and time in hospital records. To some extent, measuring this definition of cohesion depends on being able to clearly measure both the presence and absence of links between entities -- a property inherent in social network data, but less obvious in link data, which we define and discuss in the next section. In link data, a stochastic model must fill the roll of defining what comprises a concentration of links. In the next section we describe several widely used algorithms to detect cohesive groups in both types of data.

Several non-cohesive group types are also popular in sociological research, particularly structural similarity. Under this criteria, entities are grouped if their interaction patterns are similar -- that is, if they interact with the same other entities or classes of entities. The group they experience membership with in this

case is a structural role. Distinct techniques exist for discovering structurally similar nodes, block modeling (Lorrain and White, 1971) and CONCOR (Breiger *et al.*, 1975) being the most commonly used. However, some relationships exist between structurally similar and cohesive groups. Membership in a strongly cohesive group can directly constitute a structural role, as members tend to interact almost exclusively with members. Even in some cases where individuals sharing a role do not interact, a simple inversion of the data may make them detectable as a cohesive group. For example, consider a dataset collected from vendors in which each lists the clients they sell to. If we invert the association data and create a list of vendors for each client, cohesive groups in this new dataset will collect vendors who filled similar roles in the original data. In other cases, cohesion and structural similarity are almost completely orthogonal. In a strictly hierarchical organization, for example, only the most contrived algebraic convolutions of interaction data will place non-interacting mid-level managers together.

### **Detecting Cohesive Groups**

In network data, cohesive groups are commonly identified with clustering algorithms, such as hierarchical clustering or FACTIONS search (Borgatti *et al.*,

2005), which partition entities into groups maximizing some internal link density metric. Although iterative clustering and partitioning algorithms are typically queried to provide a specified number of group entities, the hierarchical structure may in some cases serendipitously provide information about groups-within-groups. Several heuristic methods which would not be described as clustering also see wide use. Girvan and Newman have introduced a heuristic method (2002) with speed advantages based on iterative removal of high-betweenness edges, which has seen several extensions and applications (Clauset *et al.*, 2004; Newman, 2004a; 2004b; Newman and Girvan, 2004).

The methods listed so far each operate on *network data*, which we defined as a 2-dimensional matrix with entries representing the existence or strength of association between each pair of interacting entities. However, some contexts are better represented by distinct interaction events (*link data*), which may carry redundant associations or simultaneous associations of more than two entities. In this paper we're concerned with "undirected" link data, represented as an unordered set of links (the *evidence*), each of which is an unordered set of entities in which each entity is assumed to have the same relation to an observation (i.e., "signed meeting roster", or "was observed in photograph"). Link data can also

be considered a type of two-mode data, with one mode as the entity and the other as links.

Data mining communities have produced several methods for extracting group entities from this type of data, including the GDA model / k-Groups algorithm (Kubica *et al.*, 2003b) and Battacharya and Getoor's iterative deduplication method (2004). These algorithms partition link data to infer groups which maximize the likelihood of observing the given data, according to a stochastic model. The fact that groups are built from links, rather than the observed entities themselves, produces the advantage that individuals may belong to more than one group. The method we introduce in this paper extends on these methods by allowing varying levels of association from entities to groups. This relaxation is intended to allow our group models to more tightly fit the data and to represent a wider variety of associative structures.

Another existing technique with some similarity to the FOG framework is Latent Dirichlet Allocation (LDA) (Blei *et al.*, 2003), a recently introduced stochastic model for machine learning mixed memberships. Airolidi *et al.* have adapted the model to examine single-mode network data (2005), yielding novel clusterings in protein-protein interaction networks (Airolidi *et al.*, 2006). The primary

distinction between FOG and relational LDA models is that LDA allows a single observed to be explained by a mixture of groups, whereas FOG assumes that a single social context is associated with a given observation.

### **Link Data from Network Data**

Link analysis and network analysis have grown out of distinct communities, despite being frequently applied to the examination of the same interaction phenomena. In many ways, grouping research has become an intersection point in which practitioners of both fields are attempting to capitalize on the strengths of the other. Link analysis researchers approach group models as an opportunity to characterize structure and dependence in interaction data which is too often analyzed as though observations were independent. Analysts who have traditionally used graph theoretic approaches to examine network data are incorporating statistical models and significance tests to improve their ability to reason about noisy data and support claims about the significance of structural characteristics in their networks. For frameworks such as FOG to see the widest use (and scrutiny), we must develop translation techniques that allow data in one format to be examined using algorithms for the other. These translations

must account for disparate data qualities emphasized in the two branches of analysis.

Since small changes in network structure can have a large impact on the graph theoretic measures used in network analysis, translations from link data to network data are designed to reduce noise as much as possible. Many network datasets begin life as something more closely resembling link data. Lists of interactions or survey responses are “flattened” into a matrix of pair wise interactions using summation, cutoffs, or reciprocation criteria depending on the interaction being studied and the network type desired. Recently, Kubica *et al.* have presented cGraph, an expectation maximization approach to detecting underlying networks (2003a).

The stochastic link analytic techniques we examined are intended to robustly handle noise given enough data. However, when our source data is limited to the information in an interaction matrix, we run the risk of amplifying any noise present when we generate additional links. We must also tackle the problem of reflecting the structural data contained in the network model in a way that link analytic algorithms can interpret. The naïve approach - interpreting each edge in network data as a single link of two entities - can be detrimental to algorithms,

including FOG-Greedy, which rely on richer link data to assist the clustering process.

In this paper we've adopted the "random tree" solution described by Kubica *et al.* (2003a), in which link data is constructed stochastically by iteratively adding to links entities which are randomly chosen from the neighbors of those already present. Figure X illustrates this process. Since nodes A and B have been visited already, the entire peripheral boundary of C, D, and E are available as our next addition. Note that C has a greater probability of selection than D or E, as there are two links from proceeding to it from our visited structure. If our matrix were weighted, rather than binary, the relative probability of visiting a node on the perimeter would be determined by summing the link weights leading to it from visited nodes. We use this technique on a weighted collation of Sampson's monastery data later in this paper.

Classification experiments using random walks as a kernel have shown that random walks can successfully represent structural features in a graph (Kashima and Tsuboi, 2004). However, to properly interpret generated data, we would like to be able to relate our link generator to a real world process that could have produced observable data. Random walks and trees in social networks have

been used in simulations as analogs to real-world processes, such as knowledge dissemination the spread of a disease (Christley *et al.*, 2005), or a search for information (Page and Brin, 1998). As such, we could feel comfortable interpreting groups found in random trees as sets of individuals who are likely to be exposed to many of the same ideas or illnesses, or look to each other for information which they lack.

## Data Sets

**Sampson Monastery.** We chose Sampson's monastery dataset (1968) as a testbed for the FOG framework because it is one of the datasets most widely discussed in social grouping literature. Sampson conducted a survey in which novice monks at a monastery ranked their compatriots according to four criteria: like / dislike, esteem, personal influence, and consistency with the creed of the monastic order. Sampson made strong arguments for several discrete social groups in the data based on direct anthropological observation. Events confirmed his observations when, during the study, novices of one group resigned or were expelled over religious differences. Samson's surveys may be the dataset that comes closest to providing social data with a labeled "ground truth" for grouping research.

Sampson's monastery is discussed in greater depth in Sampson's original (1969) dissertation, and in the December 1988 issue of the journal *Social Networks*. We compare the groups discovered by FOG to Sampson's and those presented by Reitz in that issue in his introduction of a hierarchical clustering algorithm. Like that paper, we use Breiger's (1975) collation of Sampson's data: for each of the relations "like", "esteem", "influence", and "consistency", the top three positive selections by each individual at time 3 are recorded in a relation matrix. These matrices are summed, yielding a single matrix summarizing the preferential data at that time period<sup>1</sup>. Because FOG analyzes link-based data, we pre-process this matrix to generate links using the techniques described above. The matrix in its entirety is shown below as Table 1.

	ROMUL	BONAVEN	AMBROSE	BERTH	PETER	LOUIS	VICTOR	WINF	JOHN	GREG	HUGH	BONI	MARK	ALBERT	AMAND	BASIL	ELIAS	SIMP	
	10	5	9	6	4	11	8	12	1	2	14	15	7	16	13	3	17	18	
ROMUL	10	0	1	1	0	1	0	0	0	0	0	0	0	0	0	1	0	0	0
BONAVEN	5	0	0	2	0	3	3	0	0	1	0	0	0	0	0	0	0	0	0
AMBROSE	9	0	1	0	0	2	0	2	1	2	1	0	0	0	0	0	0	0	0
BERTH	6	0	1	3	0	4	2	0	0	1	0	0	0	0	0	0	0	0	0
PETER	4	3	1	0	4	0	4	0	0	0	0	0	0	0	0	0	0	0	0
LOUIS	11	0	3	2	0	2	0	2	0	0	0	1	0	0	1	0	0	0	0
VICTOR	8	0	1	2	3	4	2	0	1	0	0	0	0	0	0	0	0	0	0
WINF	12	0	0	0	0	0	0	0	3	3	1	0	2	0	0	0	0	0	0
JOHN	1	0	1	0	0	0	0	1	4	0	1	2	0	1	0	0	1	1	0
GREG	2	0	1	0	0	0	0	1	3	4	0	0	0	3	0	0	0	0	0
HUGH	14	0	0	0	0	0	0	3	4	3	0	4	0	1	0	0	0	0	0
BONI	15	0	0	0	0	0	0	1	2	4	3	0	2	0	0	0	0	0	0
MARK	7	0	0	0	0	0	0	3	0	4	0	2	0	4	0	0	0	0	0
ALBERT	16	0	0	0	0	0	0	1	0	4	0	4	4	0	0	0	0	0	0
AMAND	13	0	4	0	0	0	3	0	0	0	0	0	3	0	0	0	0	0	1
BASIL	3	0	0	0	0	0	0	0	4	0	0	0	0	0	4	0	4	2	0
ELIAS	17	0	0	0	0	0	0	0	0	3	0	0	0	0	1	3	0	3	0
SIMP	18	0	0	0	0	0	0	0	1	4	0	0	0	0	0	3	4	0	0

Table 1. Breiger's (1975) Collation of Sampson Survey Data

**Davis, Gardner, and Gardner's Southern Women.** The southern women dataset (Davis *et al.*, 1941) lists the attendance of 18 women and 14 parties. The parties in this network are precisely the type of linking observation which FOG is designed to analyze without pre-processing. As with the monastery dataset, there exists a labeling for groups based on direct observation rather than algorithmic analyses. Davis et al. used ethnographic analysis, including surveys, to distinguish not only between the two major cliques, but three tiers of centrality within them.

A wide variety of mathematical approaches have been used to reanalyze the data. Freeman performed a comprehensive meta-analysis of 21 such studies (1992), and we analyze our results in response to some of his conclusions. We've also accepted that paper's verdict on which of two conflicting figures in the original work was correct. We reproduce that figure as table 2, below, for reference.

NAMES OF PARTICIPANTS OF GROUP I	CODE NUMBERS AND DATES OF SOCIAL EVENTS REPORTED BY <i>Old City Herald</i>													
	(1) 6/27	(2) 3/2	(3) 4/12	(4) 9/26	(5) 2/25	(6) 5/19	(7) 3/15	(8) 9/16	(9) 4/8	(10) 6/10	(11) 3/23	(12) 4/7	(13) 11/21	(14) 8/3
1. Mrs. Evelyn Jefferson.....	X	X	X	X	X	X		X	X					
2. Miss Laura Mandeville.....	X	X	X		X	X	X	X						
3. Miss Theresa Anderson.....		X	X	X	X	X	X	X						
4. Miss Brenda Rogers.....	X		X	X	X	X	X	X						
5. Miss Charlotte McDowd.....			X	X	X		X							
6. Miss Frances Anderson.....			X		X	X	X							
7. Miss Eleanor Nye.....					X	X	X	X						
8. Miss Pearl Oglethorpe.....						X	X	X						
9. Miss Ruth DeSand.....					X		X	X						
10. Miss Verne Sanderson.....							X	X	X			X		
11. Miss Myra Liddell.....								X	X	X		X		
12. Miss Katherine Rogers.....								X	X	X		X	X	X
13. Mrs. Sylvia Avondale.....							X	X	X	X		X	X	X
14. Mrs. Nora Fayette.....						X	X		X	X	X	X	X	X
15. Mrs. Helen Lloyd.....							X		X	X	X			
16. Mrs. Dorothy Murchison.....								X	X					
17. Mrs. Olivia Carleton.....									X		X			
18. Mrs. Flora Price.....									X		X			

Table 2. Southern women party attendance, reproduced from (Davis *et al.*, 1941)

## The FOG Framework

Grouping methodologies are often introduced as algorithms, although they encompass distinct models, measures, data translations, and validation schemes as well as the model-fitting algorithm itself. To minimize this confusion, we discuss FOG as a framework consisting of several components. The FOG generative model relates link interactions we observe to group entities, which are hidden. The FOG-Greedy algorithm is a simple link-clustering approach to fitting groups of the type described in the model to data. (As we will show, the algorithm does not guarantee optimality and future work may yield a fast

algorithm that finds better fits.) A separate link generation algorithm creates link data from social network data.

### Stochastic Model of Evidence Generation

Since we are trying to infer groups based on link evidence, we define our group membership relation as the tendency to be produced in observations associated with the group. We can alter the strength of the tendency to be included in observations without altering its fundamental character. We formulate the above mathematically as follows.

Consider a set of entities  $E$  and a set of groups  $G$ . Entities are elementary objects whose presence or absence is observable in a set  $L$  of links (which are sets of entities). Groups emit pieces of evidence, which consist of sets of entities which are co-observed. Groups emit with different frequencies, according to a probability distribution  $\bar{\theta}$  across groups such that  $\theta_g$ , for  $g \in G$ , is the probability that any given link was emitted by group  $g$ ;  $\sum_{g \in G} \theta_g = 1$ . Elsewhere in this paper we refer to  $\theta_g$  as the *emission prior*, since it represents our expectation that a piece of evidence will come from a specific group, prior to examining the members

observed in the link. A membership vector  $\vec{g}$ , whose entries are the probabilities that each entity is present in a link that has been emitted by group  $g$ , further describes each group. We write this as  $g_e = P(e \in l \mid g \Rightarrow l)$ , or the shorthand  $P(e \mid g)$ . We will refer to  $g_e$  as membership strength or affiliation. Figure 2 illustrates the hierarchy of objects we have defined.

When considering the likelihood that a particular group would produce a specific link, we must consider not only the probability of observing the entities present in the link but the probability of excluding those not present.

$$P(l \mid g \Rightarrow l) = \left( \prod_{e \in l} g_e \right) \left( \prod_{e \notin l} 1 - g_e \right) \quad (1)$$

The assumption that, in the emissions of a single group, members are emitted completely independently is important to maintaining that the membership relation differs only in intensity between entities. (A joint distribution would imply additional substructure.) Similarly, we assume that links are generated completely independently given the groups and their emission priors, so that the only structure exists between the groups and the entities themselves, and in the relative frequency of emission of the groups. Combining, these we can derive

the likelihood that an entire set of evidence would be produced given a grouping and an emission distribution vector. The factorial coefficient in this equation normalizes for the ordering of the link set, which is irrelevant to our model.

$$P(L | G, \vec{\theta}) = |L|! \prod_{l \in L} \sum_{g \in G} \theta_g \left( \prod_{e \in l} g_e \right) \left( \prod_{e \notin l} 1 - g_e \right) \quad (2)$$

Performance and representation precision (probabilities involved can be extremely small) demand that the above likelihood function be calculated via a log-likelihood transformation. To enable this transformation, we place the restriction that  $g_e \in [p_{\min}, p_{\max}]$ , where  $0 < p_{\min} < p_{\max} < 1$ . This ensures that a group always has some nonzero probability of emitting its least related entity, or excluding from a link even its most significant member.

Previous stochastic models of link generation have included an “error term” under which there is some small probability that a link will be emanated containing entities which do not cohabit any group. This was necessary to allow models to be fit to data without placing extreme penalties on groups which were forced to include outlying entities as equal members to more supported nodes. In FOG, a similar purpose is served by allowing weak memberships and

assuming weak universal memberships, with the advantage that we need no prior beliefs about an error rate.

### The FOG-Greedy Algorithm

With the relationship between groups and evidence described above, we can reduce group detection to an optimization problem which searches for the grouping with the greatest likelihood of generating the observed evidence:  $\arg \max_{G, \vec{\theta}} P(L | G, \vec{\theta})$ . Calculating this explicitly would be intractable, so we propose an estimation algorithm.

Since our model requires that a single group be responsible for emanating each link, we can restrict our search by considering only groups which optimally represent some partition of the data. The group with the highest probability of single-handedly generating a set of links is the one which emits each entity with probability equal to the proportion of the link set in which that entity occurs. We build groups of this sort by iteratively clustering link evidence in a way that ensures links with the greatest similarity are grouped together. For each pair of groups  $g_1, g_2$ , we consider a new group  $g_n$  that would maximize probability of

emitting the combined evidence supporting both groups ( $L_n \leftarrow L_1 \cup L_2$ ). We then calculate the ratio as a heuristic indicating the relative increase in expectation of the underlying links that this merge would cause.

$$\frac{|L_n|^2 P(L_n | g_n)}{|L_1|^2 P(L_1 | g_1) + |L_2|^2 P(L_2 | g_2)} \quad (3)$$

The pair for which this ratio is highest is merged.

The tree in Figure 3, constructed from the southern women dataset, illustrates the hierarchical clustering of evidence. Each intermediate node corresponds to a group tuned to produce evidence of the types found in the leaves below. We define a *horizon* from this tree as a set of nodes such whose children span all of the evidence, for which none is the ancestor of another. A horizon, such as the circled nodes in Figure 3, corresponds to a set of groups which account collectively for all of the observed evidence. If we choose our horizon from the bottom level, groups are tuned to very specific profiles of evidence, so that they are expected to produce any of the few links below them with relatively high probability. As we move up the tree, membership rosters for groups become more complex and the distribution of links which they produce becomes more

entropic, so that the probability of producing any particular link drops exponentially. At the same time,  $\theta_g$  values rise as we ascend the tree, since each group represents a greater proportion of the underlying evidence. Near the top, groups are overly general and fit the evidence underneath poorly, so that, even though  $\theta_g$ 's are high, the total probability of producing the evidence set is quite low.

Unfortunately, reduced  $P(l | g \Rightarrow l)$  outpaces increased  $P(g \Rightarrow l)$  over a climb of the tree, so that there is usually no optimal midpoint that would allow us to discern a “most probable number of group entities”. As such, the operator must currently specify a number of groups,  $k$ , for which to search, effectively deciding on a tolerable tradeoff between a simple model with few groups and a model which most closely fits the evidence but may in fact be over-fit.

## Results and Analysis

### Runtime

The graph in Figure 4 profiles the runtime of the FOG-Greedy algorithm on evidence generated from random groupings with varying number of entities. Memberships for each entity are chosen from a uniform distribution, and several sets of evidence with varying numbers of links are generated, according to the FOG stochastic model described previously. In total, the figure summarizes runtime on 1500 evidence sets from 150 groups. The algorithm was implemented in Python, an interpreted scripting language, and executed on an Intel Pentium 4 machine running at a speed of 3 GHz.

Number of groups (not shown) has a minimal effect on runtime, as the vast majority of calculation is performed on the lower levels of the tree, before the cutoff for number of groups is reached. The effect of the number of links ( $L$ ) being grouped dominates that of any other variable on runtime. Runtimes of FOG-Greedy fall under an  $O(L^3)$  bounding. This is expected from the fact that FOG-Greedy must examine  $O(L^2)$  candidate merges at each of  $L$  levels in the merge tree. Runtime is affected linearly by the number of entities observed, as this number determines the upper bounds on the number of calculations necessary to examine a candidate group. In practice, many of these calculations are memoized, allowing for a tighter bound than discussed.

## Fuzzy Groups in the Monastery

Since Sampson's data consists of pairwise relations, we generated link data using the random-tree technique previously discussed. Ten trees were initiated at each node, each expanding to contain three nodes, and the set was clustered using the FOG algorithm. Results are shown as a two mode (agent  $\rightarrow$  group) network in figure 5. Line thickness indicates the degree of membership. We normalized line thickness within each cluster because it is not necessarily appropriate to compare association levels between groups. This is because our link generation method required exactly three individuals in each observation, artificially deflating the average frequencies of emission in large groups and inflating it in small ones. Nodes have been manually laid out to elucidate membership categories we discuss.

Sampson identified novice 2, Gregory, as the most significant leader of the "young Turks", the liberal newcomers who would be expelled or resign in the coming drama. The members of that group are collected exactly as those affiliated with group A. Gregory's position, as both the most affiliated to the Turks and the only novice with connections to all three groups, suggests a high

degree of centrality both within the young Turks and in the network as a whole. This type of border-spanning centrality has been linked to iconoclasm, power, and stress, painting a vivid picture of factors which may have contributed to Gregory's exit. The official reasons given for his expulsion were excessive independence and arrogance. Could he have been singled out as more dangerous precisely because he had captured the attention and esteem of individuals outside the clique? Rank of affiliation to group A turns out to be a good predictor for the order in which the young Turks left the abbey: the second individual most affiliated with the Turks is Sampson's second identified leader, John (novice 1), who resigned shortly after Gregory's exit, trailed by Mark (7) and the other novices.

Sampson's third and lowest-ranking leader, Winfred (12), does not buck this trend. Although his relatively low association to group A belies his eventual identification as a leader of the Turks, it accurately reflects his placement as the last one to leave the Abbey. Winfred's undistinguished position in our plot illustrates some biases of this analytic method, as well as some peculiarities of his situation. Winfred identified strongly enough with his group that he was completely embedded: all of his incoming and outgoing connections in the survey data lie within the Turks. The result is that, although the random trees in

which he appears are exclusively tied to Group A, he simply doesn't participate in nearly as many total evidence pieces as high-betweenness boundary spanners like Gregory or John. As this shows, our evidence-generation technique could rightly be said to put a premium on individuals with high betweenness.

However, this is defensible when paired with interpretation placing it in a social context. Recall that we generated random trees in order to model iterative interaction processes within the graph, such as the spread of a rumor or the slow accumulation of individuals to a casual gathering. It is easy to imagine an individual with more diverse ties, such as Gregory, being drawn into a wider variety of gatherings. By being a prolific interactor, Gregory may well have defined the Turks to the rest of the community, without necessarily intending to or even identifying exclusively with them.

Evidence supports the distinction between Gregory's "celebrity" and Winfred's "poster child" stances. In Sampson's study, Winfred's leadership was either absent or unobserved in the presence of the two higher-profile leaders, and became clear only after their exit. Winfred's embeddedness seems to reduce his significance at the time of our analysis, but as the split widened between the Turks and the opposition, making positions like Gregory and John's untenable,

Winfred's exclusive loyalty became the crucial element of his in-group leadership.

The membership and leadership of the "loyal opposition" party are similarly gathered around Group B in our plot. Peter (4) and Bonaven (5), who were identified by Sampson as the leaders of the opposition, show the highest affinity for the group. Members described as less attached show less affinity, and one such novice shows a split allegiance to the outcast group.

The absence of any links, save Gregory, between the opposition and the Turks serves to reflect the conflict between the two groups. By contrast, the "outcasts" in group C have several members associated with other groups as well. These cases show that fuzzy memberships can help elucidate not only the complexity of an individual's allegiances, but the character of a group as exclusive or inclusive to interstitial members.

Sampson originally identified a fourth group, but we restricted our analysis to three clusters because the last was not a cohesive group fitting our definition.

Sampson does not describe the "waverers" as a set of individuals allied or interacting with one another, but as being in similar positions of doubt between

the two major groups. Additionally, previous analyses have questioned the distinction between the waverers and the loyal opposition. Our own analysis places two of them, Romul (10) and Victor (8), as weak members of the loyal opposition. From a purely structural perspective they are tied more to the loyal opposition; whatever their mental allegiance. Armand (13) is categorized as an outcast, owing less to his statements of affinity for those individuals than from Basil's (3) and Elias' (17) connections to him.

### **Fuzzy Groups Among Southern Women**

Analyses of the DGG data, including the original, have generally partitioned the women into two cliques<sup>2</sup> that intersect on a few individuals or events. We use a “spectrographic<sup>3</sup>” visualization scheme in figure 6 to present the results of a 2-clustering of the southern women in greater detail than would have been readable in the Sampson analysis. Bars of each color indicate each woman's affiliation with two groups derived from 8 and 6 of the party rosters respectively. Individuals are sorted along the X axis according to the difference in their membership levels, which maximizes the visual distance between the cliques.

---

<sup>2</sup> We use *clique* here to maintain consistency with prior work, not to indicate a graph theoretic relationship.

<sup>3</sup> So named after similarity to overlaid graphs of element density used to differentiate substances in mass spectrometry

We have also included a 2-mode network visualization for comparison to the one we presented for the Sampson data.

The results of our algorithmic approach correspond strongly to the intuitive conclusions of Davis et al. In group A, the core and primary periphery are reproduced precisely as plateaus in the membership levels. Someone attempting to fit our analysis to their mode might draw slightly different tiers for the group B, but the rough ordering of individual affiliations is the same. For both groups, the most peripheral members are seen in the center of our chart, with low levels of affiliation in both groups. Some of these members have been shown to be interstitial; for example Davis et al. report that Ruth (9) was claimed by both cliques in interviews with members. Others, such as Pearl (8) and Verne (10) were only claimed by members of the cliques to which our chart shows greater affiliation.

There are many mathematical studies of the DGG data to which the FOG-Greedy clustering correlates. We'll omit a pairwise comparison, as many of the results are significantly similar to Davis et al.'s intuitive analysis described above, and a

comprehensive meta-analysis has already been accomplished by Freeman.

Instead, we focus on FOG's contribution to one prong of that analysis: the core-periphery structure of the two cliques.

Davis et al.' describe as "core" the individuals that are seldom excluded from their clique's functions. We see that the most affiliated individuals in both groups demonstrate a propensity to appear with the other group as well. This supports the argument we proposed with the Sampson data<sup>4</sup>, that leaders of a group may either arise out of greater participation with other groups than do the less active members, such as those in DGG's "primary" and "secondary" members, or else experience more pressure to do so.

In his meta-analysis, Freeman treated core-periphery as an ordering of individuals for each group, without specifying that centrality in one group promoted distance from the other (although that was a side effect of many techniques compared). FOG results certainly fit that mode, but the juxtaposition of affiliations given above lends itself to an additional breakdown of several interactions. We can separate individuals into several modes of interaction. We have central leaders, such as the novices John or Gregory or the Southern women

---

<sup>4</sup> Since the DGG analysis is based on direct observations rather than synthetic observations from random trees, we do not have the same concern about overemphasizing centrality that we did with the novices.

Nora and Katherine. There are embedded leaders such as Winfred, Laura, or Brenna. There is a loyal second tier in each of the groups we've analyzed, and finally a set of truly interstitial individuals who participate at low levels in both groups.

From our observations of these roles so far, we might issue the prediction that a thoroughly embedded member, such as Brenda or Flora, would flourish if there were a falling out between the two groups. On the other hand, if good relations continued between the groups, our profile of an emergent leader might better fit individuals such as Ruth (9) or Helen (15): those with strong ties to one group, but some degree of participation with the other. Davis et al. do not examine conflict between these cliques and describe no events that would be telling regarding our first hypothesis. However, they completed a larger study of many cliques, in which they used interstitial members to examine relations between social classes associated with each clique. They describe a class of "on the way up" individuals, who participate in events outside their clique in order to socialize with those above them in social class.

## **Discussion and Future Work**

We set out to introduce a new quantitative way of reasoning about the complex relationship between individuals and groups, allowing varied degrees of participation in multiple groups. We proposed the FOG stochastic model, which dictates relationships between individuals, groups, and observable interactions as a generative model for link data. To make FOG a useful analysis tool, we introduced the FOG-Greedy algorithm which fits a model to existing link data. To investigate single mode network data, we implemented a simple method for generating rich multi-entity links from a pairwise network based on a simplistic simulation of interaction processes.

*Validation.* Mathematical approaches to group detection are based on the assumption that groups have a reality outside individual perceptions, which we can detect statistically. It should therefore be possible to empirically validate grouping methods on their ability to predict the outcome of processes in which groups play a role. Currently, the closest we have to this sort of predictive test is to compare our analysis to that of anthropologists like Sampson, who were able to relate their intuitive observations to unforeseen events in the social group. Can fuzzy grouping rediscover social patterns that stood out to ethnographers in the field?

In the two datasets we studied, the answer is yes. The discrete groups identified by both Sampson and the DGG team were nearly identical to the list of individuals with greatest affiliation to each group in our analysis. Additionally, substructures and leadership roles identified by the original authors corresponded strongly to the levels of affiliation we discovered. FOG sits well among a variety of mathematical approaches which have supported the original intuitive analyses. However, these have usually relied on separate techniques to distinguish groups, leaders, and internal structures. One advantage of FOG is the ability to unify these multiple levels of analysis under a simple model.

On the subject of validation, many link analytic methods, including k-Groups and iterative deduplication, have been validated from a data mining perspective by testing the ability of the method to rediscover groups from artificially generated data. We plan to conduct this type of examination when we complete a new fitting algorithm to replace FOG-Greedy.

*Interstitial Roles.* The existence of interstitial roles, where an individual retains several group affiliations, was our principle motivation for developing a fuzzy grouper. We identified many such individuals in our analysis, fitting several profiles. With great frequency, the most apparent leaders of a group had weak

ties to other groups as well, as did those members with the least affiliation to any group. The differentiation of these two roles, as well as the surprising result that most groups contained a well-embedded middle tier, would be difficult without FOG's novel properties: the combination of multiple memberships and degrees of membership. As FOG is applied to additional data sets we expect that a better understanding of individual roles based on multiple memberships will emerge. The FOG approach holds promise of providing a mathematical bases for capturing and defining some critical types of social roles not heretofore measurable.

It's worth noting that FOG did not always identify as interstitial the individuals whom we would expect. In some cases, such as with Sampson's waverers, individuals who were considered interstitial by an observer were placed in single groups by FOG. Conversely, some of the "secondary" clique members in the DGG dataset would appear to be interstitial on a reading of our charts, but were only claimed by a single clique in specific surveys conducted by Davis et al. The distinction between members who are simply weakly connected and those who fill an actively interstitial role may be beyond our level of analysis. Alternatively, noise may have been introduced in the specific data we examined, or results may have been misinterpreted by the original observers. Since analysis of interstitial

roles is a vital component of FOG, future work should investigate in depth what factors in data affect our ability to differentiate roles.

*Generating link data from networks.* Although the theory underlying the FOG model requires link data that indicates a shared context between members, we are optimistic about the ability to examine single-mode network data by generating fake data from simulated interactions. In the Sampson data, we were able to affirm existing knowledge about the monastery social groups using this approach, while generating new theories.

A crucial aspect of this analysis was to connect the final results with the assumptions under which link data was generated. Since Breiger's matrix indicated relationships between novices that could lead to interaction, we built our link generator as a simulator of social contexts that spread "infectively" through iterative interactions. This type of link increased the observation frequency of high betweenness individuals, but we might expect those individuals to be disproportionately represented in real data recording this type of interaction. Understanding this bias helped us interpret the difference between embedded and interstitial members when interpreting the role of novice Winfred (12), the last leader of the young Turks.

One potential criticism of the random tree link model is that it discards the directionality of links. Since FOG interprets only the presence or absence of an individual in a link, no distinction is drawn between individuals originating a random observation and those added subsequently. This affects the placement of individuals like Amand (13), who appeared in many interactions with individuals whose admiration or affection he did not reciprocate. One could again argue that many types of real data would have similar confusion, but it is also possible that a link model could be developed to include this information.

Networks of different relations may require different link models. In a formal communication network, such as a corporate hierarchy, where messages pass along a fixed route from source to destination, a random walk would be more appropriate than a random tree. It might be convenient to analyze 2-mode networks by simply interpreting one of the modes as links, but that decision should similarly depend on the type of relationship represented in the network.

Link generation for multi-mode networks is another direction we intend to develop to further FOG's applicability.

*Analyzing and visualizing fuzzy relationships.* Social groups with binary memberships can be analyzed by common statistical techniques. For example, when Davis et al. introduced the southern women dataset as overlapping cliques, they were able to investigate the character of each clique by taking aggregate statistics over its members. The same analysis would be non-trivial for a FOG cluster. What is the mean income of the members of a fuzzy group? The question is especially difficult because our results are intended to denote a level of participation, and not necessarily the degree to which members are representative of their group. If fuzzy groupings become a useful analytic tool, new measurements will have to be developed or adapted to properly interpret the new information given.

We've barely scratched the surface on that work in our intuitive analysis of the clusterings in this paper, but we've tried to uphold several principles in our analysis. First, membership values should not be examined independently of the context of other memberships held by the same individual and to the same group. Groups or individuals may have different average memberships, for reasons that have less to do with the actual importance of those memberships than with the nature of group events or the way data was collected or generated. Secondly, the novel strength of grouping with multiple, variable memberships is

the ability to compare several simultaneously occurring memberships in individuals. We intend use FOG to define and investigate roles that are defined in terms of multiple memberships, rather than to rehash issues of internal group structure that have been examined by other algorithms.

At this phase in our understanding of fuzzy overlapping groups, visualizations play an especially important role by influencing the types of patterns we can identify intuitively. We've presented two visualization paradigms in this paper, one indicating individuals' memberships to groups as a weighted two-mode network, and the other a spectrographic view providing all membership levels explicitly in bar chart form. As with most visualizations of overlapping clusters, placement of individuals can be difficult as the page does not have enough dimensions to represent all association patterns. We had few enough groups in both of our analyses that we were able position individuals for reasonable clarity, but this would not be true in more complicated datasets. We've experimented with several heuristics for laying out more than two groups in spectrographic figures, but more work needs to be done in this area.

*Future algorithmic adjustments.* In addition to the analytic work described above, three major improvements to the FOG-Greedy algorithm, while maintaining the

FOG stochastic model, are called for. First, while many datasets fall within the range of feasible analysis using FOG-Greedy, the time complexity of the algorithm must be improved to allow analysis of larger datasets. Secondly, fit metrics that would allow the tool to recommend a number of clusters to the analyst rather than requiring one as input need to be developed. Finally, we need to move to the analysis of streaming links, so that FOG can update its clustering predictions as more link data becomes available.

FOG represents a significant movement forward in our ability to identify groups as it enables the location of fuzzy groups. Fuzzy groups are a more natural and compelling way of thinking of human social groups. An unintended consequence of this approach is that the strength of membership in groups and the prevalence of exclusive members is diagnostic. We saw historical case study evidence that the strength of membership was valuable in predicting the willingness of actors to act with their group; e.g., in the case of the Sampson data, the strength of group membership is a good indicator of the order of leaving. We saw similar evidence that the higher the prevalence of exclusively tied individuals the higher the likelihood that the group would fission into an isolated component; e.g., in the case of the DGG group 2 is predominantly composed of exclusive members and it is the group that ultimately fissioned off.

While preliminary, and based on only two case studies, these findings are strongly suggestive. As such, we expect that the fuzzy group approach may be key to building a mathematics of emergent group phenomena.

## References

- Airoldi, E., D. Blei, E. Xing and S. Fienberg, 2005. A Latent Mixed Membership Model for Relational Data. Proceedings of the ACM Link-KDD Workshop in conjunction with ACM SIG-KDD.
- Airoldi, E., D. Blei, E. Xing and S. Fienberg, 2006. Mixed membership stochastic block models for relational data with application to protein-protein interactions. Proceedings of ENAR Annual Meetings.
- Battacharya, I. and L. Getoor, 2004. Deduplication and Group Detection Using Links. KDD Workshop on Link Analysis and Group Detection.
- Blei, D.M., A.Y. Ng and M.I. Jordan, 2003. Latent Dirichlet Allocation. Journal of Machine Learning Research 3, 993-1022.
- Borgatti, S.P., M.G. Everett and L.C. Freeman, 2005. UCINET 6. Analytic Technologies.
- Breiger, R., S. Boorman and P. Arabie, 1975. An Algorithm for Clustering Relational Data with Applications to Social Network Analysis and Comparison with Multidimensional Scaling. Journal of Mathematical Psychology 12, 328-383.
- Christley, R.M., G.L. Pinchbeck, R.G. Bowers, D. Clancy, N.P. French, R. Bennett and J. Turner, 2005. Infection in Social Networks: Using Network Analysis to Identify High-Risk Individuals. American Journal of Epidemiology 162, 1024-1031.
- Clauset, A., M.E.J. Newman and C. Moore, 2004. Finding community structure in very large networks. Physical Review E 70.
- Davis, A., B.B. Gardner and M.R. Gardner, 1941. Deep South: A Sociological and Anthropological Study of Caste and Class.
- Freeman, L.C., 1992. The Sociological Concept of 'Group': An Empirical Test of Two Models. American Journal of Sociology 98, 55-79.
- Girvan, M. and M.E.J. Newman, 2002. Community Structure in Social and Biological Networks. Proceedings of the National Academy of Sciences USA 99, 7821-7826.

Kashima, H. and Y. Tsuboi, 2004. Kernel-Based Discriminative Algorithms for Labeling Sequences, Trees, and Graphs. Proceedings of 21st International Conference on Machine Learning.

Kubica, J., A. Moore, D. Cohn and J. Schneider, 2003a. cGraph: A Fast Graph-Based Method for Link Analysis and Queries. Third Workshop on Link Analysis, Counterterrorism and Security, SIAM National Conference on Data Privacy.

Kubica, J., A. Moore and J. Schneider, 2003b. Tractable Group Detection on Large Link Data Sets. IJCAI Text-Mining and Link-Analysis Workshop.

Lorrain, F. and H.C. White, 1971. Structural Equivalence of Individuals in Social Networks. Journal of Mathematical Sociology 1, 49-80.

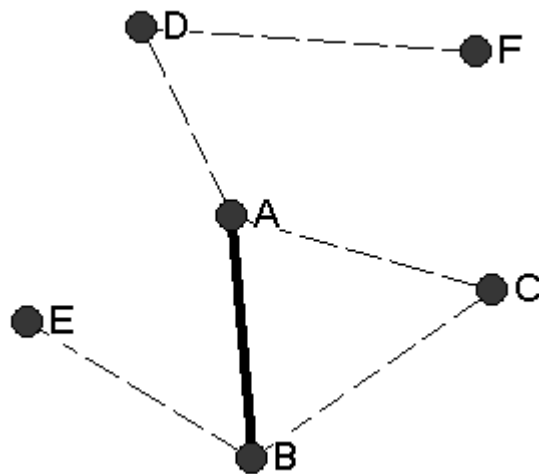
Newman, M.E.J., 2004a. Analysis of Weighted Networks. Physical Review E 70, 056131.

Newman, M.E.J., 2004b. Coauthorship networks and patterns of scientific collaboration. Proceedings of the National Academy of Sciences USA 101, 5200-5205.

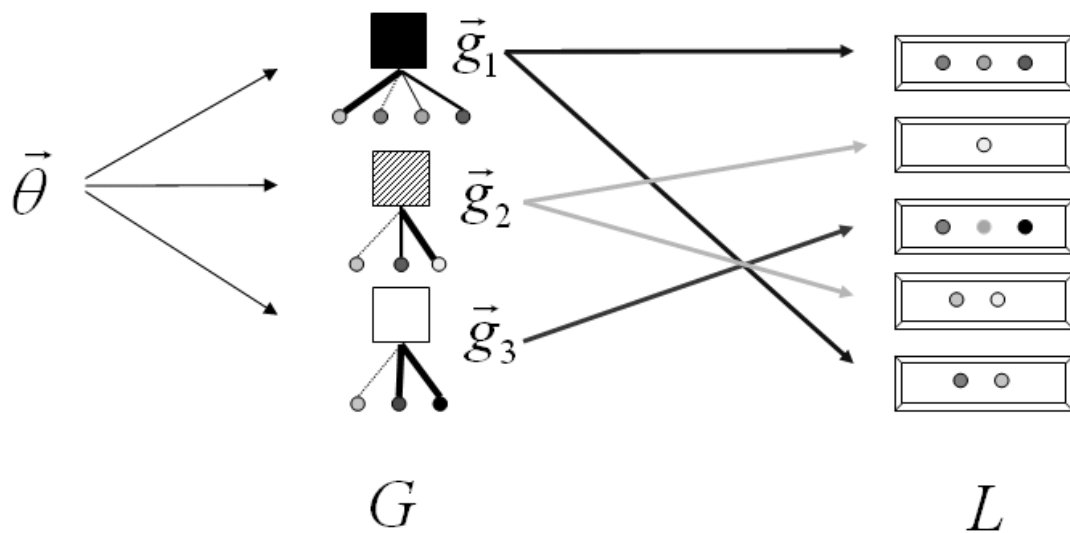
Newman, M.E.J. and M. Girvan, 2004. Finding and Evaluating Community Structure in Networks. Physical Review E 69.

Page, L. and S. Brin, 1998. The Anatomy of a large-scale hypertextual (Web) search engine. Computer Networks and ISDN Systems 30, 107-117.

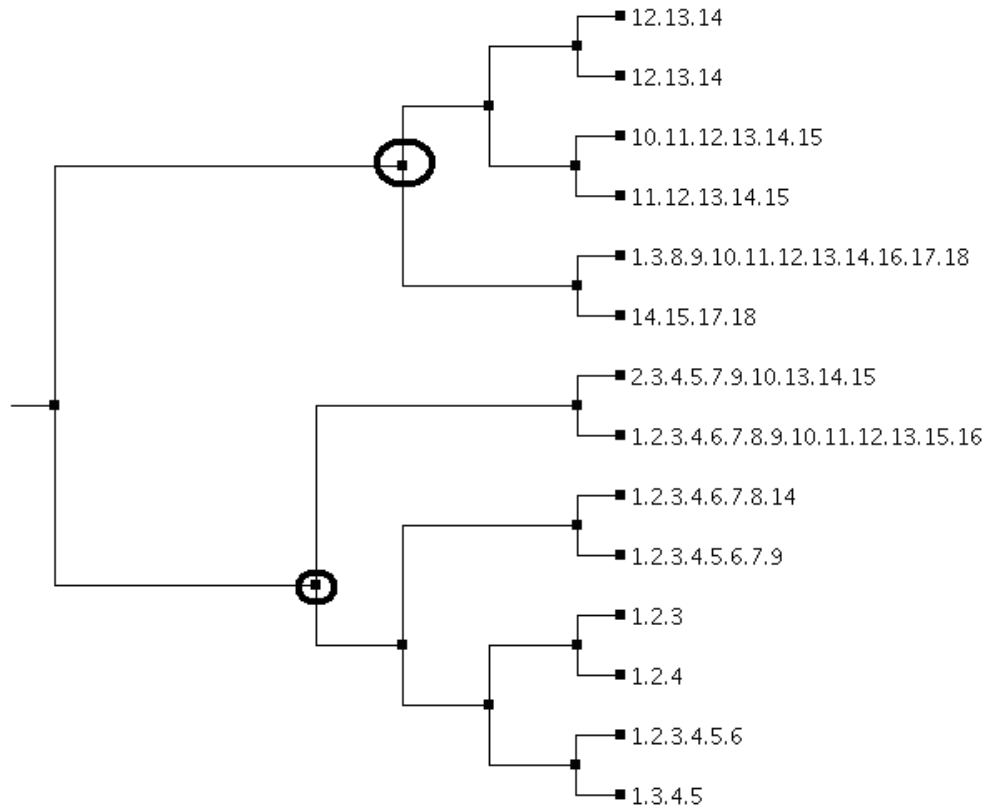
Sampson, S.F., 1968. A Novitiate in a Period of Change: An experimental and case study of social relationships. Unpublished Doctoral Dissertation.



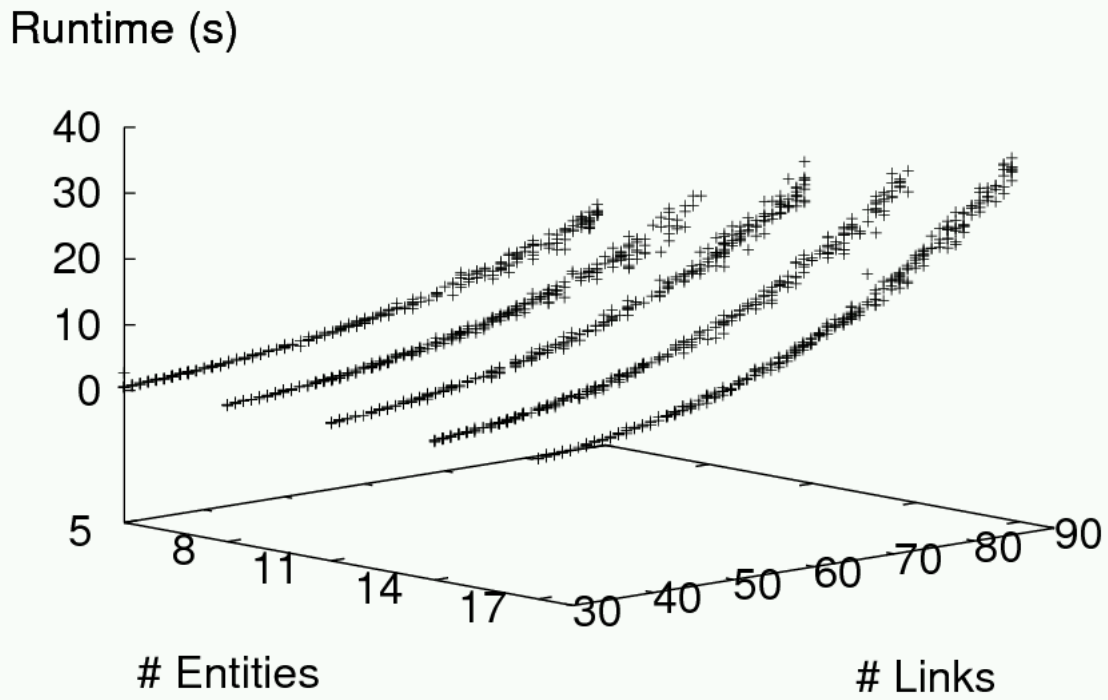
**Figure 1.** Random tree in progress.  
A and B have been visited; C, D and E are candidates.



**Figure 2.** Illustrating relationships in the FOG model.

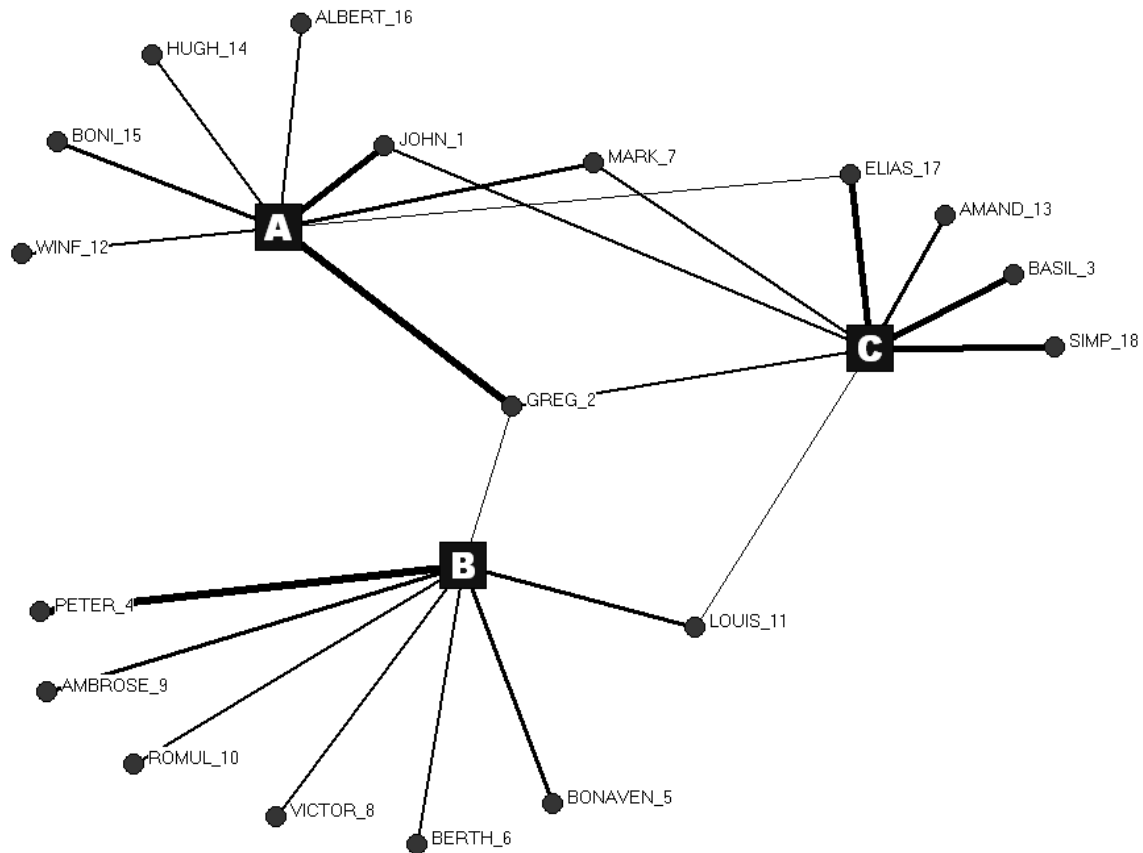


**Figure 3.** Clustering tree from the DGG Dataset

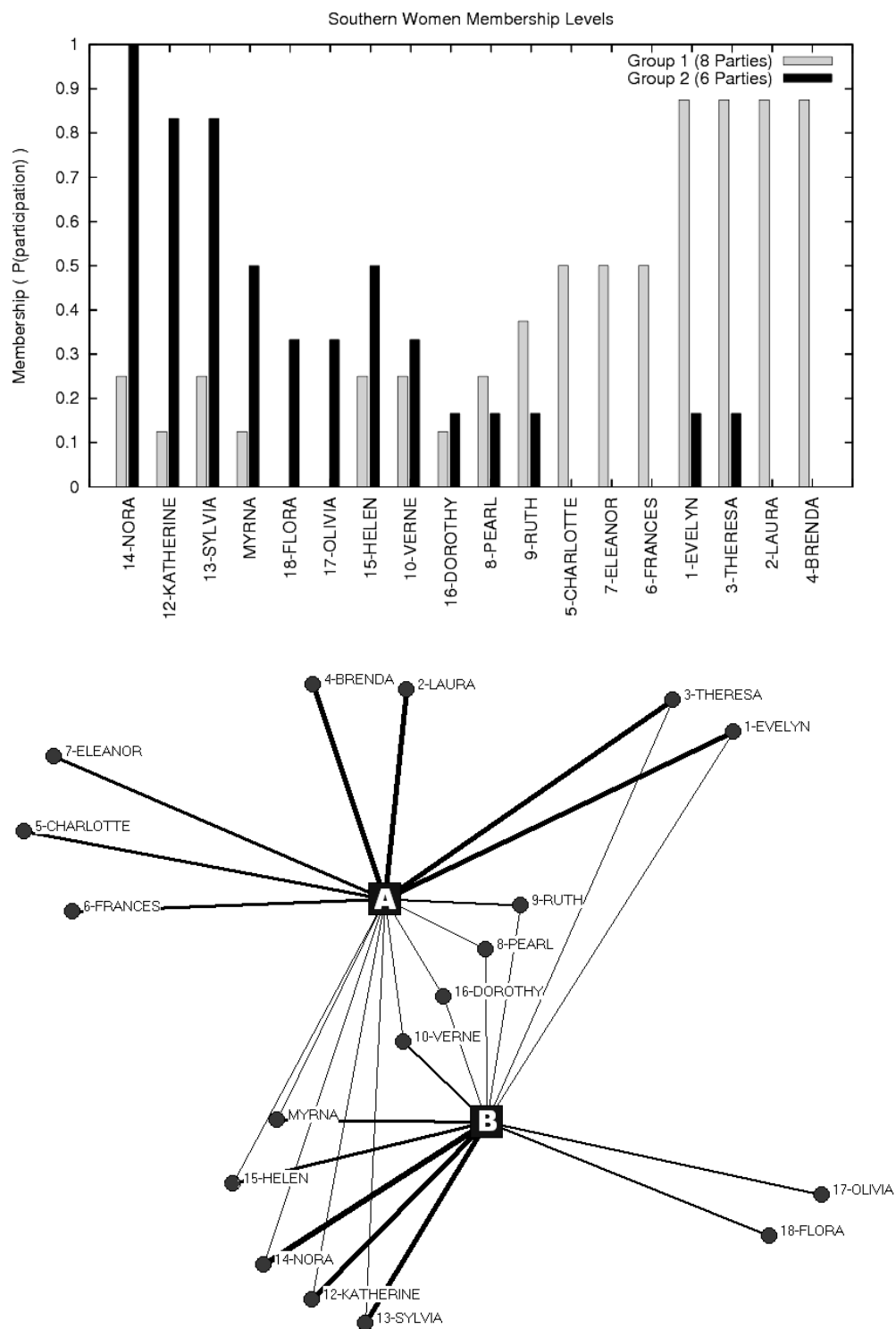


**Figure 4.** FOG Runtime v. # Entities and # Groups

|



**Fig 5.** Fuzzy Groups in Sampson's Monastery. Group A corresponds to the "young turks", group B to the "loyal opposition", and group C to the "outcasts".



**Fig 6.** FOG-Greedy 2-clustering of the DGG dataset, spectrographic (top) and network (bottom) representations