



## Abstract

Frank's (1995) clustering technique for one-mode social network data is adapted to identify positions in affiliation networks by drawing on recent extensions of  $p^*$  models to two-mode data. The algorithm is applied to the classic *Deep South* data on southern women and the social events in which they participated with results comparable to other algorithms. Monte Carlo simulations are used to generate sampling distributions to test for the presence of clustering in new data sets and to evaluate the performance of the algorithm. The algorithm and simulation results are then applied to high school students' transcripts from one school from the Adolescent Health and Academic Achievement (AHAA) extension of the National Longitudinal Study of Adolescent Health.

## 1. Introduction: Preserving the Duality of Affiliation Networks

We introduce a new algorithm for identifying positions in affiliation networks consisting of sets of actors and the events in which they participate. The algorithm represents a synthesis of Frank's (1995) network algorithm for one-mode data with recent adaptations of  $p^*$  models to two-mode data (Skvoretz and Faust 1999). Through this synthesis, the algorithm preserves the duality of the two modes of data (actors and events) that is lost in reductions to a single mode (Wasserman and Faust 1994, chapter 8). The algorithm is applied to the classic affiliation network of southern women and the social events in which they participated (Davis et al., 1944) as well as to our motivating example of high school transcripts from one school from the Adolescent Health and Academic Achievement (AHAA) extension of Adolescent Health. In future research, we will apply the algorithm to all Add Health schools (for which there is adequate transcript data), to explore links between position membership and health, educational, and work behaviors and attitudes.

Most techniques for identifying clusters or categories of actors focus on commonalities of attributes (e.g. race, age). But such approaches are essentially asocial (Coleman 1958; White, Boorman and Breiger 1976), presuming that actors will act or think similarly merely because they share common attributes. Alternatively, network approaches offer a great advancement, identifying cohesive subgroups of actors who engage in frequent direct interactions (e.g., Cartwright 1968; Collins and Raven 1968; Frank 1995) or blocks of actors who engage in structurally similar patterns of interaction (e.g., Borgatti, Everett and Shirey 1990; Borgatti and Everett 1994; Merton 1957; Nadel 1957; White, Boorman and Breiger 1976). The underlying premise is that an actor's behaviors and thoughts are related to the behaviors and thoughts of others in the actor's cluster.

The network approaches described above are defined in terms of one-mode data – who talks to whom, or who interacts directly with whom. But positions may also be defined by participation in common events or experiences (Breiger 1974) as represented in affiliation networks. Thus positions can locate actors within the social space defined by participation in events.

The events may serve as objects through which actors express pre-existing commonalities or the events may serve as experiences that socialize a set of actors (Simmel 1950; Seidman 1978). Actors who participate in common events may also have increased opportunities to interact with one another (Homans 1950; McPherson and Smith-Lovin 1982). Thus the events can serve as foci (Feld 1981) through which group phenomena may be manifest.

There are many examples of events that can form foci for social positions. Durkheim (1915) describes banging on a drum as an elemental activity that can elicit a common social position, and, as Skvoretz and Faust (1999) put it “Simmel forcefully contends that people are defined socially by the intersection of the various collectivities (family, occupation, neighborhood, voluntary organizations) to which they belong” (page 254 – see Simmel 1950, 1955). At the macro level, Laumann et al., (1978), argued that interconnections among organizations generate an important aspect of social structure for communities. Classic examples in the social network literature include attendance at community events (Davis et al 1941), membership in elite social clubs and policy groups (Domhoff 1971), and participation in corporate boards (Galaskiewicz 1985; Levine 1972). In the motivating example for this manuscript, high school transcripts reflect opportunities for students’ formal and informal interaction.

Each of the above examples of people and events can be represented as an affiliation network. Like other network data, positive values in an affiliation network establish a relation or link between units. But in an affiliation network the relations are between actors and events<sup>1</sup>. Thus actors are not represented as directly interacting with other actors, but instead are related to each other through participation in common events. Likewise, events are socially defined in that they are similar if they attract common participants, even if the events are different in content and format. Thus affiliation networks represent the duality of social experience, defined by actors and events (Breiger 1974, 1991).

There have been only a few attempts to identify social positions from affiliation networks (e.g., Bonacich and Domhoff 1981; Borgatti and Everett 1997; Breiger 1991; Doreian et al 2004; Wasserman and Faust, chapter 8). Furthermore, as noted by Bonacich and Domhoff, few techniques for analyzing such data preserve the duality of the data. That is, most algorithms operate on reductions of the data into sets of actors or sets of events, but do not simultaneously attend to information regarding actors *and* events.

Those few approaches that preserve the duality of the data have critical limitations. Bonacich and Domhoff’s theoretically sound approach based on latent class analysis (LCA) requires estimation of a large number of parameters, many of which may be near boundary values of zero (an actor in a given class is highly unlikely to participate in an event) or one (an actor in a given class is extremely likely to participate in an event). In addition, these parameters must be precisely estimated, which

---

<sup>1</sup>An affiliation network is a particular form of a two-mode network in which the modes consists of actors and events (Wasserman and Faust 1994, page 40)

typically will be possible only when there are a large number of actors relative to the number of events. In contrast, Borgatti and Everett's (1997) social network approach does not explicitly require estimation of a large number of parameters. But their clustering algorithm operates on constructed distances that are themselves modifications of the original data, introducing uncertainty that is not reflected in the interpretation of the results (cf. Doreian et al. 1994).

Most recently, Doreian et al., (2004) adapted the approach of Doreian et al., (1994) to two-mode data. Their approach has several new advantages including use of a general fit criterion and graphical representation of the data. We expand on their approach generally by drawing on the approach of Frank (1995) for one-mode data. In particular, we maximize a fit criterion defined by recent extensions of  $p^*$  mod<sup>2</sup>els to two-mode data (Skvoretz and Faust 1999). Furthermore, the algorithm we employ determines the number of positions as it iterates, eliminating the need to pre-specify the number of positions or to compare and evaluate multiple, similar, representations of the data. This is especially important because in future research we intend to apply the algorithm to 70-75 Add Health high schools for which adequate transcript data are available.

In the next section we explore constraints, opportunities and dependencies in course-taking as the research context that motivated our study. In section 3 we briefly discuss why preserving the duality of the data is important to characterize the structure among actors and events and in section 4 we present our approach. In section 5 we apply our approach to the classic southern women data and compare it to other procedures, and graphically represent the data by embedding position boundaries in a sociogram (cf. Frank 1996; Frank and Yasumoto 1996, 1998). In section 6 we use simulation studies to help evaluate evidence of clustering in the data and the performance of our algorithm (cf. Frank 1995; Krause et al., 2004). In section 7 we apply the algorithm and simulation results to transcript data reflecting high school course-taking from the recent education component of the Add Health data set. In the discussion we consider the value of our approach to identifying positions in terms of understanding social processes in schools and other institutions.

---

<sup>2</sup>For this paper, we think of structure as discernable patterns in social networks produced by constraints, opportunities and dependencies in actors' interactions.

## 2. The Motivating Problem: High School Student Course-taking

### 2.1 The AHAA Data and Course-taking

Our effort to identify positions from affiliation networks is motivated by a new project to augment Add Health with information about educational experiences and outcomes known as Adolescent Health and Academic Achievement (AHAA). The original Add Health contains data concerning friendships and health related outcomes. The longitudinal component also includes data on employment, social networks, education, and family experiences.

Drawing on transcripts, the AHAA study adds information to Add Health regarding the social and academic experiences of course-taking. Although much has been written about the influences of direct peers (Crosnoe 2002; Eccles and Barber 1999; Matsueda 1998), recent evidence suggests that adolescents may be as influenced by those with whom they would *like to be* friends as by those with whom they are already friends (Bearman and Bruckner, 1999; Harter 1999; Lightfoot 1997). And classmates may well constitute much of the potential pool of friends, as adolescents gain first intuitions regarding others' attitudes and behaviors during attendance in common courses. Thus, participants in common courses may influence the health and educational behaviors that are measured in Add Health.

Transcript data could be used to define a network relation in terms of the extent to which any two students overlapped in course-taking. But this has two important theoretical limitations. First, such a network measure reduces the duality of the transcript data to a single dimension representing potential overlap in experience, but ignoring information regarding the types of experiences (courses) that bring students together. Second, peers may well influence each other through norms (e.g., Coleman 1990) reference (Merton 1957) and identity (Ackerlof and Kranton 2002; Coleman 1961; Eckert 1989; Everhart 1983; Foley 1990; Barber et al, forthcoming), which depend on group based mechanisms. For example, members of a social position may enforce a norm of disaffection, even though that norm only indirectly benefits them by controlling the general behavior of the group (Coleman 1990; McFarland 2001). Therefore we seek to explore influences of course-taking in terms of positions of students who take similar courses rather than just in terms of the extent to which

pairs of students take similar courses. In particular, we will use the data on course-taking for a given year to characterize the social positions within a single Add Health school. We will then relate these positions to a variety of student-level demographic and educational characteristics.

## *2.2 Dependencies in Course-Taking*

There are at least three underlying mechanisms generating dependencies in course-taking that we believe should lead to clustering. These mechanisms reflect structural constraints and opportunities shaped by schools as well as by students' preferences. First, especially in the earlier grades of high school, students must fulfil academic requirements in core subjects. The average state requires four years of English, two to three years of mathematics, and two years of science (Stevenson and Schiller 1999). Because students are limited to taking five or six courses a semester, meeting these core academic requirements constrains students' opportunities to take electives. Thus students in the early grades will tend to cluster around the academic core, and course-taking tends to cluster by grade.

Second, enrollment in many courses is limited to those students who meet certain criteria such as demonstrated academic ability or having taken the appropriate prerequisites. For example, freshmen are unlikely to be allowed to take Geometry unless they have taken Algebra I in middle school (Stevenson, Schiller and Schneider 1995). Similarly, enrollment in courses like English I Honors or Below Grade Level English is usually based on test score performance or teacher recommendations (Gamoran 1992). The combination of prerequisites and constraints is consistent with the general phenomenon of tracking, a broad stratified clustering of students that is evident in middle grades and may begin as early as pre-school (Dauber, Alexander and Entwistle 1996; Entwistle Alexander and Olson 1997; Gamoran 1992; Oakes 1985; Useem 1992).

Third, despite constraints, students have some choice in which courses they take in the typical "shopping mall" high school (Powell Farrar and Cohen 1985). Independent of academic level, students may wish to perpetuate social groups established earlier by taking courses with their friends. In addition, students may identify with social categories that manifest themselves in their interests in courses (Ackerlof and Kranton 2002; Coleman 1961; Eckert 1989; Everhart 1983; Foley 1990; Barber et al, forthcoming). Thus those students identifying as working class may tend to take vocational

courses while their more affluent schoolmates may focus on academic electives such as advanced mathematics and science. While course-taking is not simply an expression of individual dispositions, these dispositions are one factor that can lead to clustering in course-taking.

Scheduling constraints can amplify existing tendencies for clustering as schools organize course offerings in terms of tracks and cohorts (Gamoran 1987; Hallinan and Sorensen 1985; Nystrand and Gamoran 1991; Oakes 1985; Oakes and Guiton 1995; Sorensen 1987). That is, because courses offered simultaneously are mutually exclusive, some courses cannot be taken together. This increases the likelihood that other courses will be taken together. For example, consider a school in which German I is offered only first period, and French I and Geometry are offered only second period. In such a school, those advanced freshman taking Geometry may well find themselves together in German I, first period. This leaves fewer places for other freshman in German, thus directing them towards French.

Given the multiple mechanisms that can generate clustering in course-taking, we believe it reasonable to seek to identify discrete clusters of courses taken together (along with the students who take them). This is in contrast to application of procedures such as principal components, correspondence analysis, or multidimensional scaling to reduce the data into a few, *continuous*, dimensions<sup>3</sup>. Recognizing that any clustering procedure may impose a discrete organization on a fluid pattern of relations, we will use simulations to test for evidence of discrete clusters.

A critical challenge for us is that in future research we intend to characterize course-taking patterns among multiple cohorts of students in 70-75 high schools that participated in Add Health and AHAA. Therefore we seek to identify clusters of students and courses using an algorithm and maximization criterion that can be objectively applied across schools. To accomplish this we need a robust procedure because extensive reliance on adhoc modifications and interpretations will compromise the validity of any results across settings. In the next section we emphasize the importance of preserving the duality of persons and events from a technical perspective and then turn to our clustering procedure which preserves duality by drawing on the recent extension of  $p^*$  models to two-mode data (Skvoretz and Faust 1999).

---

<sup>3</sup> In fact, we tried applying many of these procedures to our data with minimal success in generating an interpretable solution. For example, principal components required a reduction of the data to a single mode and correspondence analysis typically identified only one large dimension.

### 3. Limitations in Reducing Two Modes to One

Because there are many multivariate procedures and clustering algorithms that analyze a square matrix, we considered identifying positions in two-mode data by first reducing to a single mode. In particular, one-mode matrices could represent the extent to which two events were participated in by similar people, or the extent to which two people participated in similar events<sup>4</sup>. However, attempts to reduce the data to a square matrix of events or people ultimately fail to preserve the duality of the data (Wasserman and Faust, chapter 8). The key to duality is that links between units of the same mode must pass through units in the other mode (Bonacich and Domhoff 1981; Breiger 1974, 1991; Wasserman and Faust 1994). The consequence of this duality is that important structural features of the relations between the elements of one mode can only be completely understood if one simultaneously considers the way in which these same elements form relations among the elements of the other mode. For example, the relationships between courses can only be fully understood in terms of the specific students that attend the courses.

To develop a small theoretical example, define  $y_{ij} = 1$  if actor  $i$  participates in event  $j$ , 0 otherwise. Then consider the affiliation network depicted in the actor  $\times$  event matrix  $Y$  in Figure 1. The actor $\times$ actor and event $\times$ event one-mode matrix representation of  $Y$  are shown in matrices  $Y'Y$  and  $YY'$  respectively (actors are numbers and events are letters). As can be seen in the graphs below the matrices, both one-mode matrix representations of  $Y$  reveal an absence of clustering. Each pair of actors and each pair of events are linked by exactly one member of the other mode so that the relations among elements within each mode form a maximally complete graph with no internal differentiation.

---

Insert Figure 1 about here

---

---

<sup>4</sup> These could be standardized in various ways by the sizes of events or the number of events in which each actor participates (Bonacich 1972; Wasserman and Faust 1994).

However, even a cursory examination of the bipartite graph of the two-mode matrix  $Y$  reveals a clear pattern, both in relations among actors and among events. For example, actors 1, 2, 3 are all linked to each other by a single event, A, which is itself differentiated from events B, C, and D because it is the only event that actor 4 did not attend. In fact, we would argue that it is reasonable to interpret this bipartite graph as consisting of two positions. The first position is comprised of event A and all of the actors who attend event A, while the second position consists of actor 4 and all of the events attended by 4. Such a clustering of actors and events would have been very difficult to detect in the one-mode representations of the original two-mode matrix.<sup>5</sup> In the section that follows, we develop a methodology that is explicitly designed to identify positions directly from the two-mode matrix  $Y$ . As the example above illustrates, by working directly with  $Y$ , we can detect a pattern in the relations among elements of either mode that could not be identified through an analysis of a single mode.

#### 4. Adapting Frank's (1995) Algorithm for One-mode Data to Two-mode Data

To preserve the duality in two-mode data, we extend Frank's (1995, 1996) technique for identifying cohesive subgroups in one-mode data to identify positions in two-mode data. Consider the following logit model for one-mode data defined for  $x_{ii'} = 1$  if a relation is present between actors  $i$  and  $i'$ , 0 otherwise:

$$\log \left| \frac{x_{ii'} - 1}{1 - x_{ii'}} \right| = \theta_0 + \theta_1 \text{same subgroup}_{ii'} \quad , \quad (1)$$

where:

$\text{same subgroup}_{ii'} =$  1 if actors  $i$  and  $i'$  are members of the same subgroup,  
0 otherwise.

---

<sup>5</sup>In a one-mode representation of the data, ties between actors could be weighted by the size of the events they jointly attended. In the example above, actors 1,2,and 3 are all joined by a single large event, while their ties to actor 4 are exclusively through smaller events. Thus, a clustering algorithm that operated on weighted ties in a one-mode representation of these data might identify the same positions that we identified by a visual inspection of the two-mode matrix. However, such analysis would require an a priori, and somewhat subjective, definition of weights.

Thus  $\mathbf{Z}_1$  is large to the extent that membership in the same subgroup increases the odds that a relation is present. Frank then developed a clustering algorithm to maximize  $\mathbf{Z}_1$ , arguing that  $\mathbf{Z}_1$  is the term that best represents the key structural feature of the subgroup – an increase in the odds of a relation being present.

Drawing on Skvoretz and Faust (1999), model (1) can be adapted from a standard network (one-mode) to an affiliation network (two-mode) by defining the relation to apply to actors and events. Thus the model for actor  $i$  and event  $j$  is

$$\log\left(\frac{P[y_{ij}=1]}{1-P[y_{ij}=1]}\right) = \theta_0^* + \theta_1^* \textit{same position}_{ij} \quad . \quad (2)$$

where *same position<sub>ij</sub>* takes a value of 1 if actor  $i$  is assigned to the same position as event  $j$ , 0 otherwise.

Frank (1995) showed that model (1) could be expressed in a limited form of the  $p^*$  framework for one-mode data (Wassrerm and Pattison 1996). Similarly, model (2) can be linked to  $p^*$  models by applying Skvoretz and Faust’s (1999) extension of  $p^*$  models to two-mode data. In particular, Skvoretz and Faust argue that, given the  $p^*$  framework, only a single count within each mode is needed to generate a triad homogeneous model. For example, they represent how actors are linked through events (event 2-stars) and events are linked through actors (actor 2-stars), while the homogeneous assumption implies that actors are interchangeable and events are interchangeable. Skvoretz and Faust argue compellingly that their model is responsive to the theoretical motivations of affiliation networks, namely “that an actor’s involvement in a particular event may depend on other actors’ involvement in that event (captured by the event 2-star count) and that an actor’s involvement in one event may depend on his or her involvement in another event (captured by the actor 2-star count)” (page 266).

Skvoretz and Faust discuss several possible extensions of their models including effects of stars greater than 2 and actor and event specific effects (making the model nonhomogeneous). Of particular interest to us is their last extension, to include effects of event membership in the same subgroup or block (i.e., position). While these models could be used to estimate unique effects for each block, the potential multitude of parameters associated with blocks makes it necessary to adapt for maximization in clustering. In particular, model (2) incorporates two key modifications. First, the block effects are considered constant across blocks – that is that there is a single theoretical tendency for actors to prefer to attend events within their block relative to events outside their block. While this

may appear unrealistic, it does not imply that ultimately maximizing  $\mathbf{Z}_1$  will result in a constant value of  $\mathbf{Z}_1$  across blocks, as deviation from the model can occur differently across blocks. Second, by expressing actors' preferences for engaging in events within their blocks,  $\mathbf{Z}_1$  captures the effects of both actor 2-stars (event-actor-event) and event 2-stars (actor-event-actor). That is, if there are many actor 2-stars or event 2-stars then  $\mathbf{Z}_1$  can be made large by assigning members actor 2-stars or event 2-stars to the same position. Thus Skvoretz and Faust's two terms are combined in  $\mathbf{Z}_1$ . As a result,  $\mathbf{Z}_1$  directly and parsimoniously quantifies a key feature of social positions; actors exhibit a preference for participating in events within their positions relative to events outside of their positions<sup>6</sup>.

Like  $\mathbf{Z}_1$  in model (1),  $\mathbf{Z}_2$  in model (2) has a direct interpretation in terms of (the log of) an odds ratio as in table 1. The odds ratio is large to the extent that actors attend events in their position (cell D) and do not attend events outside of their positions (cell A). The odds ratio is small to the extent that actors do not attend events in their positions (cell C) and actors attend events outside of their positions (cell B). Thus maximizing  $\mathbf{Z}_2$  accommodates deviation from the model by balancing the positive values in cells A and D against those in B and C (in contrast to accommodating deviation from the model by assigning actors and events to multiple and overlapping positions).

---

Insert Table 1 about here

---

Given the objective function  $\mathbf{Z}_2$ , Frank's algorithm, *KliqueFinder*, can be adapted to two-mode data:

1) Assume that each unit  $u$  ( $u$  can be an actor or event) is assigned to one and only one position (positions are non-overlapping and comprehensive)<sup>7</sup>;

---

<sup>6</sup>We recognize that  $p^*$  models, by definition, include controls to account for dependencies induced by the network. Our goal here is not to estimate or develop a full  $p^*$  model but instead to use the  $p^*$  models to motivate the criterion to be maximized to identify positions.

<sup>7</sup>We considered allowing events to be assigned to more than one position. However, we discovered that allowing for overlapping membership in one mode of the data does not result in any meaningful reduction of the data when  $\mathbf{Z}_1$  is taken as the clustering criteria. To illustrate, consider the case in which actors 1, 2, and 3 are assigned to a position because they participated in events  $a$ ,  $b$ , and  $c$ . Now consider that actor 4 participated in events  $a$  and  $b$  but not  $c$ . If events can be assigned to more than one position then actor 4 will not be assigned to a position with actors  $a$ ,  $b$ , and  $c$  because

- 2) Calculate the change in  $Z_1^*$  that would occur if unit  $u$  were reassigned to position  $g$ , for all  $u$  and for all  $g$  containing three or more units;
- 3) Execute the reassignment that maximizes  $Z_1^*$ .
- 4) Return to step 1 unless no single change can increase  $Z_1^*$ .

The following comments apply to the algorithm:

In step 1, the algorithm is initiated by identifying a seed defined by the largest value in  $Y(Y'+I)$ . The position seed contains actors and the common events in which they participated. Therefore position seeds will be attractive to common others throughout the network. In most applications, we extend to identifying seeds of three units, which must include at least one actor and one event.

In step 2, we estimate  $Z_1^*$  by adding 500 to each cell to stabilize the procedure at boundaries when some cells approach zero (cf. Agresti 1984, page 16). This prevents solutions distinguishing most of the network from a small number of units in a nearly separate component that would be associated with a high value of  $Z_1^*$ . The value of 500 was chosen based on performance of the algorithm in simulations.

In step 3, ties between reassignments that equally increase  $Z_1^*$  are broken by assessing the extent of similarity *outside* the position between units of the same mode. Note that the application of this tie-breaking criterion is relatively rare as it applies only when the maximal increase in  $Z_1^*$  can be achieved through multiple alternatives.

---

actor  $4$  can be assigned to a unique position containing only the events actor  $4$  attends,  $a$  and  $b$ . Ultimately, then, allowing events (or actors) to be assigned to multiple positions will not result in a useful reduction of the data, unless some ad-hoc or heuristic penalty is imposed for increasing the number of positions, but this is not likely to lead to objectively defined positions across multiple analyses.

Also in step 3, each position is required to contain at least three units. Therefore new positions cannot be formed by joining a pair of units, and when one member of a position of size three is reassigned the other two are reassigned to their next optimal positions (called dissolving).

Also in step 3, if a position at any point contains only actors or only events all of the members are assigned to their next optimal position (called disbanding).

In step 4, if the algorithm stops at any point with three or more units not assigned to positions (unaffiliateds), a new seed of three is identified from the unassigned units. Each unit can be a member of a position seed only once.

The algorithm undergoes four distinct phases. In the first phase, positions are seeded and can be dissolved or disbanded as in the comments above. In the second phase unaffiliated units are attached to their optimal positions, even if this decreases  $\mathcal{Z}$ . In the third phase units are iteratively reassigned to existing positions to maximize  $\mathcal{Z}$ , but no new positions can be formed. Convergence occurs when no single reassignment can increase  $\mathcal{Z}$ . In the third phase the algorithm can be cast as a form of the Expectation-Maximization (EM) algorithm (Frank 1995; Nowicki and Snijders 2001; Snijders and Nowicki 1997). As a result the solution reaches a local, if not global maximum in the likelihood (Dempster, Larid and Rubin 1972). In the last phase units are removed from their position if they have no direct links to other members of the position (e.g., an actor who attends no events in the position).

The details of the algorithm combine with maximization of  $\mathcal{Z}$  to address the critical issue of the number of positions in the data. Because  $\mathcal{Z}$  is reduced when actors do not participate in events within their positions (cell C),  $\mathcal{Z}$  cannot be made artificially large by assigning units to a few large positions. Similarly because  $\mathcal{Z}$  is reduced if actors participate in events that are not in their position (cell B),  $\mathcal{Z}$  cannot be made artificially large by assigning units to many small positions<sup>8</sup>. Thus, maximizing  $\mathcal{Z}$  balances forces for initiating new positions with forces for

---

<sup>8</sup>This is in contrast to deviance based measures of fit which improve with the number of positions. Though these measures can be corrected by penalizing the

eliminating existing positions as the algorithm determines the number of positions “on the fly,” during the iterations. In the next section we evaluate the performance of the algorithm by comparing our results with results from the frequently analyzed two-mode Deep South data (Davis et al 1941; see Freeman’s 2003 synthesis).

## 5. Application of the Algorithm to the *Deep South* Data

Davis, Gardner and Gardner (1941) revealed how the class structure in an old southern city was manifest in women’s participation in social events (as reported in the local newspaper and in interviews). As such the data are considered a classic example of an affiliation network (Bonacich and Domhoff 1981; Doreian et al 2004; Freeman 2003; Homans 1950; Wasserman and Faust 1994, page 30).

Our algorithm identified three positions among the eighteen women and fourteen events (first reported by Davis et al. Figure 3, page 148) associated with an odds ratio of 1.28. Position 1 contains women 1 through 7 and 9 and events 1 through 8; Position 2 contains women 15, 17 and 18 and event 11; Position 3 contains women 8, 10, 11 through 14 and 16 and events 9, 10, and 12 through 14.

The positions we identified compare favorably with those reviewed by Freeman (2003) and with Freeman’s overall synthesis. In particular, the clustering of women 1 through 7 and 9 is similar to that found by application of CONCORR (Breiger 1978), algebraic topology (Doreian (1979) and G-transitivity (Freeman 1992). Interestingly, the odds ratio associated with the positions we identified is slightly greater than the odds ratio associated with Freeman’s synthesis defined by women 1 through 9 versus 10 through 18 (for which the odds ratio is 1.23).

To represent our results we extend Frank’s (1996) application of multidimensional scaling (MDS) to embed position boundaries in a sociogram. We begin by applying MDS to locate people and events within each position relative to one another. Distances between people and

---

introduction of new positions, there is no universally accepted correction. Thus use of the odds ratio adheres to model (2), but maximizes the relevant parameter from the model instead of the fit statistic.

events were defined by participation as in the raw data. Following Borgatti and Everett (1997), distances between people were defined by the number of common events within the position in which they participated, divided by the number of events in the position. Thus the range was 0 for no common events to 1 for common participation in all the events within the position. Similar calculations were applied to obtain distances between events. We then used MDS at a second level to locate the positions relative to each other, defining the distance between positions in terms of the proportion of actors in one position who participated in events in the other. Locations of actors and events were then rendered in UCINET's netdraw (chosen for desirable esthetics). Position boundaries and annotation were then added to the bitmap.

---

Insert Figure 2 about here

---

The crystallized sociogram intuitively represents many aspects of the organization of the data. Women who are central to their positions, such as Evelyn and Ruth of Position 1, are plotted in the center of their positions while others, like Myrna and Verne of Position 2, are plotted relatively peripheral to their positions. Furthermore the figure reveals that though positions 2 and 3 are separated in this analyses, they are plotted closely together, consistent with Freeman's synthesis which would have combined these two positions (except for Pearl who is moved to Position 1 in Freeman's synthesis).

Of course, a unique advantage of this two-mode crystallized sociogram is that it represents events in the same social space as people. Thus Position 1 can be thought of as centered around events on Saturday (6/27) and Sunday (4/12) while the event on Monday, March 2 is relatively marginal to Position 1. Interestingly, the crystallized sociogram reveals the importance of certain events, for example Position 3 is defined by the single event occurring on a Sunday. Distinctions in terms of events, which are not possible in interpretations that focus on only one mode of the data (for example, Freeman's synthesis is only in terms of the people, not events), may reflect differences in constraints on availability – some people may be more available on the weekends than others.

The positions define a localized social context, beyond the mere aggregate of events. Thus for example Sylvia of Position 3 may have felt something of an interloper when she

participated in the event on 3-15 (Sun) of Position 1. In contrast to Sylvia, the members of Position 1 (Brenda, Charlotte, Eleanor and Laura) who attended the event on 3-15 also commonly attended many of the other events that defined Position 1. Thus Sylvia may have been less aware of, and less likely to participate in, norms that emerged during the other events that defined Position 1. Also, Sylvia and Nora might well have identified strongly with each other on 3-15, as they were the only two of Position 3 in attendance at the Position 1 event.

Though our algorithm has identified positions similar to other analyses of the *Deep South* data, there are important questions regarding the internal validity of the positions identified by the algorithm when applied to new data. In particular, we may ask whether there is evidence of clustering in any new data and, if so, whether the algorithm identified the correct, or true, positions. To address these questions we present the results of Monte Carlo simulations in the next section.

## 6. Monte Carlo Simulations to Evaluate the Internal Validity of Identified Positions

Before applying the algorithm to empirical examples we follow Frank (1995) to explore two critical features of the internal validity of the positions identified by the algorithm. First we use a range of Monte Carlo simulations to generate a sampling distribution for  $\mathbf{Z}_1^*$  to evaluate whether there is evidence of discrete clustering in the data. Second we apply the algorithm to simulated data with *known* positions to evaluate to what extent and under what circumstances the algorithm recovers the known positions.

### 6.1 Sampling Distribution for $\mathbf{Z}_1^*$

The first question we address is whether there is evidence of clustering in the observed data. Typically, within a modeling framework, this could be addressed by testing whether  $\mathbf{Z}_1^*$  is different from zero. But any test of  $\mathbf{Z}_1^*$  based on positions identified by the algorithm is likely to lead one to reject the hypothesis that  $\mathbf{Z}_1^*$  is equal to zero. After having applied the algorithm,  $\mathbf{Z}_1^*$  will have a non-zero value,  $\mathbf{Z}_{1\text{base}}^*$ , even when the data were generated at random. Therefore the

position memberships may have been imposed on a fluid pattern of relations (Barnes, 1972; Davis, 1967; Perry, 1979).

Instead of testing whether  $\mathbf{Z}_1$  is different from zero we test whether  $\mathbf{Z}_1$  is greater than  $\mathbf{Z}_{1\text{ base}}$ . To do so, we generate a sampling distribution of  $\mathbf{Z}_{1\text{ base}}$  under the null hypothesis that participation in events is independent of membership in positions. This was accomplished by applying the algorithm to a large sample of networks generated from a model with no clustering (i.e.,  $\mathbf{Z}_1=0$ ). Because we anticipated that  $\mathbf{Z}_{1\text{ base}}$  would vary with network size, we varied the number of actors from 5 to 400 and the number of events from 5 to 170.

For each network, we generated a random vector of edges,  $y_{ij}$ , based on a draw from a Bernouli distribution with probability  $p$ , where

$$p[y_{ij}=1]=\frac{e^{\theta_0^*}}{1+e^{\theta_0^*}} \quad (3)$$

The model above assumes that the probability that actor  $i$  attends event  $j$  depends only on the overall density of the network which is determined by  $\mathbf{Z}_0$  and not by any clustering effect, (i.e.  $\mathbf{Z}_1 = 0$ ). In our simulated networks,  $\mathbf{Z}_0$  varied from -2 to 1. Our final sample consisted of 1262 networks.

After generating the data, we ran the Cliquefinder algorithm on each network and recorded the resulting value of  $\mathbf{Z}_{1\text{ base}}$ . Noting that the distribution of the log of  $\mathbf{Z}_{1\text{ base}}$  was approximately normal within conditions, we regressed  $\log(\mathbf{Z}_{1\text{ base}})$  on the size and density of the sampled network in order to obtain a predicted value of  $\log(\mathbf{Z}_{1\text{ base}})$ . The OLS estimates of the parameters and the standard error of regression are in the technical appendix.

Table 2 presents the predicted value of  $\mathbf{Z}_{1\text{ base}}$  (using the exponential to backtransform) over the range of network size and density. Note that  $\mathbf{Z}_{1\text{ base}}$  increases with network size (for which there is more information) but with negative curvature.  $\mathbf{Z}_{1\text{ base}}$  also decreases with network density emphasizing the value of the approach for sparse data (but with positive curvature). Using the prediction equation in the technical appendix, we can compare the value of  $\mathbf{Z}_1$  identified by CliqueFinder for any observed data with the sampling distribution for  $\mathbf{Z}_{1\text{ base}}$  to evaluate the hypothesis of no clustering in the data. For example, for a network of 20 actors with density .3

the predicted value of  $\mathbf{z}_{\text{base}}^*$  is 1.06 associated with an odds ratio of 2.9. This can then be compared with the observed value by constructing a z-score and then testing the statistical significance.

---

Insert Table 2 about here.

---

## 6.2 Calibrating the Performance of the Algorithm: Recovery of Known Positions

The simple hill climbing algorithm described in section 4 proceeds linearly until it identifies a local maximum in  $\mathbf{z}_1$ . But there are more sophisticated algorithms that do not proceed linearly as they seek to identify a superior local, or even the global, maximum (e.g., genetic algorithms, simulated annealing, taboo searches, etc.). But rather than explore improvements in algorithms for identifying a higher maximum of  $\mathbf{z}_1$ , we follow Frank (1995) to focus on interpreting the local maximum identified by the iterative partitioning algorithm. In particular, instead of asking: "Did the algorithm identify the global maximum?" we ask: "Are the assignments identified by the algorithm *similar to* assignments associated with the global maximum?" If the answer is yes, we can have more confidence in our interpretation of the assignments.

In order to address whether or not the algorithm recovers the true assignments associated with the global maximum, we simulated two-mode data sets based on varying values of  $\mathbf{z}_0$ ,  $\mathbf{z}_1$ , and network size (as determined by number of actors). The range of each parameter across the simulations was intended to represent plausible values in observed data sets. In particular, we varied  $\mathbf{z}_0$  from -1 to -3 by 1 (representing densities of .05 to .48) and  $\mathbf{z}_1$  from .8 to 3.5 by .27 (representing odds ratios of 2.2 to 33.1). The combinations of the parameters generated 33 cells that covered the entire experimental range and from which we could discern the extent to which there was a linear relationship between each of the network characteristics and the recovery of assignments. Within each cell we generated 15 networks, with the number of actors for each network drawn from a uniform distribution with a minimum value of 30 and a maximum value of 330.

Once the number of actors was fixed, the size of each position,  $n_g$ , was determined by adding a proportion of the range of position sizes to the minimal position size. The proportion was determined by drawing a random deviate from a beta distribution with parameters  $p = 1.5$  and  $q = .3$  (defining a unimodal distribution with positive skew):

$$n_g = \text{Minimum } n_g + (\text{Beta Deviate}[p=1.5, q=.3]) \times (\text{range of } n_g),$$

where the minimum  $n_g$  was predefined as 2, the maximum was defined as 75% of the total number of actors and the range of  $n_g = \text{maximum } n_g - \text{minimum } n_g$ .

The process was continued until the position sizes summed to the number of actors in the network. The last position size was determined when the sum of the position sizes was greater than [(network size) - (minimal position size)]. After position sizes were determined, actors were effectively assigned to position. Actors 1 through  $n_1$  were assigned to position 1, actors ( $n_1+1$ ) through ( $n_1 + n_2$ ) were assigned to position 2, and so on. The number of events assigned to each position was determined in a similar manner.

After the number of actors and events were determined and assigned to positions, we then generated a random vector,  $y_{ij}$ , based on a draw from a Bernouli distribution with probability  $p$ , where

$$p[y_{ij} = 1] = \frac{e^{\theta_0^* + \theta_1^* \text{same position}}}{1 + e^{\theta_0^* + \theta_1^* \text{same position}}}, \quad (4)$$

where  $\text{same position}_{ij} =$  1 if actor  $i$  and event  $j$  are in the same position,  
0 otherwise.

In generating a standard against which to compare the performance of the algorithm, it is possible that some actors would have been better placed (in terms of maximizing  $\mathbf{2}$ ) in another position, given the simulated data  $\mathbf{Y}$ . Perhaps this is not a concern. Ambiguous position membership may reflect the possibility that the data were unreliable in terms of representing the position memberships, or that the positions were not well defined (i.e., a small value of  $\mathbf{2}$ ), and is a characteristic of the network under which the algorithm would not perform well (rightly so).

But to recognize the above concern, we identified an alternate set of known assignments against which to evaluate the performance of the algorithm. This alternate set of assignments

were based on, but not identical to, the simulated assignments. With this alternate standard we recognized it was possible to effect one or more reassignments of units that would have increased  $\mathbf{Z}_1$ . Therefore we applied the algorithm, *using the simulated assignments as the starting assignments*. Once the algorithm converged (typically after a small number of iterations), we obtained a second set of assignments to which we will refer to as the "enhanced" known assignments.

One might argue that these enhanced assignments do not represent an informative standard against which to compare the application of the algorithm because they are based, in part, on the application of the algorithm. But it is not the general algorithm which we seek to evaluate. Iterative partitioning algorithms are known to converge to local maxima (see Everitt, 1986). Rather, we seek to evaluate the results of the algorithm against the assignments associated with the global maximum in  $\mathbf{Z}_1$ . Since the enhanced assignments are obtained based on an application of the algorithm initiated with the simulated assignments, the enhanced assignments are likely to be similar to, if not exactly, the assignments associated with the global maximum in  $\mathbf{Z}_1$ . Therefore the enhanced assignments represent a valid standard against which to compare the performance of the algorithm that is initiated without prior information.

Ultimately the researcher's focus is on the set of assignments to positions. Therefore we measure the extent of recovery of the enhanced assignments in terms of the log-odds of table 3. The log odds measure incorporates information about two types of errors. Thus the algorithm would not be measured as succeeding if it simply assigned all units to one position (resulting in no error associated with cell C\* but large error associated with B\*) or each unit to a unique position (resulting in no error associated with cell B\* but large error associated with C\*).

---

Insert Table 3 about here.

---

In order to examine the performance of the algorithm under various conditions, we regressed the log odds of recovery of the "enhanced" known assignments on a set of network characteristics. A list of variables and results from the regression can be found in the technical appendix. Generally, the algorithm performs better for large values of  $\mathbf{Z}_1$  (with negative curvature) which is sensible since  $\mathbf{Z}_1$  is the indicator of clustering. The algorithm also performs

better as  $\alpha_0$  decreases (with positive curvature), indicating the benefit of the algorithm for sparse data. For example, for  $\alpha_1 = 2.4$  and  $\alpha_0 = -2$ , the algorithm is 40 times more likely to assign an actor and event to the same position if in fact they were in the same position in the “enhanced” known assignments. Aside from the key parameters  $\alpha_1$  and  $\alpha_0$ , the algorithm also performs better for larger networks (for which there is more information), and performs worse when it must use a large number of iterations to converge, suggesting that the global maximum is difficult to identify.

We briefly summarize these results in table 4 which presents the expected log odds of recovery over a range of values from two of the strongest predictors in the regression:  $\alpha_0$  and  $\alpha_1$ . Note that the performance of the algorithm is sensitive to  $\alpha_0$  primarily at the lowest values of  $\alpha_1$  where performance is weaker. That is, when actors are much more likely to participate in events in their positions, the density does not affect recovery. Linking the results of this subsection with the results of the previous subsection, the strong effect of  $\alpha_1$  on recovery suggests that the algorithm performs best when there is more likely to be evidence of clustering in the observed data. Conversely, the algorithm performs worse when there is less likely to be evidence of clustering, but this is acceptable because one is less likely to interpret the results when there is not evidence of clustering.

---

Insert Table 4 about here.

---

We now have the tools to evaluate whether there is evidence of discrete clustering in any two mode data to which the algorithm is applied, as well as to assess the extent to which the algorithm recovers the optimal positions. We now turn to our primary substantive motivation, transcripts of high school students’ course-taking, a phenomenon that is as under-analyzed using network tools as it is common in experience (for a single counterexample, see Friedkin and Thomas 1997).

## 7. Application to High School Transcripts

As noted in the introduction, the structure of course-taking in high school is typical of other positions defined by actors and events and is, in and of itself, of considerable sociological and educational interest. Our data on high school course-taking were drawn from transcripts obtained from the Wave III respondents from a single school that participated in the new study of Adolescent Health and Academic Achievement (AHAA). In the discussion that follows, we will refer to this school as Miller, a moderate sized rural public school in the Midwest. As such, it represents an important traditional class of American high schools.

Originally, we included students who were in 9<sup>th</sup> through 12<sup>th</sup> grades in 1994-1995 because it is contemporaneous with the survey information available from Wave I of the Add Health survey, and, using the transcript weights, the Wave III data are representative of the school<sup>9</sup>. But using only one year of data resulted in less than 90 transcripts for analysis (like many of the other schools in Add Health), providing insufficient information from which to identify positions with more specificity than merely broad tracks and grades. Thus we supplemented with course-taking information from transcripts from the same students but from 1995-1996 (this second sample did not include those who graduated in 1994-1995 or who entered the school in 1995-1996 because they were unlikely to be representative of the school). Thus many students contributed two sets of data, one for each year, to the analyses. While this might violate a statistical assumption of independence, our focus was on identifying the positions and we judged course taking by a 12<sup>th</sup> grader in 1995-96 to be sufficiently unique from courses he or she took as an 11<sup>th</sup> grader in 1994-1995 to warrant inclusion as a distinct case. Few students, if any, will be taking basically the same courses two years in a row and thus be in a given position both years. An important implication of our decision is that we are essentially assuming that the configuration of positions (i.e., patterns of course enrollments in an academic year) was relatively stable from 1994-95 to 1995-96. In other words, regardless of the year, 12<sup>th</sup> graders who took Calculus were also likely to be taking an advanced science course and not taking a job

---

<sup>9</sup> We sampled from the Wave III survey because the transcript data were recently collected and only permissions from those followed up in Wave III were obtained. But data contemporaneous with the high school experience were obtained in Waves I and II.

training course. The stability of schools' master schedules and systems of pre-requisites (Delany 1991; Riehl, Pallas and Natriello, 1999) suggests that this is a reasonable assumption.

The transcript data includes entries for each course taken by each student. Entries include information on the grade, number of credits, whether taken as honors course, and a Classification of Secondary School Courses (CSSC) label. From the transcript data we then generated a set of indicators for whether student  $i$  had taken each CSSC defined course  $j$  that was offered at Miller in a given year<sup>10</sup>. In our construction of the student  $\times$  course data, we included only courses from the transcript that were taken at Miller. Courses that were transferred from other schools were not considered part of the social space of Miller and were dropped from the analyses. For the combined 1994-1995 and 1995-1996 data, we analyzed data on 149 transcript units (defined as the transcript for a given person in a given year) including approximately 41 courses.

Before interpreting the results, we note that  $\mathbf{z}^* = 8.40$  while the predicted value,  $\mathbf{z}^*_{\text{base}}$ , from simulations was 7.12. The difference of the logged values compared to the standard error (from the regression in the technical appendix) generates a z-score of 2.42 with  $p \leq .008$ . Thus, consistent with the theoretical expectations developed in section 2, there is strong evidence that there were discrete positions in course-taking at Miller. Furthermore, from the recovery simulations we predict that the odds that any pair of students or courses were assigned to the same position were increased by a factor of 15 if in fact the pair were in the same true position. Thus there is strong evidence that the positions identified by the algorithm are in fact similar to (if not exactly the same as) the true positions in the data. As shown in table 5, students and the courses they took were clustered into 9 positions composed of an estimated 15 to 172 students (numbers

---

<sup>10</sup>The data do not indicate which particular class a student attended (e.g., Algebra I fourth period with Mr. Smith), but only the course taken; large courses likely reflect multiple classes while small courses may well represent a single class. Recognizing that our theoretical motivation is based partly on the argument that participation in a common event likely increases opportunities for interaction, it would be possible to weight the initial indicators inversely proportional to the size of the course (cf. Frank 1996 for extension of odds ratio to weighted data). But for simplicity, here we confine ourselves to an analysis of dichotomous indicators of course participation (either the student took the course or did not).

of students displayed at the top of each column are weighted by the transcript weights) and 1 to 8 courses (displayed as separate rows)<sup>11</sup>.

---

Insert Table 5 about here.

---

Comparable to the conditional probabilities that would be obtained by LCA, each cell in table 5 indicates the proportion of students from each position who took each of the courses. In interpreting table 5, we recognize that students from each position may have also attended courses in a given year assigned to another position. For example, Physical Education 10 is assigned to Position 1 because 98% of the 172 students in the position took it.<sup>12</sup> However, students in three other positions also took this course, but at lower rates, and tended to take different courses than those students in Position 1 (e.g., U.S. History 1 in Position 9).

---

Insert Table 6 about here

---

Table 6 lists means of a basic set of background characteristics including gender, score on a reduced version of the Peabody picture vocabulary test, self reported likelihood of attending college and grade level based on the positions to which the students were assigned using their 1995-1996 data. Together, tables 5 and 6 clearly show evidence of the three mechanisms creating clustering in course taking – graduation requirements, systems of prerequisites, and student preferences. First, graduation requirements may well be manifest in differences in the mean grade level of students across positions. For example, Position 4 contains introductory courses in each subject area that are likely required by the state and school for graduation, and therefore are indirectly required for most of the 9<sup>th</sup> graders in the position. Thus the comprehensive common

---

<sup>11</sup>One course, “agricultural mechanics, general” was assigned to a position without any students and so was omitted from analyses

<sup>12</sup> The numbering of positions reflects their rank based on the estimated number of students in each one.

core of requirements in Position 4 essentially differentiate the 9<sup>th</sup> graders from all others in the school.

Second, systems of prerequisites are manifest in increasing evidence of tracking over the grades. Position 4 is the only predominantly 9<sup>th</sup> grade position. But tracks emerge by grade 10 with those in position 9 taking lower level courses such as Informal Geometry with those in position 8 taking more advanced courses such as Advanced Biology, Physical Chemistry, and Introductory Analysis. By the upper grades a clear hierarchy was manifest, with those in position 5 taking Organic Chemistry, Genetics, and Calculus, those in Position 3 taking the general courses English 4 and United States History 2 but few uniquely identifying electives (if any electives at all), and those in position 7 taking resource courses targeted for low level students. Thus the system of pre-requisites creates increasing differentiation, and therefore clustering, over time.

Third, student preferences also clearly influenced course-taking patterns. Positions 1 and 9 each have a mix of 9<sup>th</sup> and 10<sup>th</sup> graders and each set of courses would fit into the “general” category. (E.g., Physical Education 10 and Global Education in position 1 and Informal Geometry and English 2 in position 9). But the positions are differentiated by electives, with those in position 9 showing a unique interest in business courses such as Computer Appreciation, Accounting 2, and Typewriting 3. Thus the business courses can be thought of as clustered together because they tap an underlying latent trait of “preference for business.”

Importantly, positions are not defined solely by course experiences – they are occupied by students. Table 6 shows some evidence that positions are related to background characteristics of students. Males were more represented in position 7 including resource courses and position 8 including advanced courses. This is consistent with Hedges and Nowell’s (1995) finding that there is more variation in males’ ability than in females’. Those in position 7 also had low scores on the Peabody vocabulary test and the lowest self-reported likelihood of going to college. In contrast, those (mostly seniors) in position 5 including advanced courses such as genetics and calculus and analytic geometry had the highest Peabody picture vocabulary scores and reported the greatest likelihood of attending college. Importantly, although the course positions we identified partly reflect major divisions in student background characteristics, there remains substantial variation within positions in most of the measures of student attributes examined in this study. This suggests that the positions are capturing a unique dimension of social

experiences not previously revealed by traditional characterizations of the interaction of student background and schooling.

## 8. Discussion

Instead of thinking of social structure in terms of subgroups in which interactions are concentrated, or as blocks of actors who engage in similar patterns of interaction with other *actors*, positions identified from affiliation networks force us to focus on how events bring actors together and how actors bring events together. Actors interact while attending events, and events contribute to social organization through the actors that commonly participate in them. Our approach to identifying positions then preserves the duality in affiliation networks by adapting Frank's (1995) clustering technique for one-mode data to Skvoretz and Fausts' (1999)  $p^*$  models for two-mode data.

Ultimately our analysis of transcripts extends the conceptualization of tracking in schools by focusing on the social phenomenon of positions defined by course-taking across all subjects instead of on exposure to academic content in core subject courses (e.g., Oakes 1985). As such, we identified emergent positions defined by course-taking patterns instead of *a priori* positions defined by levels of core subjects. Furthermore, our typology is unique to the school, and more fine grained than typical tracking schemes. For example, in Miller we observed an increasing differentiation of positions, and therefore social experiences, as systems of prerequisites accumulated over grades and as students in the higher grades were less constrained in their course choices. Generally, such finer grained analyses will contribute to an understanding of the social organization of the school that can then inform how norms develop, the relative status of different activities, and many related processes alluded to in ethnographies (Coleman 1961; Eckert 1989; Everhart 1983; Foley 1990).

Of course, the ultimate usefulness of our method depends on the extent to which the positions we identify constitute a substantive and significant social phenomena; and there are several good reasons to think that they do. Kubitschek and Hallinan (1998) and Hallinan and Sorensen (1985), for example, examined the impact of academic tracking on friendship formation and found students within the same, self identified academic track (measured as "general",

“academic”, and “vocational”) were more likely to become friends. It is likely that our finer measures of social position will affect friendship formation in similar ways, although we expect that our data driven approach will result in more powerful effects of positions. This is because in many schools the segregation of students by courses undoubtedly occurs at a much finer scale than is described by a three category measure of academic track. Thus the opportunities and constraints manifest in the course offerings represent expressions of the school as formal organization that then establish venues for formal and informal interactions among students.

We are also now turning to explore how much course positions account for variance in adolescent behaviors including smoking, drinking, grade point average, skipping school, working, and religious participation (Frank et al., 2005). Interestingly, in preliminary analyses positions account for considerably more (1.5 to over 30 times) variation than schools for each behavior. Furthermore, positions have their effect even controlling for prior levels of behaviors, influence through friendship, and student composition, suggesting that positional affects are manifest through processes that occur during the school year. The issue then becomes how to develop theory and measures to tap those processes. One hypothesis is that students adopt behaviors to attract new friendships. For example, students may be more likely to increase smoking if smoking is correlated with popularity in their position.

The goal of identifying social positions from course-taking represents an example of a general goal to identify social positions from two-mode data. Like course-taking, there may be constraints and structure embedded in other affiliation networks that generate clustering. In fact, the constraints on course-taking may be especially representative of contemporary societies which contain an increasing number of mutually exclusive events (e.g., Wellman and Haythornthwaite 2002). People must choose among these events given the constraints of time, location, interest, etc., and events must compete for attendance. And it is these very conditions that create dependencies in sparse data and thus pose challenges to most extant procedures. Thus our algorithm could be used to identify sub-markets of buyers and sellers (e.g., Burt 1992), invisible colleges of researchers and their publications (e.g., White et al 2003) or sets of curricular topics covered by textbooks (e.g., Valverde et al., 2002). In each case our approach can yield important insights by identifying non-overlapping, model-based positions while preserving the duality of actors and the events in which they participate. Furthermore our simulations can be used in each case to establish the internal validity of the results.

Of course, our approach may have several limitations. First, for our motivating example we analyzed transcripts which indicate courses (e.g., Physics) but not specific classes (e.g., Physics with Ms. Smith fourth period). Ideally we would obtain information about classes that are direct measures of exposure and opportunities for interaction. But participation in common courses does increase the likelihood of exposure. Furthermore, the increase in likelihood should be greater as members of a position commonly participate in many unique courses that constrain schedules, making it more likely that they will attend the same class of the more common courses. Consider a school in which Physical Education is offered in the first three periods but English 9 Honors is offered only first period and Geometry is offered only second period. In such a school, those advanced freshman taking English 9 Honors and Geometry may well find themselves together in Physical Education first period. One also could partly address the indirect measurement of exposure by analyzing weighted data reflecting the sizes of courses, and thus the reduced probability of exposure in larger courses (Frank et al 2005).

Second, we used transcripts from two consecutive years to construct the positions in Miller school. We believe this is defensible because a student's transcript in 1995-1996 should provide information beyond what could be predicted from the same student's transcript in 1994-1995. True, some courses could be predicted to be unlikely to appear on the 1995-1996 transcript if their corresponding prerequisites had not been satisfied by 1994-1995, but the information contained in a transcript reflects student preferences beyond mere lock-step sequences. Nonetheless, optimally we would have larger samples within each school up to a full enumeration of the school.

Third, the model on which we drew was an abbreviated form of a  $p^*$  model for two-mode data. Future analyses might incorporate more parameters, while maximizing a fit criterion, although this will greatly increase computation time. Critically, although we maximize a parameter instead of a fit index, our clustering approach is very much model based.

Finally, the algorithm performs best with sparse data, which is typical of network type phenomena. But it likely will be less valuable when applied to more dense data or near continuous relationships, in which case more conventional clustering approaches may apply.

Ultimately, social positions identified from an affiliation network define unique sociological entities. These positions are not defined solely in terms of actors' attributes or

patterns of interaction. Instead, social positions represent both experiences (as events) and the others who participated in them. As such they are the settings for action and interaction, and may be both the cause and consequence of beliefs and behaviors. Thus it is critical that we develop techniques to identify and represent social positions from affiliation networks to represent an integral part of the social experience.

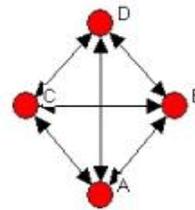
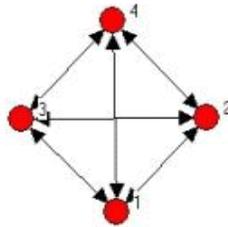
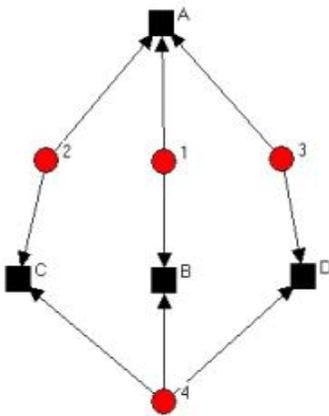
Figure 1

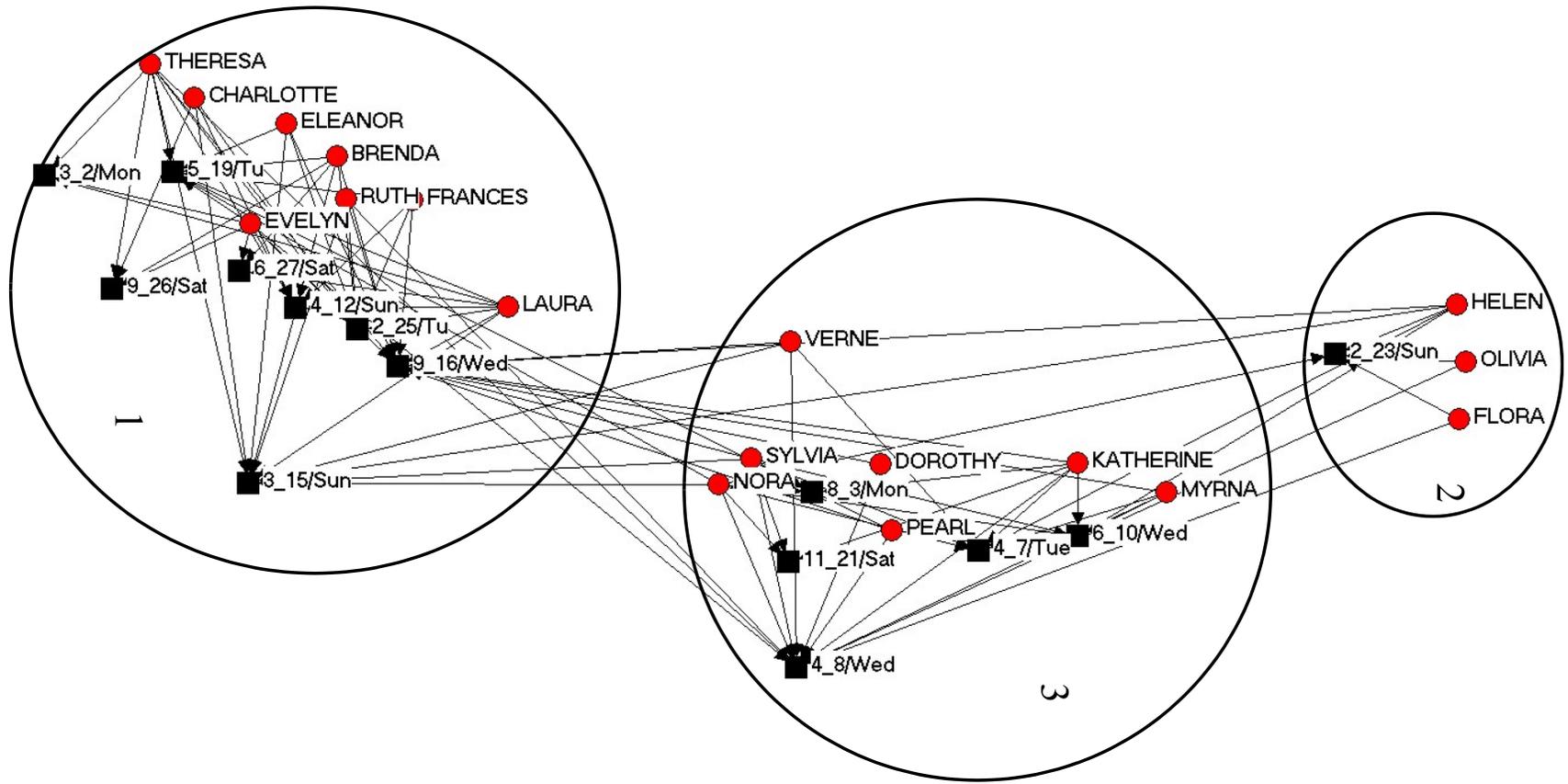
Converting a Two-Mode Matrix to One Mode

	Y			
	A	B	C	D
1	1	1	0	0
2	1	0	1	0
3	1	0	0	1
4	0	1	1	1

	Y'Y			
	1	2	3	4
1	2	1	1	1
2	1	2	1	1
3	1	1	2	1
4	1	1	1	3

	YY'			
	A	B	C	D
A	3	1	1	1
B	1	2	1	1
C	1	1	2	1
D	1	1	1	2





**Figure 2**  
 Crystallized Sociogram of Southern Women and Events

**Table 1**

**Association Between Common Position Membership and  
an Actor's Participation in an Event**

		Participation in Event	
		No	Yes
		$y_{it}=0$	$y_{it}=1$
Position Membership	Different 0	A	B
	Same 1	C	D

**Table 2**

**Predicted Value of  $\lambda_i$  (Log Odds) Under the Null Hypothesis (No Clustering Effect) by Number of Actors and Density of Network.**

	<b>Number of Actors</b>			
<b>Density</b>	20	120	220	320
0.1	1.25	1.73	2.06	2.24
0.3	0.65	1.06	1.33	1.45
0.6	0.31	0.63	0.80	0.83

**Table 3**

**Association Between Known and Observed Assignments**

		Observed Assignments	
		Different Position	Same Position
Known Assignments	Different Position	A*	B*
	Same Position	C*	D*

**Table 4**

**Predicted Recovery (in log odds) of “Enhanced” Known Assignments by Values of  $z_0$  and  $z_1$   
(from CliqueFinder Solution).**

	$z_0$		
$z_1$	-3.00	-2.00	-1.00
1.40	0.06	-0.30	-0.65
2.40	3.71	3.61	3.50
3.40	7.36	7.51	7.66



**Table 6****Student Attributes for Positions with More than Five Students.**

Position	Male	Picture		Likelihood of		Average Grade	
	%	Vocabulary		Colege		Level	
		Mean	St. D.	Mean	St. D.	Mean	St. D.
1	29.5	105.13	20.88	4.43	1.96	9.65	1.04
2	29.41	107.5	24.49	4.71	1.2	10.54	1.11
3	32.88	111.25	18.15	4.96	0.63	11.7	1.04
4	25.73	107.29	16.05	4.45	2.3	9	0
5	49.61	113.06	26.03	5	0	11.52	1.4
7	68.81	90.58	23.92	4.08	2.54	10.78	1.1
8	60.74	110.06	7.94	5	0	10.24	1.94
9	28.39	97.25	3.92	3.48	1.12	10.24	0.96

## Technical Appendix

### Predicting $\mathbf{z}_{base}$ and the Odds of Recovery Using Simulations

To evaluate evidence of clustering in the observed two-mode data, the predicted value for  $\mathbf{z}_{base}$  from simulations was:

$$\hat{\theta}_{1\ base}^* = EXP[0.549140182 + .006511461num_a + .015308725num_e - 3.463822133density - 0.011887215num_e \times density - .003187675num_a \times density - .000003951num_a \times num_e - 0.000007476num_a^2 - 0.000038274num_e^2 + 3.738069745density^2] \quad (5)$$

where  $num_a$  was the number of actors in the network,  $num_e$  was the number of events in the network and  $density$  was the density of ties in the network. Then

$$z = \frac{\log(\hat{\theta}_1) - \log(\hat{\theta}_{1\ base})}{.068277} \quad (6)$$

where .068277 was the log of the standard error of the model. A p value is then obtained from the z score.

In calibrating the performance of the algorithm, the predicted recovery from simulations was

$$\begin{aligned}
 \text{Predicted odds of recovery} = & \text{EXP}[-6.545298709 + 0.029096479 \text{num}_a + 0.130682800 \text{num}_e \\
 & - 0.716839674 \theta_0^* + 3.471847057 \theta_1^* - 0.060996971 \theta_1^{*2} + 0.255223425 \theta_1^* \times \theta_0^* \\
 & - 0.000432743 \text{num}_e^2 - 0.000027656 \text{num}_a^2 + 0.002788645 \text{num}_a \times \theta_1^* \\
 & - 0.007531484 \text{numiter1} - 0.088565601 \text{numiter2} - 0.018073363 \text{numiter3} + 0.021793975 \theta_1^* \times \text{numiter2} \\
 & + 0.056409600 \text{numposition}_a - 0.669579282 \text{numposition}_e] .
 \end{aligned} \tag{7}$$

where  $\theta_1^*$  and  $\theta_0^*$  and are the observed values after applying the algorithm to the data, *numiter1* is the number of iterations for the first phase, *numiter2* for the second phase, *numiter3* for the third phase, *numposition<sub>a</sub>* is the number of positions to which actors were assigned and *numposition<sub>e</sub>* is the number of positions to which events were assigned.

## References

- Akerlof, George A and Kranton, Rachel E (2002). Identity and Schooling: Some Lessons for Economics of Education” *Journal of Economic Literature*, Vol XL, pp. 1167-1201.
- Agresti, A. (1984). *Analysis of categorical data*. New York: Wiley and Sons.
- Albom, Mitch (1997). *Tuesdays with Morrie: An Old Man, and Young Man, and Life’s Greatest Lesson*. New York: Doubleday
- Alexander, K. L., & Pallas, A. M. (1985). School sector and cognitive performance: When is a little a little? *Sociology of Education*, 58(2), 115-128.
- Barber, B. L., Eccles, J. S., & Stone, M. R. (forthcoming). Whatever happened to the Jock, the Brain, and the Princess? Young adult pathways linked to adolescent activity involvement and social identity. *Journal of Adolescent Research*
- Bearman P, Bruckner H. Peer Influence on Adolescent GirlsÆ Sexual Debut and Pregnancy. National Campaign to Prevent Teen Pregnancy; 1999.
- Bonacich, P. (1972). Factoring and weighting approaches to status scores and clique identification. *Journal of Mathematical Sociology*, 2, 113-120.
- Bonacich, P., and Domhoff, G. (1981). Latent classes and group membership. *Social Networks*, 3, 175-196.
- Borgatti, S. P., and Everett, M. G. (1997). Network Analysis of 2-Mode Data, *Social Networks*, 19 (3) 243-269.
- Borgatti, S. P., and Everett, M. (1994). Notions of position in social network analysis. In P. Marsden (Ed.), *Sociological Methodology* (pp. 1-37). San Francisco: Jossey-Bass.
- Borgatti, S., Everett, M., and Shirey, P. (1990). LS sets, lamda sets, and other cohesive subsets. *Social Networks*, 12, 337-357.
- Borgatti, S., Everett, M.& Freeman, L. (2002). *Ucinet for Windows: Software for Social Network Analysis*. Harvard: Analytic Technologies.
- Breiger, R. L. (1974). “The Duality of Persons and Groups” *Social Forces*. 53 (2): 181-190.
- Breiger, R. (1991). *Explorations in Structural Analysis: Dual and Multiple Networks of Social Structure*. New York: Garland Press.

- Burt, Ronald S. (2002) "The Social Capital of Structural Holes," in *New Directions in Economic Sociology*, edited by Mauro F. Guillén, Randall Collins, Paula England, and Marshall Meyer. Russell Sage Foundation.
- Cartwright, D. (1968). The nature of group cohesiveness. In D. Cartwright and A. Zander (Eds.), *Group Dynamics: Research and Theory* (pp. 91-109). New York: Harper and Row Publishers.
- Clogg, C. C. & S. R. Eliason (1987). Some Common Problems in Log-Linear Analysis. *Sociological Methods and Research* 16(1): 8-44.
- Clogg, C. C. (1995). Latent Class Models. In G. Arminger, C. Clogg and M. Sobel (Eds.), *Handbook of Statistical Modeling for the Social and Behavioral Sciences*. New York: Plenum Press.
- Coleman, J. S. (1958). Relational analysis: The study of social organizations with survey methods. *Human Organization*, 17, 28-36.
- Coleman, J. S. (1961). *Adolescent Society*. Glencoe, IL: Free Press.
- Collins, R., and Raven, B. (1968). Group structure: Attraction, coalitions, communications and power. In G. Lindzey and E. Aronson (Eds.), *Handbook of social psychology*. MA: Addison-Wesley.
- Crosnoe, R. (2002). High School Curriculum Track and Adolescent Association with Delinquent Friends. *Journal of Adolescent Research*, 17, 144-168.
- Dauber, Susan L.; Alexander, Karl L., and Entwisle, Doris R. Tracking and Transitions through the Middle Grades: Channeling Educational Trajectories. *Sociology of Education*. 1996; 69(October):290-307.
- Davis, A., Gardner, B.B. and Gardner, M (1941). *Deep South*. Chicago: University of Chicago Press.
- Dayton, C. M., and Macready, G. B. (1988, March). Concomitant-Variable Latent-Class Models. *American Statistical Association*, 83, No. 401.
- Delany, Brian . Allocation, Choice, and Stratification within High Schools: How the Sorting Machine Copes. *American Journal of Education*. 1991; 99(2):181-207.
- Domhoff, G. (1971). *The Higher Circles*. New York: Random House.
- Doreian, P., Batagelj, V., and Ferligoj, A. (1994). Partitioning networks based on generalized concepts of equivalence. *Journal of Mathematical Sociology*, 19(1), 1-27.
- Durkheim, E. (1915). *The Elementary Forms of the Religious Life*. New York: The Free Press Paperback Edition.

- Eccles, Jacquelynne and Bonnie Barber. 1999. "Student Council, Volunteering, Basketball, or Marching Band: What Kind of Extracurricular Involvement Matters?" *Journal of Adolescent Research* 14(1):10-43.
- Entwisle, Doris R. and Alexander, Karl L. Entry Into School: The Beginning School Transition and Educational Stratification in the United States. *Annual Review of Sociology*. 1993; 19:401-423.
- Entwisle, Doris R.; Alexander, Karl L., and Olson, Linda Steffel. *Children, Schools, and Inequality*. Boulder, CO: Westview Press; 1997.
- Eckert, Penelope (1989). *Jocks and Burnouts: Social Categories and Identity in High School*. NY: Teachers College Press.
- Everhart, Richard. 1983. *Reading, Writing, and Resistance: Adolescence and Labor in a junior High School*. London: Routledge & Kegan Paul.
- Feld, S. L. (1981). The focused organization of social ties. *American Journal of Sociology*, 86(5).
- Foley, Douglas E. (1990). *Learning Capitalist Culture: Deep in the Heart of Texas*. Philadelphia: U of Pennsylvania Press.
- Frank, K. A. (1995). Identifying cohesive subgroups. *Social Networks*, 17, 27-56.
- \_\_\_\_\_ (1996). Mapping interactions within and between cohesive subgroups. *Social Networks*, 18, 93-119.
- Frank, K.A. and Yasumoto, J. (1998). "Linking Action to Social Structure within a System: Social Capital Within and Between Subgroups." *American Journal of Sociology* 104 (3): 642-686.
- Frank, K.A., Muller, Chandra, Schiller, Kathryn et al. (2005). Course Taking and the Social Structure of Schools. Submitted to American Sociological Association.
- Freeman, L. C. (1992). The sociological concept of 'Group': An empirical test of two models. *American Journal of Sociology*, 98(1), 152-66.
- \_\_\_\_\_. (2003). "Finding social groups: A meta-analysis of the southern women data" In Ronald Breiger, Kathleen Carley and Philippa Pattison (eds.) *Dynamic Social Network Modeling and Analysis*. Washington, D.C.:The National Academies Press, 2003.
- Friedkin NE, Thomas SL. Social Positions in Schooling. *Sociology of Education*. 1997;70(4):239-255.
- Galaskiewicz (1985). *Social Organization of an Urban Grants Economy*. New York: Academic Press.

- Gamoran, Adam. Access to Excellence: Assignment to Honors English Classes in the Transition from Middle to High School. *Educational Evaluation and Policy Analysis*. 1992; 14:185-204.
- \_\_\_\_\_(1987). The Stratification of High School Learning Opportunities. *Sociology of Education*. 60(3):135-155.
- Gaskell, Jane . Course Enrollment in the High School: The Perspective of Working-Class Females. *Sociology of Education*. 1985; 58(1):48-59.
- Guang Guo and Leah VanWey. (1999). "Family Size and Children's Intellectual Development: Is the Relationship Causal?" *American Sociological Review* 64: 169-187.
- Granovetter, M. (1974). *Getting a job: A study of contacts and careers*. Cambridge, MA: Harvard University Press.
- Hallinan, M., and Sorensen, A. B. (1985). Ability grouping and stability of instructional groups. *American Sociological Review*, 48, 838-851.
- Harter S, Fischer KW. *The Construction of the Self: A Developmental Perspective* . New York: Guilford Press. 1999.
- Hedges, L and Nowell, Amy (1995). "Sex Differences in Mental Test Scores, Variability, and the Number of High Scoring Individuals." *Science* 269: 41-45.
- Homans, G. C. (1950). *The human group*. New York: Harcourt Brace and Company.
- Kadushin, C. (1995). Friendship among the French financial elite. *American Sociological Review*, 60, 202-221.
- Krause, A., Frank, K.A., Mason, D.M., Ulanowicz, R.E. and Taylor, W.M. (2003). "Compartments exposed in food-web structure." *Nature* 426:282-285
- Kubitschek, Warren and Hallinan, Maureen (1996). Race, Gender, and Inequity in Track Assignments. *Research in Sociology of Education and Socialization* . 11:121-146.
- Laumann, E.O., Galaskiewicz, J. and Marsden, P.V. (1978). Community Structure as Interorganizational Linkages. *Annual Review of Sociology*, 4, 455-484
- Levine, J.H. 1972. "The Sphere of Influence." *American Sociological Review*. 37: 14-27.
- Lightfoot, C. (1997). *The Culture of Adolescent Risk-Taking*. New York, NY: The Guilford Press.
- Magidson, J., and Vermunt J. K. (2001) Latent Class Factor and Cluster Models, Bi-Plots, and Related Graphical Displays. In *Sociological Methodology* (pp. 223-264). Boston, MA: American Sociological Association.
- Matsueda, R. a. A., K. (1998). "The Dynamics of Delinquent Peers and Delinquent Behavior." *Criminology*, 36((2)), 269-308.

- McCutcheon, Alan L. (1987). *Latent Class Analysis*, Nebury Park, CA: Sage Publications.
- McFarland, Daniel. 2001. "Student Resistance: How the Formal and Informal Organization of Classrooms Facilitate Everyday Forms of Student Defiance." *American Journal of Sociology* 107 (3): 612-78.
- McPherson, J. M. 1982 and L. Smith-Lovin (1982) Women and Weak Ties: Differences by Sex in the Size of Voluntary Organizations." *American Journal of Sociology*, 87, 883-904.
- Merton, R. K. (1957). *Social theory and social structure*. Glencoe: Free Press.
- Mislevy, R. (1986). Recent Developments in the Factor Analysis of Categorical Variables. *Journal of Educational Statistics*, 11, Number 1, 3-31.
- Muthen, L. K. and Muthen, B. O. (1998) *M-Plus Users Guide* Los Angeles: Muthen and Muthen.
- Nadel, S. F. (1957). *The theory of social structure*. London: Cohen and West.
- Nowicki, Krzysztof, and Snijders, Tom A.B, Estimation and prediction for stochastic blockstructures. *Journal of the American Statistical Association*, 96 (2001), 1077-1087.
- Nystrand, Martin and Gamoran, Adam. Instructional Discourse, Student Engagement, and Literature Achievement. *Research in the Teaching of English*. 1991; 25:261-290.
- Oakes, J. (1985). *Keeping track: How schools structure inequality*. New Haven, CT: Yale University Press.
- Oakes, Jeannie and Guiton, Gretchen. Matchmaking: The Dynamics of High School Tracking Decisions. *American Educational Research Journal*. 1995; 32(1):3-33.
- Powell, A., Farrar, E., and Cohen, D. (1985). *The shopping mall high school: Winners and losers in the educational marketplace*. Boston: Houghton Mifflin.
- Putnam, R. D. (2000). *Bowling Alone: The Collapse and Revival of American Community*. New York: Simon and Schuster.
- Quiroz, P., Gonzalez, N., & Frank, K. A. (1996). Carving a niche in the high school social structure: Formal and informal constraints on participation in the extra curriculum. In A. Pallas (Ed.), *Research in Sociology of Education and Socialization* (pp. 93-120). Greenwich, CT: Jai Press Inc.
- Riegle-Crumb, C. (2003). "Strategies for Success: Transitions from High School to Post-secondary" Paper presented at symposium for the American Education Research Association Annual Meeting. Chicago Illinois.

- Riehl, Carolyn, Aaron M. Pallas, and Gary Natriello. 1999. "Rites and Wrongs: Institutional Explanations for the Student Course- Scheduling Process in Urban High Schools." *American Journal of Education* 107:116-54.
- Schneider, Barbara; Swanson, Christopher B., and Riegle-Crumb, Catherine. Opportunities for Learning: Course Sequences and Positional Advantages. *Social Psychology of Education*. 1998; 2:25-53.
- Seidman, S. a. F., B.L. (1978). A graph-theoretic generalization of the clique concept. *Journal of Mathematical Sociology*, 6, 139-154.
- Simmel, G. (1950). *The sociology of George Simmel*. Glencoe, Illinois: Free Press.
- Simmel, G. (1955). *Conflict and the web of group affiliations* (K. Wolff, Trans.). Glencoe, Ill: Free Press.
- Skvoretz, J. and K. Faust. (1999) "Logit models for Affiliation Networks" *Sociological Methodology*, 29, 253-278.
- Snijders, T.A.B. & Nowicki, K., Estimation and prediction for stochastic block models for graphs with latent block structure. *Journal of Classification*, 14 (1997), 75 - 100.
- Stevenson, D. L. and Schiller, K. S. (1999) "State Education Policies and Changing School Practices: Evidence from the National Longitudinal Study of School, 1980-1993." *American Journal of Education* 107:261-288.
- Stevenson, David Lee; Schiller, Kathryn S., and Schneider, Barbara. Sequences of Opportunities for Learning. *Sociology of Education*. 1994; 67(3):184-198.
- Sorensen, Aage. The Organizational Differentiation of Students in Schools as an Opportunity Structure. Hallinan, Maureen T. *The Social Organization of Schools: New Conceptualizations of the Learning Process*. New York: Plenum; 1987; pp. 103-130.
- United States Department of Education (2001). *The 1998 High School Transcript Study User's Guide and Technical Report*, NCES 2001-477, by Stephen Roey, Nancy Caldwell, Keith Rust, Eyal Blumstein, Tom Krenzke, Stan Legum, Judy Kuhn, Carl Walsberg. And, Jaqueline Haynes. Project Officer: Janis Brown. Washington, DC.
- Useem, Elizabeth L. Getting on the Fast Track in Mathematics: School Organizational Influences on Math Track Assignment. *American Journal of Education*. 1992; 100(3):325-353.
- Valverde, Gilbert A., Leonard J. Bianchi, William H. Schmidt, Curtis C McKnight, and Richard G. Wolfe. 2002. *According To the Book: Using TIMSS To Investigate The Translation Of Policy Into Practice In The World Of Textbooks*. Dordrecht, The Netherlands: Kluwer Academic Publishers.
- Vermunt, J. K. (1997) *LEM: A General Program for the Analysis of Categorical Data*. unpublished user's manual.
- Vermunt, J. K. and Magidson, J. (2000) *Latent Gold: User's Guide*. Belmont, MA: Statistical Innovations.

- Wasmus, A., Kindel, P., Mattussek, S., and Raspe, H. (1989). Activity and Severity of Rheumatoid Arthritis in Hannover/FRG and in One Regional Referral Center. *Scandinavian Journal of Rheumatology*, 79, 33-44.
- Wasserman, S., and Faust, K. (1994). *Social networks analysis: Methods and applications*. New York: Cambridge University.
- Wasserman, S., & Pattison, P. (1996). Logit models and logistic regressions for univariate and bivariate social networks: I. An introduction to Markov graphs. *Psychometrika*, 61(3), 401-426.
- Wellman, B., and Haythornthwaite, C. (2002). *The Internet in Everyday Life*. Oxford: Blackwell.
- White, Howard, Wellman, Barry and Nazer, Nancy (2003). Does Citation Reflect social structure: Longitudinal Evidence from an interdisciplinary research group.” *Journal of the American Society for Information Science and Technology* Volume 55, Issue 2 , Pages 111 - 126
- White, H., Boorman, S., and Breiger, R. (1976). Social structure from multiple networks. *American Journal of Sociology*, 81(4), 730-781.